


AUTHOR QUERY FORM

	Journal: International Journal of Forecasting Article Number: 4653	Please e-mail your responses and any corrections to: E-mail: corrections.eseo@elsevier.river-valley.com
---	--	--

Dear Author,

Please check your proof carefully and mark all corrections at the appropriate place in the proof. **It is crucial that you NOT make direct edits to the PDF using the editing tools as doing so could lead us to overlook your desired changes.** Rather, please request corrections by using the tools in the Comment pane to annotate the PDF and call out the changes you would like to see. To ensure fast publication of your paper please return your corrections within 48 hours.

For correction or revision of any artwork, please consult <http://www.elsevier.com/artworkinstructions>.

Any queries or remarks that have arisen during the processing of your manuscript are listed below and highlighted by flags in the proof.

Location in article	Query / Remark: Click on the Q link to find the query's location in text Please insert your reply or correction at the corresponding line in the proof
<u>Q1</u>	Your article is registered as a regular item and is being processed for inclusion in a regular issue of the journal. If this is NOT correct and your article belongs to a Special Issue/Collection please contact c.muttram@elsevier.com immediately prior to returning your corrections.
<u>Q2</u>	Please confirm that given names and surnames have been identified correctly and are presented in the desired order and please carefully verify the spelling of all authors' names.
<u>Q3</u>	Correctly acknowledging the primary funders and grant IDs of your research is important to ensure compliance with funder policies. We could not find any acknowledgement of funding sources in your text. Is this correct?
<u>Q4</u>	We have followed BiBTeX database for reference list. Please check, and correct if necessary. <div style="border: 1px solid black; padding: 10px; margin-top: 10px; color: red;"> Please check this box or indicate your approval if you have no corrections to make to the PDF file </div>

Thank you for your assistance.



Contents lists available at ScienceDirect

International Journal of Forecasting

journal homepage: www.elsevier.com/locate/ijforecast

Q1 A brief history of forecasting competitions

Q2 Rob J. Hyndman

Department of Econometrics & Business Statistics, Monash University, Clayton VIC 3800, Australia

ARTICLE INFO

Keywords:
 Evaluation
 Forecasting accuracy
 Kaggle
 M competitions
 Neural networks
 Prediction intervals
 Probability scoring
 Time series

ABSTRACT

Forecasting competitions are now so widespread that it is often forgotten how controversial they were when first held, and how influential they have been over the years. I briefly review the history of forecasting competitions, and discuss what we have learned about their design and implementation, and what they can tell us about forecasting. I also provide a few suggestions for potential future competitions, and for research about forecasting based on competitions.

© 2019 International Institute of Forecasters. Published by Elsevier B.V. All rights reserved.

1. Introduction

Prediction competitions go back millennia; for example, rival diviners in ancient Greece competed to predict the future more accurately (Raphals, 2013, p. 124). However, the history for general time series forecasting (i.e., predicting the future of regularly observed data over time) is much more limited, going back only about 50 years. In fact, it wasn't until computers were widely available that it became feasible for forecasting competitions to be held at all.

Time series forecasting competitions have been a feature of the *International Journal of Forecasting* and the *Journal of Forecasting* ever since the journals were founded in the early 1980s. This strong emphasis on large-scale empirical evaluations of forecasting methods, and the need to compare newly proposed methods against existing state-of-the-art methods, has played a large part in pushing researchers to develop new methods that can be shown to work in practice (Fildes & Ord, 2002).

Researchers who are new to forecasting are often surprised to learn how controversial such competitions were when they were first conducted about 50 years ago. I review this controversy in Section 2. The influential series of Makridakis competitions are discussed in Section 3, and other forecasting competitions are described in Section 4. Finally, I provide a few comments

on the future of forecasting competitions, and research about forecasting competitions, in Section 5. I do not cover forecasting competitions that are not based around time series data.

2. Early controversy

The earliest forecasting competitions were between methods rather than people. Given the communication tools available at the time, it was not feasible to conduct large-scale forecasting competitions involving many entrants spread around the world. Thus the first few competitions involved individual researchers comparing the accuracy of several methods applied to multiple time series. I only include the first two of these. From 1980 onwards, my scope is restricted to competitions involving multiple entrants.

2.1. Nottingham studies

The earliest non-trivial study of time series forecast accuracy was probably that conducted by David Reid as part of his PhD at the University of Nottingham (Reid, 1969). Building on his work, Paul Newbold and Clive Granger then conducted a study of forecast accuracy involving 106 time series (Newbold & Granger, 1974). Although they did not invite others to participate, they did start the discussion as to what forecasting methods are the most accurate for different types of time series. They presented their ideas to the Royal Statistical Society, and the subsequent

E-mail address: Rob.Hyndman@monash.edu.

<https://doi.org/10.1016/j.ijforecast.2019.03.015>

0169-2070/© 2019 International Institute of Forecasters. Published by Elsevier B.V. All rights reserved.

discussion reveals some of the erroneous thinking of the time.

One important feature of their results was the empirical demonstration that forecast combinations improve the accuracy. A similar result had been demonstrated as far back as Francis Galton in 1907 (Wallis, 2014), yet one discussant (GJA Stern) stated

“The combined forecasting methods seem to me to be non-starters ... Is a combined method not in danger of falling between two stools?”

Maurice Priestley, later to become the founding and long-serving Editor-in-Chief of the *Journal of Time Series Analysis*, said

“The authors’ suggestion about combining different forecasts is an interesting one, but its validity would seem to depend on the assumption that the model used in the Box-Jenkins approach is inadequate—for otherwise, the Box-Jenkins forecast alone would be optimal”.

This reveals a view commonly held (even today) that there is some single model that describes the data generating process, and that the job of a forecaster is to find it. This seems patently absurd to me – real data come from processes that are much more complicated, non-linear and non-stationary than any model we might dream up – and George Box himself famously dismissed it saying, “All models are wrong but some are useful”.

There was also a strong bias against automatic forecasting procedures. For example, Gwilym Jenkins said

“The fact remains that model building is best done by the human brain and is inevitably an iterative process”.

Perhaps Jenkins was reflecting the widely-held view that the type of intuitive thinking and extensive experience that was typically involved in model building cannot be represented by an algorithm or mathematical model. Subsequent history has shown that to be untrue provided that enough data are available and the model is flexible enough to capture the variations seen in real data.

Of course, human judgment still has value in forecasting, as was demonstrated by Petropoulos, Kourentzes, Nikolopoulos, and Siemsen (2018), who show that combining judgment with statistical models can lead to statistically significant improvements in forecast accuracy.

3. The Makridakis competitions

3.1. Makridakis and Hibon (1979)

Five years later, Spyros Makridakis and Michèle Hibon put together a collection of 111 time series and compared many more forecasting methods. They also presented the results to the Royal Statistical Society. The resulting paper (Makridakis & Hibon, 1979) seems to have caused quite a stir, and the discussion published along with the paper is entertaining, and at times somewhat shocking.

Maurice Priestley was in attendance again, and still clinging to the view that there was a true model waiting to be discovered:

“The performance of any particular technique when applied to a particular series depends essentially on (a) the model which the series obeys; (b) our ability to identify and fit this model correctly and (c) the criterion chosen to measure the forecasting accuracy”.

Makridakis and Hibon replied,

“There is a fact that Professor Priestley must accept: empirical evidence is in *disagreement* with his theoretical arguments”.

Many of the discussants seem to have been enamoured with ARIMA models.

“It is amazing to me, however, that after all this exercise in identifying models, transforming and so on, that the autoregressive moving averages come out so badly. I wonder whether it might be partly due to the authors not using the backwards forecasting approach to obtain the initial errors”. – W.G. Gilchrist

“I find it hard to believe that Box-Jenkins, if properly applied, can actually be worse than so many of the simple methods”. – Chris Chatfield

At times, the discussion degenerated to questioning the competency of the authors:

“Why do empirical studies sometimes give different answers? It may depend on the selected sample of time series, but I suspect it is more likely to depend on the skill of the analyst ... these authors are more at home with simple procedures than with Box-Jenkins”. – Chris Chatfield

Again, Makridakis & Hibon responded:

“Dr Chatfield expresses some personal views about the first author ... It might be useful for Dr Chatfield to read some of the psychological literature quoted in the main paper, and he can then learn a little more about biases and how they affect prior probabilities”.

3.2. M-competition

In response to the hostility and the charge of incompetence, Makridakis and Hibon followed up with a new competition involving 1001 series. This time, anyone could submit forecasts, making this the first true forecasting competition (where multiple people could submit entries) as far as I am aware. They also used multiple forecast measures to determine the most accurate method.

The 1001 time series were taken from demography, industry and economics, and ranged in length between 9 and 132 observations. All of the data were either non-seasonal (e.g., annual), quarterly or monthly. Curiously, all of the data were positive, which made it possible to

compute mean absolute percentage errors, but was not really reflective of the population of real data.

The results of their 1979 paper were largely confirmed. The four main findings (taken from [Fildes, Hibon, Makridakis, & Meade, 1998](#)) were:

1. Statistically sophisticated or complex methods typically do not produce more accurate forecasts than simpler ones.
2. The ranking of the performances of the various methods varies according to the accuracy measure being used.
3. The accuracy of the combination of various methods outperforms the individual methods being combined, on average, and does well in comparison with other methods.
4. The performances of the various methods depend on the length of the forecasting horizon.

Remarkably, the best performing method overall was “DSES”, which used a classical multiplicative decomposition ([Hyndman & Athanasopoulos, 2018](#)), with simple exponential smoothing for forecasting the seasonally adjusted data and a seasonal naive method for forecasting the seasonal component. The two forecasts were then combined. This extremely simple, and somewhat ad hoc, approach outperformed the best that experienced academic researchers could produce.

The paper describing the competition ([Makridakis et al., 1982](#)) had a profound effect on forecasting research. It caused researchers to:

- focus their attention on what models produced good forecasts, rather than on the mathematical properties of those models;
- consider how to automate forecasting methods;
- be aware of the dangers of over-fitting; and
- treat forecasting as a different problem from time series analysis.

These now seem like common-sense to forecasters, but they were revolutionary ideas in 1982. Even today, I often have to explain to other academics why forecasting is not just an application of time series analysis.

3.3. M2-competition

A few years later, a second competition was run ([Makridakis et al., 1993](#)); it used only 29 series, but with much richer contextual information, and run in real-time. Given the small sample size and the use of additional information, few general conclusions about time series forecasting methods could be drawn.

3.4. M3-competition

In 1998, Makridakis and Hibon ran their third competition, intending to take into account new methods that had been developed since their first competition nearly two decades earlier. They wrote

“The M3-Competition is a final attempt by the authors to settle the accuracy issue of various time series methods... The extension involves the inclusion of more methods/researchers (in particular in the areas of neural networks and expert systems) and more series”.

It is brave of any academic to claim that their work is “a final attempt”!

This competition involved 3003 time series, all taken from the fields of business, demography, finance and economics, and ranging in length between 14 and 126 observations. Again, the data were all either non-seasonal (e.g., annual), quarterly or monthly, and all were positive. Twenty-four entries were received (some from the organizers).

In their published results, [Makridakis and Hibon \(2000\)](#) claimed that the M3 competition upheld the findings of their earlier work; however, the results did not provide evidence supporting the first finding (that simple methods outperform more complicated methods). The two best methods were not obviously “simple”, and the Box-Jenkins’ ARIMA models did much better than in the previous competitions.

The top-performing entry was the “Theta” method, which was described by [Assimakopoulos and Nikolopoulos \(2000\)](#) in a highly complicated and confusing manner. Later, [Hyndman and Billah \(2003\)](#) showed that the Theta method was equivalent to an average of a linear regression and simple exponential smoothing with drift. Thus, it turned out to be relatively simple after all, but Makridakis and Hibon could not have known that in 2000.

The other method that performed extremely well in the M3 competition was the commercial software package ForecastPro. The algorithm used is not public, but enough information has been revealed that we can be sure that it is not simple. The algorithm selects between an exponential smoothing model and an ARIMA model based on some state space approximations and a BIC calculation ([Goodrich, 2000](#)).

The ForecastPro team also submitted an entry that involved using an automatic algorithm to select an ARIMA model (labelled BJ-automatic in the competition), and it did much better than in any previous competitions, and also better than the Holt-Winters’ method. It seems that the tendency to over-fit ARIMA models had been addressed in the 20 years since the first competition (ForecastPro uses the BIC to penalize over-parametrized models).

Even after more than 20 years of forecasting competitions, the M3 competition still generated controversy. One issue was around the statistical significance of the results ([Stekler, 2001](#)), and even whether it made sense to perform inference on the results when there is no well-defined population of possible time series. The M3 data constitute a convenience sample, and even if it can be established that method A produces statistically significantly better forecasts than method B, it is not clear what population of time series that conclusion applies to. This issue is explored by [Spiliotis, Kouloumos, Assimakopoulos, and Makridakis \(2018\)](#), building on the work of [Kang, Hyndman, and Smith-Miles \(2017\)](#).

A second area of concern was the potential for cheating. In particular, there was partial revelation of the nature of the test sets partway through the competition (Goodrich, 2001), which allowed competitors to adjust their methods on the basis of the test set. While allowing a sequence of entrants with feedback on performance is now a standard and valuable feature of many prediction competitions (Athanasopoulos & Hyndman, 2011), it must be done carefully to avoid methods being tailored to fit the test set.

There were also many calls for extensions to the competition, including

- evaluating prediction intervals (Goodrich, 2001);
- including higher frequency data such as weekly and daily data (Goodrich, 2001);
- evaluating multivariate forecasting models (Granger, 2001);
- evaluating the reasons behind the differences between methods (Hyndman, 2001), and
- identifying the circumstances in which substantial improvements are achievable (Fildes, 2001).

Finally, one issue that has arisen since the competition is the reproducibility of the results. Although the test data and the submitted forecasts are all available publicly, the computed accuracy scores do not match those in the published paper (Hyndman, 2018). This is part of a broader reproducibility problem in statistics (Peng, 2011), and forecasting in particular (Boylan, Goodwin, Moham-madipour, & Syntetos, 2015; Evanschitzky & Armstrong, 2010).

3.5. M4-competition

The most recent competition in this series organized by Spyros Makridakis is the M4 competition,¹ comprising 100,000 time series. This competition involved a few new features relative to the M and M3 competitions:

- Weekly, daily and hourly data were included, along with annual, quarterly and monthly data.
- Participants were invited to submit prediction intervals as well as point forecasts.
- There was a strong emphasis on reproducibility, with competitors being required to post their code on Github.

Extensive discussions of the results of this competition are available in other papers in this issue of the *International Journal of Forecasting*, and so will not be repeated here.

4. Other competitions

4.1. Santa Fe competitions

Parallel to the series of Makridakis competitions, mathematicians and physicists were also interested in forecasting, and a group ran their own competition at the

Santa Fe Institute, with entries closing in January 1992 (Gershenfeld & Weigend, 1993). Only six time series were included, but these were much longer than those discussed above, ranging in length from 1000 observations to 34,000 observations. Three of the series were from the physical or medical sciences, one was musical, and one was numerically simulated. Only one series (currency exchange data) was from the same domains as those used in other competitions. This choice of series led to a much greater focus on computationally-intensive nonlinear algorithms. With such a small sample, though, no general conclusions were possible.

4.2. KDD cup

Since 1997, the data mining community has held an annual competition known as the KDD cup,² organized by the Association for Computing Machinery's Special Interest Group on Knowledge Discovery and Data Mining (ACM SIGKDD). It has generated attention on a wide range of prediction problems and associated methods, and occasionally has involved time series forecasting. For example, KDD Cup 2018³ was a time series competition, with participants being asked to predict air pollution levels for two cities, Beijing, China, and London, UK. The competition ran for 31 days, and forecasts were to be made on each day, for the following 48 h.

4.3. Neural network competitions

There was only one submission in the M3 competition that used neural networks, but it did relatively poorly. To encourage additional submissions, Sven Crone organized a subsequent competition (the NN3) in 2006 involving 111 of the monthly M3 series. Over 60 algorithms were submitted, although none outperformed the original M3 contestants. The paper describing the competition results (Crone, Hibon, & Nikolopoulos, 2011) was not published until 2011.

This supports the forecasting community's consensus that neural networks (and other highly non-linear and nonparametric methods) are not well suited to time series forecasting, due to the relatively short nature of most time series. The longest series in this competition was only 126 observations long, and that simply is not enough data to fit a good neural network model.

There were some follow-up competitions,⁴ including the NN5 and NNGC1 competitions, but none of the results have ever been published.

4.4. Kaggle time series competitions

Kaggle was the first online platform dedicated to data mining competitions; it was established in 2010 by Australian economist Anthony Goldbloom. Most Kaggle competitions have involved cross-sectional prediction or

¹ <https://www.m4.unic.ac.cy/>.

² <http://kdd.org/kdd-cup>.

³ <https://www.kdd.org/kdd2018/kdd-cup>.

⁴ <http://www.neural-forecasting-competition.com/>.

classification, although a few have involved time series forecasting.

One of the earliest Kaggle competitions was on tourism forecasting, organized by George Athanasopoulos and me. It involved forecasting 793 monthly, quarterly and annual time series, all associated with tourism. The best methods were described in papers published by the *International Journal of Forecasting* (Athanasopoulos, Hyndman, Song, & Wu, 2011). Part of the competition (which was not open to outside competitors) explored the advantage of using explanatory variables (such as CPI and prices) over purely time series models. Surprisingly, the purely time series forecasts were better than the models that used explanatory variables, even ex post (with a perfect knowledge of the future explanatory variables). It is most likely that this was due to the relationships between the response and predictor variables changing over time.

In 2017, Oren Anava and Vitaly Kuznetsov organized a Web traffic⁵ competition on Kaggle. The task here was to forecast the future web traffic for approximately 145,000 Wikipedia articles. The winning method used a recurrent neural network, trained on the entire data set, with autoregressive predictors, along with some meta-data features such as the country and agent used, along with time series features such as autocorrelations. It is interesting to consider why this worked, whereas neural networks failed in the M3 and NN3 competitions. One likely factor is that the data in this competition were relatively homogeneous, and the network could be trained across all series. In both the M3 and NN3 competitions, different networks were used for each of the relatively short time series.

Two of the great benefits of the Kaggle platform (and others like it) are that it provides a leaderboard and allows multiple submissions. This has been found to lead to much better results as teams compete against each other over the duration of the competition (Athanasopoulos & Hyndman, 2011). Of course, there is a danger that the forecasts will over-fit the test data, so several precautions are taken to limit that possibility. First, only a few submissions per day are allowed, making it difficult for teams to gain enough information to fit the test data too closely. Second, the accuracy statistics that are available to the contestants during the competition are computed on only a small subset of the actual test data, with the final accuracy, computed on the remaining part of the test set, not being available until after the completion of the competition.

4.5. Global energy forecasting competitions

The GEFCom series of competitions⁶ were organized by Tao Hong in conjunction with the *International Journal of Forecasting*, and comprised three competitions, each involving several parts. They have been influential in focusing research attention on important issues in energy forecasting, and in providing benchmark methods and data on which to compare new methods. One feature of

the competitions has been that the winning entries must post complete code implementing their methods, and submit an article to the *International Journal of Forecasting* describing their approach.

The 2012 competition (Hong, Pinson, & Fan, 2014) had five main aims:

1. to improve the forecasting practices of the utility industry;
2. to bring together state-of-the-art techniques for energy forecasting;
3. to bridge the gap between academic research and industry practice;
4. to promote analytics in power and energy education; and
5. to prepare the industry to overcome the forecasting challenges introduced by the smart grid technologies and renewable integration needs.

It comprised two tracks: hierarchical load and wind power, where the contestants were asked to forecast sections of missing data and (in the case of load) one week of future data. Thus, for the most part, it was not truly forecasting, and it looked at point prediction accuracy only.

The 2014 competition (Hong et al., 2016) included four tracks: load forecasting, price forecasting, wind forecasting and solar forecasting, and was the first competition to require entrants to submit complete probability distributions rather than simply point forecasts. These were submitted in the form of 99 percentiles for each forecast period. The competition used a rolling forecast process, mimicking what happens in most real forecasting tasks, where forecasts needed to be made each week for 15 consecutive weeks. A total of 581 contestants participated, from 61 countries. The probabilistic forecasts were assessed using quantile scoring.

The 2017 competition (Hong, Xie, & Black, 2019) involved (1) hierarchical probabilistic forecasting of the hourly electricity load for ISO New England, with both zone and total load forecasts being required; and (2) probabilistic load forecasting of 183 delivery point meters of a US utility. This time, only the deciles for each forecast distribution were submitted, and again quantile scoring was used for evaluation. There were 177 entrants in total, and most of the resulting papers are yet to be published.

5. Some thoughts on the future

There is no doubt that forecasting competitions have played an important role in advancing our knowledge of what forecasting methods work. However, there has been much less research done on the conditions under which different methods work well. The publication of the large database of time series and forecasts that are associated with the M4 competition now provides researchers with an opportunity to explore this latter issue in considerable detail. The data and forecasts from the M4 competition are available in the R package *M4comp2018* (Montero-Manso, Netto, & Talagala, 2018).

⁵ <https://www.kaggle.com/c/web-traffic-time-series-forecasting>.

⁶ <http://gefcom.org>.

The data and forecasts from the M3 competition are available in the *MComp* package (Hyndman, Akram, Bergmeir, & O'Hara-Wild, 2018), which also includes data (but not forecasts) from the M competition. The *Tcomp* package contains data from the Kaggle tourism competition (Ellis, 2018), and data from the NN3, NN5, NNGC1 and GEFCom2012 competitions are provided in the *tscompdata* package (J.Hyndman, 2018). It would be helpful for other forecasting competition data (and the submitted forecasts) to be made publicly available too, in the interests of promoting research around the efficacy of different forecasting methods. For example, what are the current limits of forecast accuracy in different application areas, and how have these varied over time? Given access to all of the forecasting entries in a competition, can the winning approach be beaten using a combination of entries? We can begin to explore these questions using the M4 data, but it would be better for all forecasting competitions to provide such data as a public good.

A nice side-effect of some competitions is that they create a benchmark data set with well-tested benchmark methods. This has worked well for the M3 data, for example, and new time series forecasting algorithms have been being compared to these published results for many years. The publication of the M4 data has now provided a useful and up-to-date set of benchmarks against which to test new forecasting methods. It could reasonably be argued that no research paper proposing a new general time series forecasting method should be published unless it can be shown to perform at least as well as existing methods on some well-defined subset of the M4 data.

However, over-studying a single benchmark data set means that methods will eventually over-fit the published test data. I suspect that this has happened with the M3 data over the past 20 years, and it is likely to happen with the M4 data too, despite its much larger size. Therefore, a wider range of benchmarks is desirable, and these need to be updated regularly. Consequently, there can never be a "final forecasting competition".

One feature of both the Makridakis competitions and the GEFCom competitions that has meant they have had a substantial impact on forecasting practice is their association with the *International Journal of Forecasting*. Competitions that have been subject to careful academic scrutiny, where the data and results have been the focus of subsequent research papers, have rightfully led to changes in forecasting practice and have set benchmarks for future methodological developments in forecasting.

In surveying the last 40 years of forecasting competitions, it is apparent that the objectives of many competitions were unclear, perhaps even to the organizers. This has led, for example, to confusion over the domain to which the results apply: if a method does well in one competition, what can be said about how well it will do on data from a specific company? Ideally, future competitions should be clear about the domain to which they apply (by carefully defining the population of time series from which the sample is drawn).

The frequent use of percentage errors also makes the results somewhat difficult to apply in general. It is known

that the minimization of absolute percentage errors tends to support severe underforecasts (Gneiting, 2011). It would be better if future competitions used objective criteria that are based on well-recognized attributes of the forecast distribution. For example, minimizing MASE (Hyndman & Koehler, 2006) leads to estimates of the median of the forecast distribution, as it is equivalent to minimizing the mean absolute error.

There are many features of time series forecasting that have not been studied under competition conditions, and future competitions will surely seek to address these.

For example, few time series competitions have explored the forecast distribution accuracy (as distinct from the point forecast accuracy). The M4 competition is the first general competition to make a start in this direction by measuring the prediction interval accuracy, but it is much richer to measure the entire forecast distribution. The GEFCom2014 and GEFCom2017 competitions involved forecast distributions in the context of energy forecasting, but it would be preferable to see the forecast distributions being used in all future forecasting competitions.

New competitions will also need to select accuracy metrics that are appropriate to the forecasting task, easy to explain and understand, and independent of the scale of the data. MAPEs, and related criteria, clearly are not appropriate when the data contain zeros, or are on an interval rather than a ratio scale. Hyndman and Koehler (2006) addressed this issue by introducing the MASE, but it is only relevant for point forecasting. Winkler scores (Winkler, 1972) have been used widely for interval forecasting, but are not scale-free. The scaled version of Winkler scores that is used to assess the interval accuracy in the M4 competition seems rather ad hoc, and its properties are unknown. In any case, Askanazi, Diebold, Schorfheide, and Shin (2018) argue that comparisons of interval forecasts are problematic in several ways and should be abandoned for density forecasts. Probabilistic forecasts such as densities can be evaluated using proper scoring rules, and scale-free rules such as log density scores are available. However, there is a need for textbooks to introduce clear and accessible explanations of how to compute them and what they tell us, as such scoring rules are foreign to many who practice forecasting.

No competition to date has involved large-scale multivariate time series forecasting. While many of the time series in past competitions are probably related to each other, this information has not been provided. Again, the GEFCom competitions have been ground-breaking in this respect too, by requiring true multivariate forecasts to be provided for the energy demand in different regions of the US; however, the winning entries do not appear to have exploited the cross-correlations between regions. Some further thoughts on the design of multivariate forecasting competitions are provided by Fildes and Ord (2002).

I know of no large-scale forecasting competition that has been conducted for finance data (e.g., stock prices or returns), yet this would seem to be of great potential interest.

The frequencies of the time series studied in competitions have changed from only annual, quarterly and monthly data in the competitions up to the M3 competition, to more frequent observations such as hourly, daily, and weekly in some subsequent competitions. As data are collected more frequently using automatic sensors and scanning devices, the need for higher frequency forecasts likewise increases. The use of higher frequency data also brings with it some new challenges, including multiple seasonal patterns (De Livera, Hyndman, & Snyder, 2011), irregularly spaced observations, and the need to generate forecasts rapidly.

There has also been little focus on the conditions under which explanatory variables are useful when forecasting. The tourism forecasting competition included explanatory variables which proved to be unhelpful. The GEFCom competitions have all involved weather data as useful predictors for energy forecasting, but clearly these predictors also need forecasting, and it may be better not to use them at all over long forecast horizons.

The *International Journal of Forecasting* has been at the forefront of the research surrounding forecasting competitions since its inception, and the suggestions made here should provide enough ideas to fuel research in this space for many years to come.

References

- Askanazi, R., Diebold, F. X., Schorfheide, F., & Shin, M. (2018). On the comparison of interval forecasts. *Journal of Time Series Analysis*, 39(6), 953–965.
- Assimakopoulos, V., & Nikolopoulos, K. (2000). The theta model: a decomposition approach to forecasting. *International Journal of Forecasting*, 16, 521–530.
- Athanasopoulos, G., & Hyndman, R. J. (2011). The value of feedback in forecasting competitions. *International Journal of Forecasting*, 27(3), 845–849.
- Athanasopoulos, G., Hyndman, R. J., Song, H., & Wu, D. C. (2011). The tourism forecasting competition. *International Journal of Forecasting*, 27(3), 822–844.
- Boylan, J., Goodwin, P., Mohammadipour, M., & Syntetos, A. (2015). Reproducibility in forecasting research. *Int. J. Forecasting*, 31(1), 79–90.
- Crone, S. F., Hibon, M., & Nikolopoulos, K. (2011). Advances in forecasting with neural networks? empirical evidence from the NN3 competition on time series prediction. *International Journal of Forecasting*, 27(3), 635–660.
- De Livera, A. M., Hyndman, R. J., & Snyder, R. D. (2011). Forecasting time series with complex seasonal patterns using exponential smoothing. *Journal of the American Statistical Association*, 106(496), 1513–1527.
- Ellis, P. (2018). Tcomp: Data from the 2010 Tourism Forecasting Competition. <https://CRAN.R-project.org/package=Tcomp>.
- Evanschitzky, H., & Armstrong, J. (2010). Replications of forecasting research. *Int. J. Forecasting*, 26(1), 4–8.
- Fildes, R. (2001). Beyond forecasting competitions. *International Journal of Forecasting*, 17, 556–560.
- Fildes, R., Hibon, M., Makridakis, S. G., & Meade, N. (1998). Generalizing about univariate forecasting methods: further empirical evidence. *International Journal of Forecasting*, 14, 339–358.
- Fildes, R., & Ord, K. (2002). Forecasting competitions: their role in improving forecasting practice and research. In M. Clements, & D. Hendry (Eds.), *A Companion To Economic Forecasting* (pp. 322–353). Oxford, Blackwell.
- Gershenfeld, N. A., & Weigend, A. S. (1993). The future of time series. In A. S. Weigend, & N. A. Gershenfeld (Eds.), *Time Series Prediction: Forecasting the Future and Understanding the Past* (pp. 1–70). Reading, Mass., USA: Addison-Wesley.
- Gneiting, T. (2011). Making and evaluating point forecasts. *Journal of the American Statistical Association*, 106(494), 746–762.
- Goodrich, R. L. (2000). The forecastpro methodology. *International Journal of Forecasting*, 16(4), 533–535.
- Goodrich, R. L. (2001). Commercial software in the M3-competition. *International Journal of Forecasting*, 17, 560–565.
- Granger, C. W. J. (2001). Comments on the M3 forecast evaluation and a comparison with a study by stock and watson. *International Journal of Forecasting*, 17, 565–567.
- Hong, T., Pinson, P., & Fan, S. (2014). Global energy forecasting competition 2012. *International Journal of Forecasting*, 30(2), 357–363.
- Hong, T., Pinson, P., Fan, S., Zareipour, H., Troccoli, A., & Hyndman, R. J. (2016). Probabilistic energy forecasting: Global energy forecasting competition 2014 and beyond. *International Journal of Forecasting*, 32, 896–913.
- Hong, T., Xie, J., & Black, J. (2019). Global energy forecasting competition 2017: hierarchical probabilistic load forecasting. *International Journal of Forecasting*. To appear.
- Hyndman, R. J. (2001). It's time to move from 'what' to 'why'. *International Journal of Forecasting*, 17, 567–570.
- Hyndman, R. J. R vs Autobox vs ForecastPro vs ... <https://robjhyndman.com/hyndsight/show-me-the-evidence/> (visited on 21/12/2018).
- Hyndman, R. J., Akram, M., Bergmeir, C., & O'Hara-Wild, M. (2018). Mcomp: Data from the M- Competitions. Version 2.8. <https://CRAN.R-project.org/package=Mcomp>.
- Hyndman, R. J., & Athanasopoulos, G. (2018). *Forecasting: principles and practice* (2nd ed.). Melbourne, Australia: OTexts.
- Hyndman, R. J., & Billah, M. B. (2003). Unmasking the theta method. *International Journal of Forecasting*, 19(2), 287–290.
- Hyndman, R. J., & Koehler, A. B. (2006). Another look at measures of forecast accuracy. *International Journal of Forecasting*, 22(4), 619–688.
- Hyndman, R. (2018). tscompdata: Time series data from various forecasting competitions. <https://github.com/robjhyndman/tscompdata>.
- Kang, Y., Hyndman, R. J., & Smith-Miles, K. (2017). Visualising forecasting algorithm performance using time series instance spaces. *International Journal of Forecasting*, 33(2), 345–358.
- Makridakis, S. G., Andersen, A., Carbone, R., Fildes, R., Hibon, M., Lewandowski, R., Joseph Newton, H., Parzen, E., & Winkler, R. L. (1982). The accuracy of extrapolation (time series) methods: results of a forecasting competition. *Journal of Forecasting*, 1(2), 111–153.
- Makridakis, S., Chatfield, C., Hibon, M., Lawrence, M., Mills, T., Ord, K., & Simmons, L. F. (1993). The M2-competition: A real-time judgmentally based forecasting study. *International Journal of Forecasting*, 9(1), 5–22.
- Makridakis, S. G., & Hibon, M. (1979). Accuracy of forecasting: an empirical investigation (with discussion). *Journal of the Royal Statistical Society, Series A*, 142, 97–145.
- Makridakis, S. G., & Hibon, M. (2000). The M3-Competition: results, conclusions and implications. *International Journal of Forecasting*, 16(4), 451–476.
- Montero-Manso, P., Netto, C., & Talagala, T. (2018). M4comp2018: Data from the M4-Competition. Version 0.1.0. <https://github.com/carlanetto/M4comp2018/>.
- Newbold, P., & Granger, C. W. (1974). Experience with forecasting univariate time series and the combination of forecasts (with discussion). *Journal of the Royal Statistical Society, Series A*, 137(2), 131–165.
- Peng, R. D. (2011). Reproducible research in computational science. *Science*, 334(6060), 1226–1227.
- Petropoulos, F., Kourantzis, N., Nikolopoulos, K., & Siemsen, E. (2018). Judgmental selection of forecasting models. *Journal of Operations Management*, 60, 34–46.

- 1 Q4 Raphals, L. (2013). *Divination and prediction in early China and ancient*
 2 *Greece*. Cambridge, UK: Cambridge University Press.
 3 Reid, D. J. (1969). *A comparative study of time series prediction techniques*
 4 *on economic data*. Ph.D. thesis, Nottingham, UK: University of
 5 Nottingham.
 6 Spiliotis, E., Kouloumos, A., Assimakopoulos, V., & Makridakis, S. (2018).
 7 Are forecasting competitions data representative of the reality?
 8 Tech. rep. 3/18, Forecasting and Strategy Unit, National Technical
 9 University of Athens.
 10 Stekler, H. (2001). The M3-competition: the need for formal statistical
 11 tests. *International Journal of Forecasting*, 17, 576–577.

- Wallis, K. F. (2014). Revisiting francis galton's forecasting competition. *Statistical Science*, 29(3), 420–424.
 Winkler, R. L. (1972). A decision-theoretic approach to interval estimation. *Journal of the American Statistical Association*, 67(337), 187–191.

Rob J. Hyndman is a Professor of Statistics and Head of the Department of Econometrics & Business Statistics at Monash University, Australia. He was Editor-in-Chief of the *International Journal of Forecasting* from 2005-2018.

UNCORRECTED PROOF