



**Progress in Forecasting with Neural Networks? Empirical Evidence from the NN3 competition**

Journal:	<i>International Journal of Forecasting</i>
Manuscript ID:	INTFOR_1010219
Manuscript Type:	Original Article
Keyword:	Time series, Forecasting competitions, Comparative studies, Neural networks, Automatic forecasting
Abstract:	<p>This paper reports the results of the NN3 competition, a replication of the M3 competition, and an extension towards methods of neural networks (NN) and computational intelligence (CI), assessing if progress has been made in the past 10 years since M3. Two masked subsets of 111 and 11 empirical time series of monthly M3 industry data were chosen, controlling for multiple data conditions of forecasting horizons (short/medium/long), time series length (short/long) and data patterns (seasonal/non-seasonal). Relative forecasting accuracy was assessed using the metrics of M3, and latter extensions of scaled measures. The NN3 competition attracted 59 submissions from NN, CI and statistics, making it the largest competition of CI on time series data. Its main findings include: (a) only one NN outperformed Dampen Trend, but more outperformed AutomatANN of M3; (b) Ensembles of CI-approaches performed very well, (c) a novel statistical method outperformed all statistical and CI-benchmarks, and (d) for the most difficult subset of short and seasonal series, a methodology employing Echo State Neural Networks outperformed all, highlighting the ability of NN to handle complex data and routes for future research.</p>

# Progress in Forecasting with Neural Networks? Empirical Evidence from the NN3 competition

## Abstract

This paper reports the results of the NN3 competition, a replication of the M3 competition, and an extension towards methods of neural networks (NN) and computational intelligence (CI), to assess if progress has been made in the past 10 years since M3. Two masked subsets of 111 and 11 empirical time series of the M3 monthly industry data were chosen, controlling for multiple data conditions of forecasting horizons (short / medium / long), time series length (short / long) and data patterns (seasonal / non-seasonal). Relative forecasting accuracy was assessed using the metrics of M3, and latter extensions of scaled measures. The NN3 competition attracted 59 submissions from NN, CI and statistics, making it the largest competition of CI on time series data. Its main findings include: (a) only one NN outperformed Dampen Trend on sMAPE, but more outperformed Automat ANN of M3; (b) Ensembles of CI-approaches performed very well, (c) a novel statistical method outperformed all statistical and CI-benchmarks, and (d) for the most difficult subset of short and seasonal series, a methodology employing Echo State Neural Networks outperformed all others, highlighting the ability of NN to handle complex data beyond prior expectation, showing at the same time the way for future research.

**Keywords:** Time-series Forecasting; Empirical evaluation; NN3 competition; Artificial Neural Networks; Computational Intelligence;

## 1. Introduction

“*Neural networks: forecasting breakthrough or passing fad?*” Chatfield wondered back in (1993); and as of today the question remains largely unanswered. On the one hand, if we consider only the number of publications of artificial Neural Networks (NN) the answer would indicate a breakthrough: motivated by their proven theoretical properties of non-parametric, data driven universal approximation of any linear or nonlinear function, the last two decades have witnessed over 5000 publications in academic journals and conference proceedings on forecasting with NNs across a wide range of disciplines (Sven F. Crone & Preßmar, 2006). In a recent series of surveys on forecasting publications, Fildes et al. note that while the last 25 years have seen rapid developments in forecasting across a broad range of topics, computer intensive methods such as NNs contribute the single largest area of publications in Operational Research (2008) and one of the top 4 in forecasting journals (2006). Their growth in prominence appears to be easily justified: a wealth of publications indicate the competitive or even superior performance of NN, from early publications on single benchmark time series such as the popular airline passenger data (Faraway & Chatfield, 1998; Kolarik & Rudorfer, 1994; Tang & Fishwick, 1993), to subsets of established benchmarks from previous forecasting competitions (Sharda & Patil, 1992; Foster, Collopy, & Ungar, 1992; Hill, O'Connor, & Remus, 1996). Adya and Collopy (1998) found eleven studies that met the criteria for a valid and reliable empirical evaluation, and in 8 of these (73%) NNs were more accurate. However, their evaluative review of the experimental design and the implementation of NNs also raised concerns on the validity and reliability of the results in 37 of 48 studies (77%). For novel algorithms that are not evaluated following a rigorous experimental design, the results from an ex post evaluation (where the test data is known to the authors) may not be sufficiently reliable, but require an objective, unbiased ex ante evaluation in order to determine their true empirical accuracy under varying data conditions.

If, on the other hand, we considered only the empirical post-sample accuracy demonstrated by NNs, a different answer to Chatfield’s question arises. In contrast to their optimistic publications, NNs have failed to provide objective evidence of their forecasting accuracy in large scale empirical evaluations in the form of forecasting competitions. The most renowned empirical investigation conducted - the M3

1  
2  
3  
4  
5  
6 competition (S. Makridakis & Hibon, 2000) - indicated a comparatively poor performance from a single  
7  
8 contestant. Consequently, the performance of NNs on batch forecasting fell far short of their presumed  
9  
10 potential.

11  
12 In contrast, forecasting competitions conducted in computer science and machine learning (e.g.  
13  
14 the Santa Fe (Weigend, 1994) or EUNITE competition (Suykens & Vandewalle, 1998a)) attracted a large  
15  
16 number of NN and CI algorithms. Although these demonstrated preeminent performance of NNs,  
17  
18 algorithms were often not evaluated against statistical methods, using only a single time series (and time  
19  
20 origin) or a small set of heterogeneous time series. This ignored evidence within the forecasting field on  
21  
22 how to design valid and reliable empirical evaluations (see, e.g., Fildes, Hibon, Makridakis, & Meade,  
23  
24 1998), severely limiting the validity and reliability of their findings. As a consequence of the poor  
25  
26 experimental designs, the forecasting community largely ignored the findings.  
27  
28

29  
30 The discrepancy between NNs' preeminent theoretical capabilities, their promising accuracy in  
31  
32 publications on known datasets and some real world application, in contrast to the lack of empirical  
33  
34 accuracy in large scale ex ante evaluations has raised serious concerns in the forecasting domain on their  
35  
36 adequacy for forecasting. As a consequence, Chatfield (quoted in Armstrong, 2006) suspects a positive  
37  
38 bias in NN publications due to a "file-drawer problem" of negative results, leading Armstrong (2006) to  
39  
40 conclude that too much research effort is being devoted to this method. However, to date this scepticism  
41  
42 is founded only on a single contestant entering the last large scale evaluation of automatic forecasting.  
43

44  
45 In order to explore the persisting gap between the theoretical capabilities and empirical accuracy  
46  
47 of NNs, we conducted a forecasting competition to provide valid and reliable empirical evidence on  
48  
49 accuracy, evaluate and disseminate potential progress in modelling NNs and to determine the conditions  
50  
51 under which different algorithms perform well. Our motivation to conduct yet another competition draws  
52  
53 upon the same argument as the original M-competition by (S. Makridakis, et al., 1982): a full decade has  
54  
55 passed since the start of the M3 competition, a decade that has seen the development of extended NN  
56  
57 paradigms, theoretical advances in methodologies on specifying NNs and a range of novel computer  
58  
59 intensive algorithms in CI becoming available for forecasting. In addition, it has seen substantial progress  
60

1  
2  
3  
4  
5  
6 in information technology (IT) that may facilitate the application of earlier algorithms and novel additions  
7  
8 in large scale forecasting competitions that were infeasible before due to limited computational resources.  
9  
10 As new alternatives exist, choices on selecting and using appropriate forecasting methods need to be  
11  
12 revisited.  
13

14  
15 To evaluate progress in NNs, and to allow a comparison to the original M3-contestants over time,  
16  
17 we utilised a subset of 111 monthly industry time series taken from the previous M3-data for which all  
18  
19 original predictions were available. The dataset contains a balanced sample of seasonal and non-seasonal,  
20  
21 short and long time series in order to evaluate the conditions under which an algorithm performs well.  
22  
23 The competition was open to all methods of NNs and CI. To limit biases we also allowed novel statistical  
24  
25 methodologies and newer software releases to participate as benchmarks. NN3 attracted 59 submissions,  
26  
27 making it the largest competitions in CI and forecasting to date. Results were evaluated using multiple  
28  
29 error metrics, including the original symmetric mean absolute percent error (sMAPE), mean absolute  
30  
31 scaled error (MASE) as proposed by Hyndman and Koehler (2006) and two non-parametric tests  
32  
33 employed by Koning et al. (2005) in a follow-up analysis of the M3-data: Analysis of the Mean (ANOM)  
34  
35 and Multiple Comparisons to the Best method (MCB). In short, we attempted to take into consideration  
36  
37 all recommendations on how to conduct a valid and reliable empirical evaluation while balancing effort  
38  
39 and resources to maximise the number of submissions in order to receive a more representative sample of  
40  
41 algorithms. As the competition followed the original design of the M3, it was launched under the name  
42  
43 *NN3 competition*. This paper summarises its findings, discusses the results of the experiments and  
44  
45 implications for future research.  
46  
47

48  
49 The rest of the paper is structured as follows: section two discusses prior forecasting competitions  
50  
51 in forecasting and CI, with a discussion on the relevance of forecasting competitions to derive empirical  
52  
53 evidence, guidelines for their setup and contrasts the findings of major competitions in forecasting and CI  
54  
55 in order to justify the rationale for another one; As competitions in CI have not followed similar designs,  
56  
57 the best practices derived in the experimental design of forecasting competitions are explored in more  
58  
59 detail in order to disseminate them to a interdisciplinary readership. Sections three and four describe the  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

setup and the results of the empirical evaluation, taking these best practices into consideration. Section five provides a brief discussion of the most important findings followed by the conclusions and implications for future research.

## 2. Evidence from Competitions in Forecasting and Computational Intelligence

In the absence of a universal (theoretical or empirical) dominance of a single ‘best method’, competitions are an established means to provide objective evidence on the empirical ex ante accuracy of forecasting methods, and to guide rational choices between algorithms and methodologies for a given set of data conditions. Forecasting competitions have received substantial attention and initiated stimulating discussions within the academic forecasting community, opening up new areas of academic research (e.g. model selection and evaluation) and leading to best practices on valid and reliable competitions and experimental designs (Ord, Hibon, & Makridakis, 2000) An overview and discussion of the impact of empirical evaluations may be found in Fildes et al. (Fildes & Ord, 2002; Fildes & Makridakis, 1995). In contrast, competitions on time series prediction conducted in other domains, including computer science, machine learning, engineering and CI, have largely pursued different experimental designs that have ignored best practices on conducting competitions, limiting their validity and reliability. In order assess the empirical evidence provided in each field to date, and contrast the lack of dissemination of algorithms, applications and best-practices across both domains, we briefly summarize existing competitions in forecasting and CI and provide an overview in table 1.

In forecasting research, a series of competitions have been conducted that have received substantial interest. Drawing upon criticism on earlier competitions on time series data (Reid, 1969, 1972; Newbold & Granger, 1974; Groff, 1973; S. Makridakis & Hibon, 1979) Makridakis et al. conducted a series of enlarged forecasting competitions where experts could submit predictions of their preferred algorithms, starting with the M-Competition on two datasets of 1001 or 111 time series, which took into account suggestions made at a meeting at the Royal Statistical Society (S. Makridakis, et al., 1982). A reduced subset of data was offered to allow participation of algorithms which required time and cost

1  
2  
3  
4  
5  
6 intensive manual tuning through experts (e.g. ARIMA models required more than one hour per time  
7 series). The subsequent M2-competition (Spyros Makridakis, et al., 1993) focussed on non-automatic,  
8 real time judgmental forecasts of 23 time series, and hence holds less relevance for our quantitative  
9 competition design. None of the earlier competitions attracted any submissions of NNs or CI methods, as  
10 these did not emerge until the late 1980s, e.g. in the case of NNs through the (re-)discovery of the  
11 Backpropagation algorithm (Rumelhart, Hinton, & Williams, 1994) and others such as CART (Breiman,  
12 1984) although algorithms of Fuzzy Logic (Zadeh, 1965) and Evolutionary Computation (Fogel & Fogel,  
13 1994) had already been developed. In 1998 the popular M3-Competition evaluated the accuracy of 24  
14 algorithms on 3003 univariate empirical time series of historical data (S. Makridakis & Hibon, 2000), the  
15 largest dataset ever used in such a competition. The time series were selected from various domains of  
16 micro- and macroeconomic, industrial, financial and demographic activity, and from different time  
17 frequencies (yearly, quarterly and monthly data) in order to cover a wide variety of time series structures  
18 and different data conditions. All methods were implemented by academic experts and commercial  
19 software providers, leading to the most representative ex ante evaluation of forecasting methods to date.

20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100  
101  
102  
103  
104  
105  
106  
107  
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161  
162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000

Across all time series, two methods on average outperformed all other methods: the software expert system Forecast Pro using automatic model selection and parameterisation of Exponential Smoothing (ES) and ARIMA models (Goodrich, 2000), and Theta, a decomposition approach combining exponential smoothing and regressing around a dampened trend line (Assimakopoulos & Nikolopoulos, 2000). Further statistical analysis by Koning *et al.* (Koning, et al., 2005) has revealed statistically significant results for a group of four methods, also including Rule Base Forecasting (Adya, Armstrong, Collopy, & Kennedy, 2000) and Comb S-H-D an equal weighted combination of single, Holt's linear trend and Gardner's dampened trend ES methods (computed by Hibon) in the top performers.

Despite initial interest by various CI researchers, only one competitor finally submitted results to the competition using a NN methodology (Balkin & Ord, 2000). However, their fully automated methodology AutomatANN showed only moderate performance in comparison to the majority of the twenty statistical approaches, and was not ranked among the top performers (Table 15, S. Makridakis &

1  
2  
3  
4  
5  
6 Hibon, 2000). The low participation of CI-approaches had been attributed to the high computational costs  
7  
8 in building and parameterising these methods for each time series, but also the lack of methodologies that  
9  
10 would allow automation beyond manual tuning by a human expert. The poor performance however was  
11  
12 neither expected, nor explained sufficiently.  
13

14  
15 The conclusions drawn from prior M-competitions(S. Makridakis, et al., 1982; Spyros  
16  
17 Makridakis, et al., 1993) have been confirmed in the M3-competition (see S. Makridakis & Hibon, 2000),  
18  
19 verified through follow-up studies (see, e.g., Fildes, 1992) and extended to provide additional insights  
20  
21 (Fildes, et al., 1998): (H1) *the characteristics of the data series are an important factor in determining*  
22  
23 *relative performance between methods*, (H2) *Accuracy of a method depends upon the length of the*  
24  
25 *forecasting horizon*, (H3) *Relative performance ranking of methods varies with the accuracy measure*;  
26  
27 (H4) *Sampling variability of performance measures renders comparisons based on single time series*  
28  
29 *unreliable; comparisons based on multiple time origins are recommended* (H5) *Accuracy of a*  
30  
31 *combination of predictions performs well and often outperforms the individual methods*; (H6)  
32  
33 *Sophisticated methods do not necessarily provide more accurate forecasts than simpler ones*; and.  
34  
35 Consequently, valid competitions have developed a rigorous design, including the use of a representative  
36  
37 number of time series, rolling origin design, the use of multiple robust error metrics, the comparison  
38  
39 against established (statistical) benchmark algorithms, and the analysis of the data conditions under which  
40  
41 a method performs well (Tashman, 2000) in order to obtain valid and reliable results. Conclusion H6  
42  
43 seems of particular relevance, as NNs and other computer intensive methods - just as sophisticated  
44  
45 statistical algorithms like ARIMA before them - do not guarantee enhanced forecasting performance  
46  
47 merely by their proven capabilities or theoretical features, and requires evaluation against simpler  
48  
49 benchmarks. No competitions of similar scale have been conducted since (with the exception of the MH-  
50  
51 competition on transportation data of varying time frequency, conducted in 2007 by Hibon, Young and  
52  
53 Scaglione, and Athanasopoulos et al. (2009) on tourism forecasting (which is rather a conventional  
54  
55 empirical study, as no call for participation was issued), of which results have not been published yet.  
56  
57 This leaves the M3 as the last large scale evaluation in the forecasting domain to date, and explains the  
58  
59  
60



1  
2  
3  
4  
5  
6 impact and prominence of the disappointing yet unchallenged results of NN in empirical forecasting,  
7  
8 based upon the single entry of the only CI-contestant AutomatANN (Balkin & Ord, 2000).  
9

10  
11 On the other hand, the findings of the M3 were not representative of NN. Despite a myriad of  
12  
13 published NN methodologies, only a single methodology was evaluated, limiting the representativeness  
14  
15 of the results for the class of NN (which encompasses a variety of feedforward and recurrent  
16  
17 architectures) and for CI as a whole. Also, the M3 attracted no interest from the computer sciences,  
18  
19 engineering or machine learning community where CI and other approaches of artificial intelligence had  
20  
21 been advanced for years, introducing a sample selection bias of algorithms and dissemination of results  
22  
23 (an omission caused by disseminating the CfP only through media of the International Institute of  
24  
25 Forecasters (IIF), i.e. IJF and the ISF conference). Consequently, the poor performance of a single NN-  
26  
27 approach in M3 cannot be considered representative of the whole class of algorithms.  
28

29  
30 Furthermore, almost a decade has passed since M3, so that results may no longer reflect today's  
31  
32 capabilities of NN. Evidence for substantial theoretical progress in NN exists, both in forecasting single  
33  
34 time series (Preminger & Franck, 2007; de Menezes & Nikolaev, 2006; Terasvirta, van Dijk, & Medeiros,  
35  
36 2005) and on representative sets of empirical time series (see, e.g., Liao & Fildes, 2005; Zhang & Qi,  
37  
38 2005) applying methodologies for fully automated applications. These have not yet been evaluated in an  
39  
40 objective empirical competition. Lastly, today's computational power is far superior to that available in  
41  
42 1997, when automated NNs were run for the M3 competition, also apparent in the expanding community  
43  
44 regularly applying computationally intensive methods, which may enable a much wider participation.  
45  
46 Thus, the results of the M3 may no longer be deemed representative. However, in the absence of more  
47  
48 recent forecasting competitions its critical findings to NN stand unchallenged.  
49

50  
51  
52  
53 Outside the forecasting domain, competitions have been equally popular to determine the  
54  
55 predictive accuracy of algorithms, and many more recent than the M3. Regular data mining competitions  
56  
57 have been conducted, albeit focussed on classification tasks, including the annual competitions at the  
58  
59 KDD conference, attracting over 1,000 contestants in 2008, or the recently closed Netflix competition  
60

1  
2  
3  
4  
5  
6 (www.netflixprice.com) in predicting movie choices, attracting 44,014 submissions (by awarding a price-  
7 money of 1 Million US\$). As in forecasting, competitions on classification with CI generally follow a  
8 rigorous experimental design, adhere to established best practices for a valid and reliable experimental  
9 designs, and often address sophisticated modelling questions such as if domain knowledge allows better  
10 predictive accuracy than agnostic prediction (Guyon, Saffari, Dror, & Cawley, 2008) or to what extent in  
11 sample performance can be generalised for out of sample accuracy (Cawley, Janacek, Haylock, &  
12 Dorling, 2007).

13  
14  
15  
16  
17  
18  
19  
20  
21 In contrast, only a small number of competitions in the CI-domain were dedicated to time series  
22 data, often of small scale. A discussion of all competitions and their contributions is beyond this paper,  
23 but we will outline those most influential to identify fundamental differences in the experimental design.  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
The *Time Series Prediction and Analysis Competition*, organised by Weigend and Gershenfeld (1994)  
under the auspices of the Santa Fe Institute, was the first dedicated competition in CI to evaluate the  
ability of NNs in forecasting using a variety of exclusively nonlinear datasets recorded along time.  
Datasets were highly heterogeneous and included time series of a physics experiment recording  
oscillations and structural breaks of a NH5-Laser, and multivariate time series of tick-by-tick currency  
exchange rates, astrophysical data of light fluctuations from a white star, physiological data from a patient  
with sleep apnea, and music from Bach's last (unfinished) Fuge. given the heterogeneity of data  
conditions, datasets consisting only of a single time series, that most participants predicted only one of the  
datasets from a single origin (instead of - at least - all series), and that no statistical benchmarks were  
evaluated, the comparative work undertaken remains rudimentary and does not provide sufficient  
evidence to draw conclusions on the nonlinear algorithms' accuracy (S. Makridakis, 1994). The lack of  
validity seems particular disappointing considering that the authors were aware of the design and findings  
of the M-competitions, and that the late Clive Granger served on the competition's advisory board (Jadiz,  
1995).

The largest CI-competition on time series to date was organised by Suykens and Vanderwalle  
(Suykens & Vandewalle, 1998a) in 2000 for the European Network on Intelligent Technologies for Smart

1  
2  
3  
4  
5  
6 Adaptive Systems (EUNITE, [www.eunite.org](http://www.eunite.org) no longer online) attracted 24 submissions from 16  
7  
8 contestants (only 29% of those 56 registered to compete). It evaluated the accuracy in predicting one time  
9  
10 series of maximum electrical load 31 days into the future from a single origin, using 2 years of half-  
11  
12 hourly electrical load data. Data was provided by the Eastern Slovakian Electricity Corporation including  
13  
14 explanatory variables of past temperature and holidays. The best contestant used Support Vector  
15  
16 Regression (Chen, Chang, & Lin, 2004) to outperform contestants of CI and one 'statistical' contender  
17  
18 using regression on decomposed time series components. Although all algorithms were published in a  
19  
20 monograph (Sincák, Strackeljan, Kolcun, Novotný, & Szathmáry, 2002) it has received limited attention  
21  
22 outside the electrical load literature.  
23  
24

25 A range of smaller competitions has been run at conferences on computational intelligence, often  
26  
27 without publications, including the Competition on Artificial Time Series (CATS) using multiple samples  
28  
29 of synthetic data (Lendasse, Oja, Simula, & Verleysen, 2007), held at the 2004 IEEE International Joint  
30  
31 Conference on Neural Networks (IJCNN), the Predictive Uncertainty Competition at the 2006 IJCNN  
32  
33 (Cawley, et al., 2007) on environmental data, the 2003 and 2006 Business Intelligence Cup on predicting  
34  
35 time series of sugar and retail sales, organised by Richard Weber held at the IEEE Latin-American  
36  
37 Summer School on Computational Intelligence (EVIC), the 2001 ANNEXG competition on river stage  
38  
39 forecasting (Dawson, et al., 2005) held at the 2002 BHS National Hydrology Symposium (the 2005 re-run  
40  
41 attracted no competitors) and the KULeuven competition by Suykens and Vandewalle (Suykens &  
42  
43 Vandewalle, 1998b) held at the International Workshop on Advanced Black-Box Techniques for  
44  
45 Nonlinear Modeling in 1998 on synthetic data (see also McNames, Suykens, & Vandewalle, 1999). Table  
46  
47 1 provides a structured summary of prior time series competitions, both in forecasting and CI, and  
48  
49 contrasts differences in the experimental design in both domains to assess their contributions.  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Table 1: Competitions design in Forecasting and Computational Intelligence

Competition Name	Data Properties			# of algorithms		Conditions evaluated <sup>▲</sup>			
	# of series	# of observ.	Data type	Statistics	NN & CI	Multiple Metrics*	Multiple horizons*	Multiple types*	Time Frequency <sup>♦</sup>
Empirical Accuracy <sup>1</sup>	111	?	univariate	22	0	X	X	X	Y, Q, M
M1 <sup>2</sup>	1001	15 to 150	univariate	24	0	X	X	X	Y, Q, M
M3 <sup>3</sup>	3003	20 to 144	univariate	24	1	X	X	X	Y, Q, M
MH/Transport <sup>4</sup>	278	19 to 1502	univariate	3-10	1	X	X	X	Y, Q, M W, D, H
Santa Fe	6	1,000 to 300,000	2 univariate 4 multivariate	0	14	-	-	-	Synthetic
KULeuven <sup>6</sup>	1	2000	Univariate	0	17	-	-	-	Synthetic
2001 EUNITE <sup>7</sup>	1	35,040	multivariate	1	24	-	-	-	30m
ANNEXG <sup>8</sup>	1	1,460	Multivariate	0	12	-	-	-	360m
BI Cup 2003 <sup>9</sup>	1	365	Multivariate	0	10	-	-	-	D
CATS 2005 <sup>10</sup>	1	4,905	univariate	0	25	-	-	-	Synthetic
Predictive Uncertainty, <sup>11</sup>	4	380 to 21,000	1 univariate 3 multivariate	0	20	-	-	-	Synthetic D, 3D
BI Cup 2006 <sup>12</sup>	1	1325	1 univariate	0	?	-	-	-	15m

<sup>1</sup> (S. Makridakis & Hibon, 1979); <sup>2</sup> (S. Makridakis, et al., 1982); <sup>3</sup> (S. Makridakis & Hibon, 2000); <sup>4</sup> unpublished; <sup>5</sup> (Weigend, 1994); <sup>6</sup> (Suykens & Vandewalle, 1998a); <sup>7</sup> (Suykens & Vandewalle, 1998b); <sup>8</sup> (Dawson, et al., 2005); <sup>9</sup> unpublished; <sup>10</sup> (Lendasse, et al., 2007); <sup>11</sup> (Cawley, et al., 2007); <sup>12</sup> unpublished

\* X indicates the use of multiple error metrics, multiple forecasting horizons and data types, - indicates the absence

♦ Y=yearly data; Q = quarterly data; M=monthly data; W=weekly data; D=daily data; H=Hourly data; m=minutes; ? indicates non disclosed information

Few similarities emerge, bare one: both domains favour and evaluate almost exclusively its preferred family of algorithms: forecasting competitions evaluate only statistical methods (and expert systems configuring these), with the exception of the sole NN contender at the M3 and the unpublished MH-competitions, while CI-competitions have not evaluated statistical algorithms at all.

Another commonality shared across domains remains the evaluation across a single time origin (due to the time consuming nature of stepwise release of data), with the exception of CATS, which provided 5 origins of a continuous time series simultaneously (and the M2 not mentioned here). However, as the test data was withheld as 5 gaps of one time series, contestants used heuristics to connect end and beginning of the series in addition to fore- and backcasting in the sense of imputing consecutive missing

1  
2  
3  
4  
5  
6 values rather than forecasting. Furthermore, as the data was synthetic and not empirically motivated, the  
7  
8 competition holds limited relevance to assess accuracy in time series prediction.  
9

10 More noticeably, differences and discrepancies in the design of all CI-competitions become  
11 evident, which seriously impairs their contribution. While all competitions of the forecasting domain used  
12 representative sample sizes of hundreds or even thousands of time series, CI-competitions evaluated  
13 accuracy only on a single time series. Those competitions that evaluate multiple time series, such as the  
14 Santa Fe and the Predictive Uncertainty competition, did so in distinct groups with only one series per  
15 category, again limiting any generalisation of their findings. Had the same algorithm been used across  
16 multiple similar series, datasets or competitions, it would have permitted somewhat replicable results. In  
17 contrast, the same authors applied different methodologies for each dataset, even within a competition,  
18 leading to distinctly different models and preventing any comparisons.  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

29 Also, none of the CI-competitions compare results against established benchmark methods, be  
30 they Naive methods, simple statistical benchmarks or non-statistical methods in the same family of  
31 algorithms, e.g. a simple single layer MLP with default parameters to compete against a more  
32 sophisticated architecture. We therefore conclude, that the recommendations on the design of empirical  
33 evaluations developed in forecasting, in particular the use of multiple similar time series, have been  
34 ignored by the CI community. Makridakis' (2000) original criticism holds: just as theoretical statisticians  
35 before them, researchers in NN have concentrated efforts on building more sophisticated models without  
36 regard to assessing their accuracy and objective empirical verifications, successfully ignoring the strong  
37 empirical evidence of the M-competition and the ground rules it has laid out on how to assess them. This  
38 substantially limits the validity and reliability of the evidence from CI-competitions, which can not  
39 challenge the authority of the earlier M3-competition which failed to show benefits of NNs in terms of  
40 accuracy.  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53

54 With competitions in both domains limited in their coverage of algorithms, results of the M3  
55 competition not, or at least no longer representative of CI, and more recent competitions in CI unreliable,  
56 the gap between the theoretical capabilities and empirical accuracy of NNs remains unexplored. In order  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 to evaluate potential progress in the development of NN and CI-approaches, a new competition seemed  
7  
8 most suitable to provide valid and reliable empirical evidence on their accuracy, the conditions under  
9  
10 which different algorithms perform well, and to disseminate potential progress in modelling NNs. For the  
11  
12 sake of consistency, it seemed a natural choice to use the original setup and a homogeneous subset of the  
13  
14 M3 data in the form of a replication, which will be discussed in detail in the next session.  
15  
16  
17  
18

19 In reviewing table 1, we note a further discrepancy in data conditions explored. Forecasting  
20  
21 competitions have been focussed exclusively on low time series frequencies of yearly, quarterly, or at  
22  
23 most monthly data in a univariate context. Although this adequately reflects the persisting theme of  
24  
25 operational forecasting set out by Makridakis' series of M-competitions (S. Makridakis, et al., 1982), it  
26  
27 does not allow us to generalise these findings to unseen data conditions, in particular towards high-  
28  
29 frequency data of weekly, daily, hourly or shorted time intervals (and the resulting longer time series) on  
30  
31 which NN have been evaluated (represented in the dissimilar data conditions of CI competitions). It  
32  
33 seems Armstrong's (2006) criticism of NN is not only limited in evidence due to a single contestant of  
34  
35 M3, but more importantly limited due to a substantial omission of empirical data conditions, for which -  
36  
37 following his arguments - no evidence exists. As the omitted conditions represent those on which NN are  
38  
39 regularly - and successfully - employed in practice, it might also yield an explanation for the simultaneous  
40  
41 scepticism and euphoria re NN between disciplines, and provide the motivation to close the gap in more  
42  
43 representative competitions on novel data conditions.  
44  
45

### 46 **3. Design and Organisation of the *NN3 competition***

#### 47 **3.1 Objectives**

48  
49  
50 Following the rationale provided, we sought to explore the current forecasting performance of  
51  
52 NN and CI methods. The M-competitions explicitly focussed on a particular set of data conditions, which  
53  
54 Makridakis coined in the context of forecasting for operations. To assess progress in relation to M3, we  
55  
56 will keep this tradition and constrain our competition to the operations context of monthly industry data,  
57  
58 although other conditions for forecasting with NNs exist that may promise different results.  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

The NN3 competition was designed both to replicate and to extend the M3 competition. As a replication, the NN3 will utilise the data, experimental setup and original submissions from M3, and evaluate the working hypothesis of earlier competitions (see section 2) to challenge or confirm prior findings. In addition, NN3 represents an extension towards more methods/researchers from NN and CI in order to assess progress in accuracy and to overcome limitations in the M3 representativeness. Previous forecasting competitions have led to an established 'research methodology' for a systematic, valid and reliable design of future competitions, which we attempted to follow. We will briefly review these design choices, the datasets and conditions, accuracy metrics, methods and benchmarks, and the process in which NN3 was conducted, in order to allow verification of the experimental design and to disseminate this knowledge to the CI community.

### 3.1 Datasets, Working Hypothesis and Data Conditions

The M3 dataset yielded substantial insights, but proved challenging for CI-methods: the sample of 3,003 time series was large given the computational resources available in the 1990s, and the heterogeneity of time series frequencies and data domains required multiple candidate methodologies (and human intervention at many stages), which limited automation and may have prevented many computationally intensive methods to participate. In order to attract a representative number of contestants and algorithms we sought to limit the number of time series used and heterogeneity of data conditions (and resulting insights), yet enough to derive reliable results. A set of 111 time series was selected randomly from the M3 subset of monthly industry data, representative of the M-competition's original focus of forecasting for operations. Data from a single time frequency was chosen in order to limit the competition's complexity to a single methodology for monthly data. We also hoped that a sample would further mask the origin of the NN3 competition data and prevent biases in results through prior knowledge.

Four working hypothesis were considered in the evaluation. To determine the degree of automation or manual tuning required, and to address prevailing concerns on the computational demands of predicting a large number of time series with NN, we allowed participation on two (disguised) datasets

1  
2  
3  
4  
5  
6 of different size. Contestants were asked to predict either a reduced dataset of 11 or a complete set of 111  
7 time series (which included the reduced set) as accurately as possible. As a fully automated methodology,  
8 as required for operational forecasting, could be applied to datasets of larges size, more submissions on  
9 the reduced dataset would indicate the need for manual tuning, limitations of the automation or extremely  
10 computational intensive approaches, indicating the need for further research in methodologies.  
11  
12  
13  
14  
15

16  
17 A second working hypothesis seeks to assess the relative accuracy of NN and statistical  
18 approaches on longer forecasting horizons, where statistical algorithms have outperformed NN in past  
19 studies (Hill, et al., 1996). Each contestant is required to predict multiple forecast  $y_{t+h}$  of  $h = (1, \dots, 18)$   
20 steps into the future, which are later analysed for short (1-3 months), medium (3-12 months) and long  
21 (13-18 months) forecasting horizons to assess differences in the results.  
22  
23  
24  
25  
26

27 Two further working hypothesis address the data conditions under which different methods  
28 perform well. First, following the widespread belief that NNs are data hungry and require long time  
29 series, balanced stratified samples were taken by time series length  $n$ , resulting in 50 long ( $n > 100$ ) vs. 50  
30 short time series ( $n < 50$ ). Second, to evaluate recent publications which conclude that NNs cannot  
31 forecast seasonal time series (Nelson, Hill, Remus, & O'Connor, 1999; Zhang & Qi, 2005; Curry, 2007)  
32 stratified samples were taken to reflect time series patterns of 50 seasonal vs. 50 non-seasonal time series  
33 (as per the original M3 classification). Series with structural breaks in the test set were manually  
34 identified and excluded.  
35  
36  
37  
38  
39  
40  
41  
42  
43

44 The sample sizes were guided by the objective to derive (statistically) valid and reliable results  
45 for each data condition from as small a dataset as possible, which created a lower bound of 25 time series  
46 in each cell (i.e. short-seasonal, long-seasonal, short-non-seasonal, and long-non-seasonal), resulting in  
47 100 series as a core to the complete set. The reduced dataset contained 11 time series which we classified  
48 as difficult to forecast, 4 of which were seasonal and the remaining 7 non-seasonal (including outliers,  
49 structural breaks), to identify if non-automated methodologies are capable of predicting across different  
50 data conditions. Table 1 summarises the time series conditions of both datasets.  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



**Table 1.** NN3 Datasets with data conditions of time series length and seasonality.

	Complete Dataset		Reduced Dataset		SUM
	Short	Long	Normal	Difficult	
Non-Seasonal	25 (SN)	25 (LN)	4 (NN)	3 (DN)	57
Seasonal	25 (SS)	25 (LS)	4 (NS)	-	54
SUM	50	50	8	3	111

The conditions within the reduced dataset were not meant to be statistically explored due to the limited amount of series (3-3-4) that would result in insignificant results. Nonetheless, findings on the reduced dataset would still yield results of increased validity and reliability in comparison to previous CI competitions using only a single time series.

### 3.2 Evaluation and Error Metrics

In order to evaluate the performance of the NN3 submissions, and to ensure consistency with the results of the M3-competition, we employed three metrics used in the M3 competition (all averaged across all series): Symmetric Mean Absolute Percent Error (sMAPE), Median Relative Absolute Error (MdRAE) and Average Ranking (AR) (S. Makridakis & Hibon, 2000). We also estimated two non-parametric tests proposed by Koning *et. al.* (2005) in a follow up analysis: Analysis of the Mean (ANOM) and Multiple Comparisons to the Best (MCB), both using the average ranking as criterion. Finally, in order to align with the current literature, we have calculated the Mean Absolute Scaled Error (MASE) as proposed by Hyndman and Koehler (2006). In order to ensure a consistent computation of errors, we collaborated with one of the original investigators of the M3 competition who computed all metrics as in the original competition.

Average sMAPE was announced beforehand as the metric used to determine the “winner”, in order to allow those CI-methods capable of using non squared error cost functions to align their approaches with the final criterion (see, e.g., the discussion with Zellner (1986) following the M3). Despite the shortcomings of sMAPE (Goodwin & Lawton, 1999), it was chosen as it served as the

1  
2  
3  
4  
5  
6 primary criterion in the M3 competition and to make the NN3 results accessible to practitioners, where  
7  
8 the predominant error metric is MAPE. As the NN3 time series contained no zero, negative or small  
9  
10 values, we anticipate only limited potential for bias.  
11

### 12 **3.2 Methods and Benchmarks**

13  
14 The competition invited contributions from all areas of CI, including all NN paradigms and  
15  
16 architectures, Support Vector Regression, Fuzzy Logic, Evolutionary and genetic algorithms, and hybrid  
17  
18 methods utilising any kind of CI. In an attempt not to bias results towards novel NN-methods, we also  
19  
20 allowed novel statistical methodologies and newer software releases to be evaluated as benchmarks,  
21  
22 further extending the representativeness of NN3.  
23  
24

25 We personally invited experts in statistical forecasting methods and commercial software vendors  
26  
27 in order to ensure participation of the latest releases of those that had performed well in the original M3-  
28  
29 competition, but with limited success. We are grateful for submissions from Eric Stellwagen of Business  
30  
31 Forecasting Systems, BFS, applying the latest version of the expert system ForecastPro (B03), and Dave  
32  
33 Reilly of Autobox (B05), applying the latest version of the expert system for ARIMA- and transfer  
34  
35 function modelling, and Tucker McElroy who submitted prediction of the Census X12 method (B6).  
36  
37

38 In order to assess progress in NN modelling since the M3, the NN3 submissions needed to be  
39  
40 compared to the original M3 submission of AutomatANN (Balkin & Ord, 2000, B00). Given the identical  
41  
42 experimental setup and data taken from M3, the collaboration with one of the original investigators of the  
43  
44 M3 competition allowed us to retrieve the 111 original predictions submitted to M3 and compare them  
45  
46 directly with those of the NN3 contestants. Further to AutomatANN, five statistical Benchmarks used in  
47  
48 the M3 were recalled, including the Naïve-1 method (B04), three variants of Brown's single (B14), Holt's  
49  
50 linear trend (B15) and Gardner's dampened trend ES (B16), and their combination to Comb S-H-D (B17).  
51  
52 Predictions for Theta (B7) were recomputed by the organisers, using an identical setup as in the M3  
53  
54 competition.  
55

56 In addition, we computed various CI-benchmarks to provide additional levels of comparison of  
57  
58 the entries, including a Naïve Support Vector Regression approach (S. F. Crone & Pietsch, 2007, B01)  
59  
60

1  
2  
3  
4  
5  
6 and a Naïve Multilayer Perceptron (B02), which replicate novice model building mistakes as a lower  
7  
8 bound. A novel NN-extension of the successful Theta-Method named Theta-AI (B08) by Nikolopoulos  
9  
10 and Bougioukos that determined optimal nonlinear weights for the Theta-lines was withdrawn in order  
11  
12 not to bias results, as it was based on Theta that was known a-priori to perform well on the NN3 data.  
13

### 14 **3.4 Process of organising the Competition**

15  
16 The competition design and feasibility was pre-tested in a small scale trial competition using two  
17  
18 time series ( held at the 2005 ISF, San Antonio, USA), which facilitated feedback by 9 contestants and  
19  
20 external experts, including a panel of IIF judges for a grant. The NN3 competition was first announced at  
21  
22 the ISF 2006 in Santander and was open for eight months from October 2007 to May 2008. Contestants  
23  
24 were required to submit predictions and a full description of their methodology, which are both published  
25  
26 on the competition website in order to facilitate replication. Contestants could withhold their identity prior  
27  
28 to disclosing the final results, in order to limit negative publicity for software vendors and participants.  
29  
30 Following submission, each methodology was classified to distinguish CI-contenders eligible to “win” the  
31  
32 competition (identified by consecutive IDs C01-C59) from other submissions that would serve as  
33  
34 benchmarks: CI benchmarks (B00-B02), statistical benchmarks including forecasting packages (B03-  
35  
36 B08), novel statistical methods submitted as benchmarks (B09-B13), and the original ES variants of M3  
37  
38 (B14-B17). Some contributors requested to withhold their identity, but their results are included in the  
39  
40 tables with their original submission IDs to ensure consistency with previously disclosed results.  
41  
42  
43

44 In order to limit sample selection biases in the participation through timing, location and audience  
45  
46 of the conference where the competition was run, multiple special sessions were advertised and conducted  
47  
48 at conferences throughout 2007, and across the domains of forecasting, computational intelligence and  
49  
50 electrical engineering, data mining and machine learning. These included the 2007 International  
51  
52 Symposium on Forecasting (ISF'07), in New York, USA, the 2007 IEEE International Joint Conference  
53  
54 on Neural Networks (IJCNN'07) in Orlando, USA, and the 2007 International Conference in Data Mining  
55  
56 (DMIN'07) in Las Vegas, USA. The call for papers was disseminated via various email-lists, websites,  
57  
58 online communities and newsletters across disciplines.  
59  
60

## 4. Results of the NN3 competition

### 4.1 Results on the complete dataset

The competition attracted 46 contestants using NN and CI-methods and 17 benchmark methods, making it the largest empirical evaluation in NN, CI and forecasting to date.

All 46 contenders submitted predictions for the reduced set of the 11 time series, but only 22 contenders predicted the 111 time series of the complete set. With less than half of the contestants (47%) able to predict more than 11 series, it provides evidence that the need for manual tuning and human intervention still dominates most methodologies. This reflects our experience, both in academia and practice, and is supported by the lack of commercial CI-software for automatic time series forecasting.

Table 2 presents the names of the NN3 contestants, a summary of the algorithm and a consecutive ID (assigned during the competition) that provided forecasts for the 111 series of the complete dataset. A discussion of all submissions is not feasible within the scope of this paper, so we will limit our discussion to those methods that have stood out on some or all of the data conditions we analysed. Detailed descriptions of all methodologies, including the 24 contenders that provided forecasts only for the 11 series of the reduced dataset, are available on the NN3 competition website [www.neural-forecasting-competition.com](http://www.neural-forecasting-competition.com) for further reading.

**Table 2.** NN3 participant IDs, names and method descriptions for the complete dataset of 111 series

Code	Classification	Name	Description
C03	NN/CI Contender	Flores, Anaya, Ramirez, Morales	Automated Linear Modeling of Time Series with Self Adaptive Genetic Algorithms
C11	NN/CI Contender	Perfilieva, Novak, Pavliska, Dvorak, Stepnicka	Combination of two techniques: fuzzy transform and perception-based logical deduction
C13	NN/CI Contender	D'yakonov	Simple kNN-Method for Times Series Prediction
C15	NN/CI Contender	Isa	Growing fuzzy inference neural network
C17	NN/CI Contender	Chang	K-Nearest-Neighbor and Support-Vector Regression
C20	NN/CI Contender	Kurogi, Koyama, Tanaka, Sanuki	Using First-Order Difference of Time Series and Bagging of Competitive Associative Nets
C24	NN/CI Contender	Abou-Nasr	Recurrent Neural Networks
C26	NN/CI Contender	de Vos	Multi-Resolution Time Series Forecasting Using Wavelet Decomposition
C27	NN/CI Contender	Illies, Jäger, Kosuchinas, Rincon, Sakenas, Vaskevcius	Stepping forward through echoes of the past: forecasting with Echo State Networks
C28	NN/CI Contender	Eruhimov, Martyanov, Tuv	Windowed wavelet decomposition and Gradient Boosted Trees
C30	NN/CI Contender	Pucheta, Patino, Kuchen	Neural Networks-Based Prediction Using Long and Short Term Dependence in the Learning Process
C31	NN/CI Contender	Theodosiou, Swamy	A hybrid approach: Structural Decomposition, Generalised Regression Neural Networks and Theta model
C36	NN/CI Contender	Sorjamaa, Lendasse	A non-linear approach (Self-Organized Maps) combined with a linear one (Empirical Orthogonal Functions)
C37	NN/CI Contender	Duclos-Gosselin	Fully-recurrent neural network learned with M.A.P (Bayesian), Leventberg and genetic algorithm.
C38	NN/CI Contender	Adeodato, Vasconcelos, Arnaud, Chunha, Monteiro	Multilayer Perceptron Networks
C44	NN/CI Contender	Yan	Multiple-Model Fusion for Robust Time-Series Forecasting
C46	NN/CI Contender	Chen, Yao	Ensemble Regression Trees
C49	NN/CI Contender	Schliebs, Platel, Kasabov	Quantum Inspired Feature Selection and Neural Network Models
C50	NN/CI Contender	Kamel, Atiya, Gayar, El-Shishiny	A Combined Neural Network/Gaussian Process Regression Time Series Forecasting System
C51	NN/CI Contender	Papadaki, Amaxopolous	Dynamic Architecture for Artificial Neural Networks
C57	NN/CI Contender	Corzo, Hong	Global neural networks ensembles with M5 prime model trees
C59	NN/CI Contender	Beliakov & Troiano	Time series forecasting using Lipschitz optimal interpolation
B00	NN/CI Benchmark	Automat NN	Original M3 Submission for the M3 competiiton by Balkin & Ord
B01	NN/CI Benchmark	Naive SVR	A naïve Support Vector Regression forecasting approach by Crone & Pietsch
B02	NN/CI Benchmark	Naive MLP	A naïve Multiple Linear Perceptron by Crone
B03	Stat. Benchmark	ForecastPro	The Expert method of the renown software. Version XE 5.0.2.6.
B04	Stat. Benchmark	Naive	The Naïve method without any seasonality adjustment
B05	Stat. Benchmark	Autobox	Forecast provided directly from David Reily with Version 6.0 of the software (June 2007)
B06	Stat. Benchmark	Census - X12 ARIMA	<a href="http://www.census.gov/ts/x12a/v03/x12adocV03.pdf">Official Census method, prepared by McElroy ( www.census.gov/ts/x12a/v03/x12adocV03.pdf)</a>
B07	Stat. Benchmark	Theta	Exponential Smoothing with Decomposition, prepared by Nikolopoulos, Version TIFIS CM3 1.0
B14	Stat. Benchmark	Single ES	Original M3 Benchmark for the M3 competiiton as programmed by Dr M Hibon
B15	Stat. Benchmark	Holt ES	Original M3 Benchmark for the M3 competiiton as programmed by Dr M Hibon
B16	Stat. Benchmark	Dampen ES	Original M3 Benchmark for the M3 competiiton as programmed by Dr M Hibon
B17	Stat. Benchmark	Comb S-H-D ES	Original M3 Benchmark - Equal Weighted combination of Single, Holt and Dampen Expl Smoothing
B09	Stat C	Wildi	An Adaptive Robustified Multi-Step-Ahead Out-Of-Sample Forecasting Combination Approach
B10	Stat C	Beadle	Composite Forecasting Strategy Using Seasonal Schemata
B11	Stat. Contender	Lewicke (Parabolic Systems )	Paracaster Software -fitting equations consisting of a trend plus a series of sinusoidal error terms
B12	Stat. Contender	Hazarika	Decomposition to Random Sequence Basis Functions and a Temperature-Dependent SOFTMAX Combiner
B13	Stat. Contender	Njimi, Mélard	Automatic ARIMA modelling, using TSE-AX
C103	NN/CI Benchmark	Ensemble of Best 3 NN/CI	Equal Weighted combination of C27, C03, C46 prepared post competition by Hibon
C105	NN/CI Benchmark	Ensemble of Best 5 NN/CI	Equal Weighted combination of C27, C03, C46 ,C50, C13 prepared post competition by Hibon

Table 3 shows the results on the complete dataset as average sMAPE, MdRAE, MASE and AR across 111 time series and 18 forecasting horizons. Relative ranks by error measure are given across all methods and for CI contestants alone (NN C).

**Table 3.** NN3 errors and ranks of errors on the complete dataset

		Average Errors				Rank across all methods				Rank across NN C				
		sMAPE	MdRAE	MASE	AR	sMAPE	MdRAE	MASE	AR	sMAPE	MdRAE	MASE	AR	
B09	Wildi	14.84	0.82	1.13	17.3	1	1	1	1	-	-	-	-	Stat C
B07	Theta	14.89	0.88	1.13	17.8	2	3	1	2	-	-	-	-	Stat B
C27	Illies	15.18	0.84	1.25	18.4	3	2	11	4	1	1	4	1	NN C
B03	ForecastPro	15.44	0.89	1.17	18.2	4	4	3	3	-	-	-	-	Stat B
B16	DES	15.90	0.94	1.17	18.9	5	14	3	6	-	-	-	-	Stat B
B17	Comb S-H-D	15.93	0.09	1.21	18.8	6	5	7	5	-	-	-	-	Stat B
B05	Autobox	15.95	0.93	1.18	19.2	7	11	5	7	-	-	-	-	Stat B
C03	Flores	16.31	0.93	1.20	19.3	8	11	6	8	2	5	1	2	NN C
B14	SES	16.42	0.96	1.21	19.6	9	16	7	12	-	-	-	-	Stat B
B15	HES	16.49	0.92	1.31	19.5	10	9	16	9	-	-	-	-	Stat B
C46	Chen	16.55	0.94	1.34	19.5	11	14	18	9	3	7	9	3	NN C
C13	D'yakonov	16.57	0.91	1.26	20.0	12	7	12	15	4	3	5	6	NN C
B00	AutomatANN	16.81	0.91	1.21	19.5	13	7	7	9	5	3	2	3	NN B
C50	Kamel	16.92	0.90	1.28	19.6	14	5	13	12	6	2	6	5	NN C
B13	Njimi	17.05	0.96	1.34	20.2	15	16	18	18	-	-	-	-	Stat C
C24	Abou-Nasr	17.54	1.02	1.43	21.6	16	26	27	25	7	14	16	14	NN C
C31	Theodosiou	17.62	0.96	1.24	20.0	17	16	10	15	8	8	3	6	NN C
B06	Census X12	17.78	0.92	1.29	19.6	18	9	14	12	-	-	-	-	Stat B
B02	nMLP	17.84	0.97	2.03	20.9	19	19	37	19	-	-	-	-	NN B
C38	Adeodato	17.87	1.00	1.35	21.2	20	22	20	20	9	11	10	9	NN C
C26	de Vos	18.24	1.00	1.35	21.7	21	22	20	27	10	11	10	15	NN C
B01	nSVR	18.32	1.06	2.30	21.6	22	29	38	25	-	-	-	-	NN B
C44	Yan	18.58	1.06	1.37	21.2	23	29	23	20	11	15	13	9	NN C
C11	Perfilieva	18.62	0.93	1.57	20.1	24	11	32	17	12	5	19	8	NN C
C37	Duclos	18.68	0.99	1.30	21.5	25	20	15	24	13	9	7	13	NN C
C49	Schliebs	18.72	1.06	1.37	21.9	26	29	23	28	14	15	13	16	NN C
C59	Beliakov	18.73	1.00	1.36	21.4	27	22	22	23	15	11	12	12	NN C
C20	Kurogi	18.97	0.99	1.31	21.3	28	20	16	22	16	9	8	11	NN C
B10	Beadle	19.14	1.04	1.41	22.1	29	28	25	30	-	-	-	-	Stat C
B11	Lewicke	19.17	1.03	1.43	21.9	30	27	27	28	-	-	-	-	Stat C
C36	Sorjamaa	19.51	1.13	1.42	22.5	31	33	26	31	17	18	15	17	NN C
C15	Isa	20.00	1.12	1.53	23.3	32	32	31	33	18	17	18	19	NN C
C28	Eruhimov	20.19	1.13	1.50	23.2	33	33	30	32	19	18	17	18	NN C
C51	Papadaki	22.60	1.27	1.77	25.0	34	35	34	35	20	20	21	20	NN C
B04	Naive	22.69	1.00	1.48	24.2	35	22	29	34	-	-	-	-	Stat B
B12	Hazarika	23.72	1.34	1.80	25.6	36	36	35	37	-	-	-	-	Stat C
C17	Chang	24.09	1.35	1.81	26.3	37	37	36	38	21	21	22	22	NN C
C30	Pucheta	25.13	1.37	1.73	25.3	38	38	33	36	22	22	20	21	NN C
C57	Corzo	32.66	1.51	3.61	26.9	39	39	39	39	23	23	23	23	NN C

Stat.C = statistical contender; NN C = NN/CI contender; Stat.C = statistical contender; NN C = NN/CI contender

Has progress been made, within CI and in comparison to statistical methods? Regardless of accuracy, the ability of 22 contestants to predict a large number of time series with CI indicates

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

unsurprising progress, both on the development of methodologies that facilitate automation and / or increased computational resources.

On accuracy, the top 10 algorithms indicate some progress in accuracy, but not enough to proclaim Chatfield's breakthrough for NN. Unsurprisingly, the top contenders of the M3 monthly data are also ranked high on this subset: Theta (B07), ForecastPro (B03), Autobox (B05) and the ES variants DES (B16), Comb S-H-D (B17), SES (B14) and HES (B15). However, some new innovators have joined the best performers. These algorithms will be briefly introduced as they have not been published elsewhere.

Had it not been a competition tailored to CI, Wildi's (B09) new statistical benchmark method would have won the competition, across all error metrics and against the tough competition of the 'winners' of monthly M3 data. The prototype methodology extends a traditional adaptive state-space approach by discounting errors by their distance to the forecast origin exponentially, estimating multiple step-ahead out-of-sample errors (instead of 1-step ahead in sample), a winsorised squared error loss function, and forecast combinations by building  $h$  separate models for each forecasting horizon  $h = 1, 2, \dots, 18$ , their hyperparameters optimised for each  $h$ , and combining the 18 predictions using the median. A monograph of the algorithm is under preparation.

More in line with the competition theme, the team of Illies, Jäger, Kosuchinas, Rincon, Sakenas and Vaskevcius (C27) ranked 3rd across all methods and provided the best results of all CI-contenders. The methodology employs Echo State Networks (ESN), a novel paradigm of recurrent neural networks with random connections in a reservoir of hidden neurons. Each time series was categorised into 6 clusters by time series length (ignoring different data properties in each cluster and the unrelated nature of most series, that was unknown to the contestants) and decomposed into its time series components using X-12-ARIMA. 500 ESNs with reservoir sizes of 45 to 110 neurons were trained on pooled clusters of time series for each component, and their predictions per time series first recombined across components, and then averaged across all 500 ESNs using the mean. The approach successfully outperformed all statistical benchmarks with the exception of Theta, the top-performer of the M3 monthly data, which constitutes a substantial achievement and progress in CI-model building.

1  
2  
3  
4  
5  
6 Three more CI-contenders outperformed AutomatANN and breached the top-10: Flores et. al.  
7 (C03), who ranked 2nd for CI and 8th overall, employ a self adaptive genetic algorithm (using  
8 conventional crossover and mutation on a fixed population of 100 individuals evolved over 500  
9 generations) to specify an SARIMA model form, parameter bounds and parameters for each time series.  
10  
11 Chen and Yao (C46) employ an ensemble of 500 CART-regression trees built upon bootstrap sampling of  
12 the data and random subspace sampling of features. D'yakonov (C13) used a simple k-nearest neighbour  
13 (k-NN) method with flexible window size conditional on the time series length.  
14  
15  
16  
17  
18  
19

20  
21 The original CI-benchmark, Balkin & Ord's Automat NN (B00) is ranked 5th within all  
22 submitted CI-contenders, outperforming 16 (72%) of the 22 new submissions. Considering that  
23 AutomatANN was automated to run over 3003 series of different frequency, not just 111 monthly series,  
24 and that it was developed a decade ago, it has legitimated its representative performance of NN on  
25 monthly data. However, that 4 (18%) of the submitted CI-approaches outperforming AutomatANN,  
26 shows progress in research through Illies et al. (C27), Flores et al. (C03), Chen et al. (C46) and  
27 D'yakunov (C13). In addition, many CI-contenders achieve accuracy only marginally lower than  
28 AutomatANN, indicating that in contrast to the time of M3 now many researchers are capable of  
29 predicting large scale a level of accuracy similar to AutomatANN.  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39

40 Despite enhanced performance of a few, the field of submissions in CI remains wide, and many  
41 fail to outperform basic CI-benchmarks of a naive multilayer perceptron (B02) or SVR (B01). Some  
42 methods even fail to outperform the Naive benchmark (B04), indicating the need of enhanced  
43 understanding of empirical evaluations for internal benchmarking (ideally prior to a potentially  
44 embarrassing competition performance).  
45  
46  
47  
48  
49

50 It should be noted though, that statistical approaches - simple or complex - are also not a panacea:  
51 other novel statistical contenders such as X-12 (B06) B10 and B11 perform at best average, with B12  
52 even failing to outperform the Naive (B04). Also, weaker contestants of the M3 were not included as  
53 benchmarks, biasing the perception of relative ranking of CI contenders and benchmarks to the  
54  
55  
56  
57  
58  
59  
60



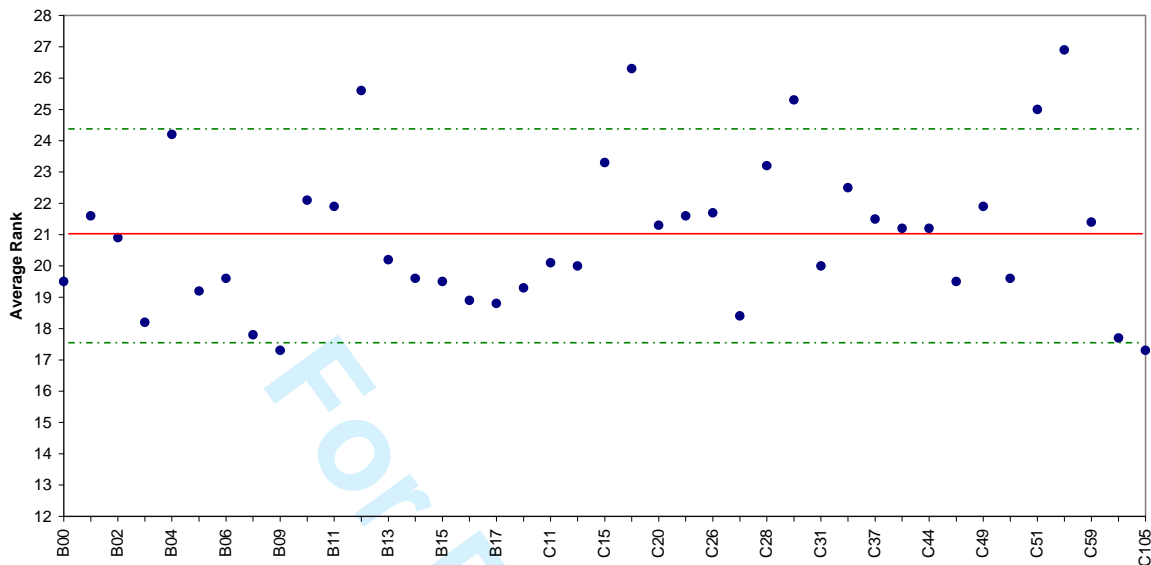
1  
2  
3  
4  
5  
6 disadvantage of NN; in fact, many contestants had outperformed established methods of the M3, but we  
7  
8 were most interested in progress at the top of the field in comparison to AutomatANN.  
9

10 As for the M3-competition, where Hibon computed Comb S-H-D as a novel contender, we sought  
11 to assess the accuracy of combining heterogeneous CI-algorithms. Following the submissions, two  
12 ensembles were created, combining the forecasts of the top three (C27, C03, C46) and the top five (C27,  
13 C03, C46, C13, C50) CI-methodologies using the arithmetic mean. Both CI benchmarks performed  
14  
15 outstandingly well: with an sMAPE of 14.89 the ensemble of the top 3 CI-algorithms would have ranked  
16  
17 overall third - tied with Theta (B07) and better than Echo State Neural Networks (C27). Even more  
18  
19 convincing, with a sMAPE of 14.87 the Ensemble of the top 5 (C105) would have ranked 2nd only to  
20  
21 Wildi (B09), outperforming Theta and all other methods. Although this ex-post combination of best  
22  
23 methods does not represent a valid "ex ante" accuracy (it may be overcome with a quasi-ex ante model  
24  
25 selection), it underlines once more the potential of combining heterogeneous predictions. While Illies' et.  
26  
27 al (C27) performance obviously contributed significantly to the performance of the two CI-ensembles, the  
28  
29 combination increases accuracy beyond that of each individual contender, an effect well documented (in  
30  
31 addition to a second benefit of decreased error variance). More importantly, by including 5 instead of the  
32  
33 top 3 CI-algorithms, essentially introducing more inferior forecasts into an ensemble, the overall accuracy  
34  
35 was increased even further. Therefore, it seems, further increases in accuracy are feasible for CI using  
36  
37 simple means, well documented in the forecasting domain.  
38  
39  
40  
41  
42

#### 43 44 **4.2 Significance of the findings** 45

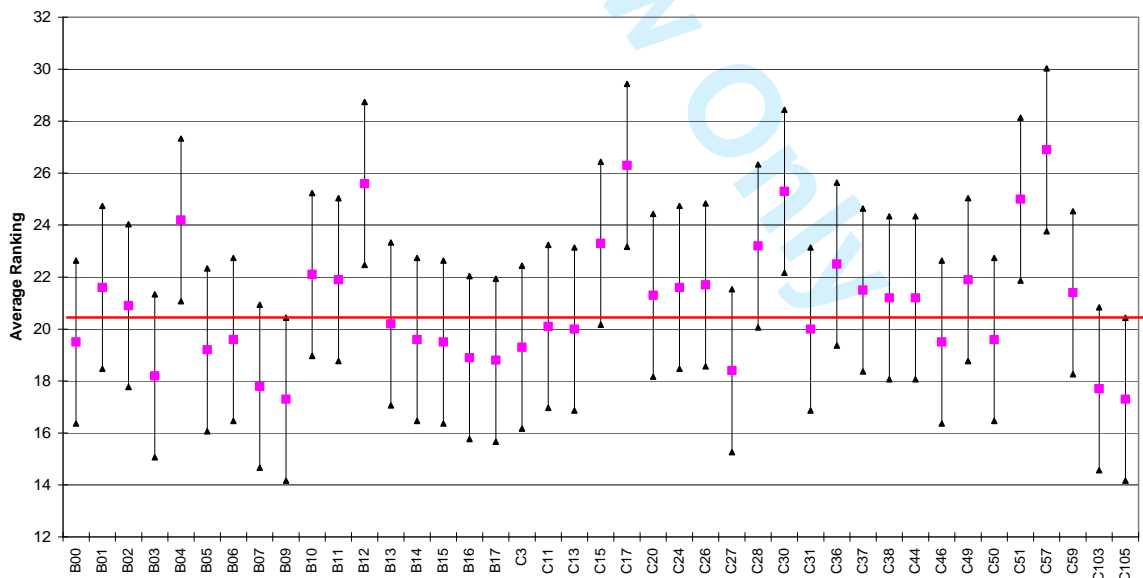
46 Regardless of the recent and vivid discussion about the statistical significance of published  
47 accuracy results within the forecasting community (Armstrong, 2007a, 2007b; Goodwin, 2007) we  
48 computed two non-parametric tests, replicating the analysis by Koning *et. al.* (2005) on the M3: Analysis  
49 of the Mean (ANOM) and Multiple Comparisons to the Best method (MCB), both based upon average  
50 ranks of 41 methods (including both CI-ensembles) over 111 series and 18 horizons (see fig. 1 and 2).  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Figure 1. Analysis of Means on the complete dataset



For ANOM, only the ensemble of the top 5 (C105) and the methodology by Wildi (B09) prove to be statistically significant better than average. On the other side, four CI approaches (C17, C30, C51 and C57) and one statistical contender (B12) perform significantly worse than average.

Figure 2. Multiple Comparisons with the Best on the complete dataset



MCB confirms similar findings as ANOM: the ensemble of the top 5 (C105) and Wildi (B09) are identified as the two best approaches; in comparison to them, the same four CI approaches (C17, C30, C51, C57) and a statistical contender (B12) plus the Naïve (B04) are significantly worse than the best.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Despite limited differences in statistical significance, it is worth mentioning that even a small accuracy gain, e.g. of 1%, is often amplified in operational benefits, and could result in manifold savings in safety stocks. Thus accuracy results in term of average metrics should never be ignored, as they are often operationally significant (Syntetos, Nikolopoulos, & Boylan, 2010). As an indication of the tests' limitations, the Theta method - previously better than other algorithms in the competition - is no longer significantly better, indicating the sensitivity of the test to sample size and structure (as of all tests), adding further to the discussion of tests.

#### 4.3 Analysis of Data Conditions

Next, we analyse the data conditions under which different algorithms perform well. As it is infeasible to present all 24 tables of ranking per error measure and data subset, Table 4 summarizes the results of the top five performers for both the complete and reduced datasets (111 and 11 series), and for the conditions of long and short time series lengths (50 series each), seasonal and non-seasonal time series patterns (50 series each), and the combination of both conditions (25 series each). Table 5 shows the top 5 performers by sMAPE across different forecasting horizons. In order to facilitate replication and external analysis of the results, all tables of sMAPE (Tables 6-15), MdRAE (Tables 16-20), MASE (Tables 21-25), AR for all methods and for CI contenders separately (Tables 26-27), and ANOM (Table 28) and MCB based upon AR (Table 29) will be provided online on the competition and journal homepage.

**Table 4.** NN3-competition results across data conditions on sMAPE, MdRAE, and MASE

Error Metrics	Complete dataset (incl.reduced)	Data Conditions					Combined Data Conditions			
		Reduced dataset	Short	Long	Seasonal	Non-seas.	Short		Long	
							Non-seas	Seasonal	Non-seas	Seasonal
# of series	111	11	50	50	50	50	25	25	25	25
sMAPE	<i>B09</i>	B05	<b>C27</b>	B03	<i>B09</i>	<b>C105</b>	<b>C27</b>	<i>B09</i>	B03	B03
	<b>C105</b>	B03	<b>C105</b>	<i>B09</i> &B16	<b>C105</b>	B07	<b>C105</b>	<b>C27</b>	B16	<i>B09</i>
	<i>B07</i> & <b>C103</b>	<b>C44</b>	<b>C103</b>	-	<b>C103</b>	<b>C103</b>	<b>C103</b>	<b>C105</b>	<i>B07</i> & <i>B09</i>	B17
	-	B07	<i>B09</i>	B14	B07	B03	B17	<b>C103</b>	-	B14
	<b>C27</b>	<b>C59</b>	B07	B07	<b>C27</b>	B14	B13	<b>C50</b>	<u>B00</u>	B16
MdRAE	<b>C105</b>	C38	<b>C27</b>	B03	<i>B09</i>	<b>C105</b>	<b>C27</b>	<b>C27</b>	<u>B00</u>	B16
	<i>B09</i> & <b>C103</b>	<b>C105</b>	<i>B09</i>	<i>B09</i> &B15	<b>C27</b>	<u>B00</u>	<b>C105</b>	<i>B09</i>	<i>B03</i> & <i>B09</i>	B14&B17
	-	<b>C11</b>	<b>C105</b>	B16&B17	<b>C103</b> & <b>C105</b>	<b>C27</b>	<i>B09</i> &B14	<b>C103</b> & <b>C105</b>	-	-
	<b>C27</b>	<b>C103</b>	<b>C50</b> & <b>C103</b>	-	-	<i>B09</i> & <b>C50</b>	-B17	-	B16& <b>C105</b>	B03
	B07	B03	-	-	B07	-	-	B05	-	B07
MASE	<b>C105</b>	<i>B05</i> & <b>C59</b>	<b>C105</b>	B14	B09	B14	<i>B09</i>	<b>C27</b>	B14	B14-B17
	<i>B07</i> & <i>B09</i>	-	<b>C27</b> & <b>C103</b>	B16	B07	<b>C105</b>	<b>C103</b> & <b>C105</b>	<b>C103</b> & <b>C105</b>	<u>B00</u>	-
	-	B03	-	B07	<b>C105</b>	<u>B00</u>	-	-	B16	<i>B09</i> &B16
	<i>B03</i> &B16	<b>C44</b>	<i>B09</i>	B17& <b>C105</b>	<b>C103</b>	<i>B07</i> &B16	<b>C27</b>	<i>B07</i> -B17	B04	-
	-	<b>C18</b>	<b>C50</b>	-	<b>C27</b>	-	<b>C50</b>	-	<b>C105</b>	B03

Bold: CI-contenders, *Italics*: Statistical contenders, Normal: Benchmarks, Underlined: AutomatANN M3 benchmark

On the complete dataset (first column of table 4), the ranking of all algorithms in is identical to results provided in table 3, identifying the top performers of NN3 by sMAPE, id est: Wildi (B09), ensemble of top 5 CI (C105), Theta (B07) in a draw with the ensemble of top 3 CI (C103 ) and Illies et al. (C27). In comparison, different algorithms performed well on the reduced dataset of 11 harder to forecast time series: the statistical expert system Autobox (B05) ranks 1st by sMAPE, playing out its strengths in modeling pulse interventions, level shifts, local time trends and seasonal pulses. ForecastPro (B03) ranks 2nd and Theta (B07) ranked 4th. Two new CI-contestants enter the top 5 of the reduced dataset: Yan (C44), ranked 3rd on sMAPE across all methods and 1st for CI-methods, employs three sets of 18 Generalized Regression NNs per time series, each trained separately to predict for a forecasting horizon  $h = 1, 2, \dots, 18$  with three distinct parameters settings, recombining the predictions to one trace forecast, and combining the predictions of the three architectures, hence called 'multiple model fusion'.

On MdRAE other CI-contenders enter the top 5, Adeonato et al. (C38) using ensembles of 15 MLPs and Perfilieva (C11), forecasting using fuzzy transforms, indicating that the results on only a few

1  
2  
3  
4  
5  
6 series do not yield the same level of reliability across error measures as on the complete set. It does  
7  
8 however show the potential that specialised statistical and CI-algorithms tuned (or robust) to particular  
9  
10 time series properties can outperform other approaches, but at the same time questions the ability of these  
11  
12 CI methodologies to generalise on larger datasets than the ones they were originally tailored to.  
13

14  
15 Next, we analyse the results across the data conditions of time series length and seasonality.  
16  
17 Wildi's (B09) new statistical approach ranks well under all data conditions and metrics, with the  
18  
19 exception of short & non-seasonal series on sMAPE, indicating that some of its success derives from  
20  
21 capturing seasonality well (1<sup>st</sup> for all metrics). Variants of ES (B14, B15, B16 and their combination B17)  
22  
23 make frequent appearances on long & seasonal time series, indicating that the decomposition approach  
24  
25 used for M3 – *DeSeasonalise + Extrapolate + ReSeasonalise* - works competitively. Similarly, the expert  
26  
27 system ForecastPro (B03), which selects amongst these methods, outperforms them on long series of both  
28  
29 seasonal and non-seasonal data, confirming that industry still does well to rely this family of methods for  
30  
31 these typical data conditions. The related Theta (B07) appears on all aggregate conditions but not its  
32  
33 combinations, verifying its robustness across many data conditions by a consistent level of accuracy, but  
34  
35 not winning any particular category.  
36  
37

38  
39 For CI, multiple CI-contenders enter the top 5 on different conditions, while the M3 benchmark  
40  
41 AutomatANN (B00) is absent across all categories and metrics (with the exception of sMAPE on  
42  
43 *Long+Non-seasonal* data). In the light of earlier research, the most striking result of NN3 comes in the  
44  
45 *Short+Non-seasonal* subset, judging by recent publications one of the most difficult conditions for CI-  
46  
47 methods. Echo State Networks by Illies et. al (C27) achieved the *colpo grosso* and won this category as  
48  
49 well as that of the broader 50 short series, we speculate as an effect of training on pooled clusters of time  
50  
51 series. CI ensembles C103 and C105 performed equally well across data conditions of short & seasonal  
52  
53 and short & non-seasonal series, ranking 2nd / 3rd and 3rd / 4th respectively, but less so across long series  
54  
55 with and without seasonality (unsurprising as C27 was contained in them). From the remaining CI  
56  
57 competitors only Kamel (C50) made an appearance in the *Short+Seasonal* category, combining MLPs  
58  
59 with Gaussian Process Regression.  
60

These results challenge prior beliefs in NN-modeling that a significant amount of historic observations are a prerequisite for sufficient initialization, training, validation, evaluation and generalisation of CI approaches (see, e.g., Haykin, 1999). Furthermore, across time series patterns more CI are ranked highly on seasonal data than on non-seasonal data, a second fundamental contradiction to prior research which had identified problems in predicting seasonal time series with NNs and proposed prior deseasonalisation (e.g., Zhang & Qi, 2005). While these results reveal no insight into the reasons for the increased performance, it does demonstrate that novel CI-paradigms can yield competitive performance beyond their traditional application domain and that systematic replications of earlier studies should be conducted to challenge prior findings. However, the majority of CI approaches is absent across datasets and conditions, one the one hand demonstrating consistent results, but on the other indicating that only few algorithms have the capacity to perform well.

**Table 5.** NN3 results of sMAPE across short, medium, long and all forecasting horizons

Error Metrics	Complete dataset (incl.reduced)	Combined Data Conditions				
		Reduced dataset	Short		Long	
			Non-seas	Seasonal	Non-seas	Seasonal
# of series	111	11	25	25	25	25
Short ( $h = 1-3$ )	B07 <i>B09</i> B03-C105 - C103	C20 <i>B10</i> C08 B03 C59	C105 C27 C50 <i>B09</i> C59	C27 <i>B09</i> B00 B05 C50	B07 B03 B16 B06 B17	B16-B17 - B03 B14 B15
Medium ( $h = 4-12$ )	C105 <i>B09</i> C103 B07 C27	C44 C50 C46 B07 B05	C27 B17 C105 C103 B14	B09 C50 C105 C27 C103	B03 C3 B07 B16 B14	B09 B06 B03 B16 B17
Long ( $h = 13-18$ )	C103 B07 C105 <i>B09</i> C27	B05 C38 B03 C18 C59	C27 C46 C103 B13 B14	B09 C27 C103 B07 C105	C105 C103 B09 C13 B00	B17 B14 B03 B07 B16
All ( $h = 1-18$ )	B09 C105 B07-C103 - C27	B05 B03 C44 B07 C59	C27 C105 C103 B17 B13	B09 C27 C105 C103 C50	B03 B16 B07-B09 - B00	B03 B09 B17 B14 B16

Bold: CI-contenders, *Italics*: Statistical contenders, Normal: Benchmarks, Underlined: AutomatANN M3 benchmark

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Results across forecasting horizons seem to confirm earlier findings by Hill et al (Hill, et al., 1996): for short term forecasting, methods of ES (B07 and B09) appear to perform best, but with an increasing forecasting horizon the CI approaches take the lead, although it remains unclear whether the contribution stems from the forecast combinations in ensembles, or the underlying method's performance increasing with the horizon.

However, for CI the accuracy achieved across horizons show surprising consistency. On the complete dataset those contenders ranked highly overall are also ranked consistently amongst the top 5 across all horizons of short, medium and long forecasts, with only minor changes in ranks. This is confirmed across data conditions, where the relative performance remains consistent across different horizons: CI-methods perform well for short time series with and without seasonality across all forecasting horizons, in particular Illies' (C27) and the ensembles C105 and C103. Similarly for long time series, methods of ES perform consistently well across all horizons, again without significant changes in ranks. The only noticeable change appears for long & non-seasonal data, where ES dominates on short, and CI on long horizons. It stands to argue, that results across horizons for a particular data subset, remain more stable than expected given prior findings. As an example, Wildi's (B09) approach, which is optimised specifically for multiple horizons of a trace forecast, performs consistently across all horizons for short & seasonal time series, as was intended by the algorithm.

## 5. Discussion

The NN3 competition contributed empirical evidence in the tradition of the M-competitions, with a particular emphasis on extending the findings of M3 towards a current and complete range of CI-methods. In that the NN3 has succeeded, attracting contestants from all mayor paradigms, including feed-forward and recurrent NN, Fuzzy Logic, Genetic Algorithms and Evolutionary Computation and hybrid systems. In addition, the results of this replication and extension of M3 allow us to evaluate the six hypothesis of the original M-competition (see section 2), and to determine if the findings conform to established wisdom, or add novel insights to the body of knowledge. First, we will review hypothesis H2,

1  
2  
3  
4  
5  
6 H1 and H3 as they allow us to assess the similarity of M3 and its replication, and allow a verification of  
7  
8 the NN3 competition design.  
9

10  
11  
12 (H1) *'Data characteristics determine relative performance?'* The results of N3 across data conditions  
13  
14 (table 4) confirm that of earlier M3. Data characteristics substantially influence the relative  
15  
16 performance of algorithms in statistics and CI alike. Here, NN3 contributes to the discussion by  
17  
18 providing objective evidence that NNs are capable to predict seasonal time series (in contrast to,  
19  
20 e.g., Zhang & Qi, 2005), and to predict short time series (in contrast to, e.g., Hill, et al., 1996)  
21  
22 accurately, contrary to prior publications and indicating the need for further research.  
23

24  
25 (H2) *'Accuracy depends upon forecasting horizon?'* Across forecasting horizons (table 5), relative  
26  
27 performance varies, different methods perform best in different horizons, confirming the findings  
28  
29 of M3. Also, the efficacy of CI-methods in comparison to statistical methods increases for longer  
30  
31 forecasting horizons, as identified in prior studies (Hill, et al., 1996) However, for the best  
32  
33 algorithms in CI the accuracy remained almost consistent for increasing forecasting horizons.  
34  
35 Further research is needed to determine if those methods incorporating trace errors in their  
36  
37 modelling (e.g. B09, C44) can overcome this limitation, as first indications seem to suggest.  
38

39  
40 (H3) *'Performance ranking varies with metric?'* NN3 rankings based upon sMAPE, MdRAE, MASE  
41  
42 and AR result in a different relative performance of algorithms, across all datasets and data  
43  
44 conditions (see table 4). However, many methods in the upper deciles of the field perform  
45  
46 consistently well on multiple metrics, and vice versa, building increasing confidence on their  
47  
48 relative performance and predictive capabilities.  
49

50  
51  
52 Next, we will review H5 and H6 which consider the relative accuracy of algorithms, the main topic  
53  
54 of this extension of the M3 competition.  
55  
56  
57  
58  
59  
60



1  
2  
3  
4  
5  
6 (H5) *Combinations outperform individual methods?* Reviewing common properties of the top  
7  
8 performers (table 4), the success of combinations stands out. With the exception of the five original  
9  
10 submissions to the M3 (ForecastPro, Autobox, SES, DES, and HES) the three leading statistical  
11  
12 methods in the top 10 use forecast combinations (most notably Wildi (B09) across all conditions,  
13  
14 Comb S-H-D (B017 for long series, and also Theta (B07) which essentially employs a weighted  
15  
16 forecast combination of linear trend and ES). Also, with the exception of Flores (C03) all CI-  
17  
18 methodologies in the top 10 employ forecast combinations (Illies (C27), Chen (C46), Ensemble of  
19  
20 the top 5 (C105), and Ensemble of the top 3 NN (C103)). C105 and C103 dominate our results, but  
21  
22 also signify the effect of increasing coverage in the ensemble (i.e. heterogeneity of the base  
23  
24 learner), which warrants more research effort across disciplines. As sophisticated 'ensembles' in the  
25  
26 form of boosting, bagging, arcing etc. are more widespread in CI classification than in statistical  
27  
28 modelling and time series prediction in particular, we see potential for cross-disciplinary research.  
29  
30

31  
32  
33 (H6) *Sophisticated methods are not better than simpler methods?* Seeing the majority of CI approaches  
34  
35 have failed to outperform Simple ES (B14), and four performed worse than Naïve (B04) (see  
36  
37 Tables 6-15) we could not disagree. However, NN3 has introduced a novel univariate method, and  
38  
39 provided the evidence of its capability to outperform established statistical benchmarks, including  
40  
41 the respective winners of the monthly M3 data (Dampen ES, Theta and ForecastPro) and all CI-  
42  
43 contenders to date. Although the algorithm by Wildi (B07) is statistical in nature and not based  
44  
45 upon CI, the method cannot be classified as anything other than complex, combining various  
46  
47 innovations in estimation and model selection to automatically tune it to the data. This stands in  
48  
49 conflict to H6 and to common belief that complex methods do not significantly outperform simple  
50  
51 ones. Similarly, NN3 provide evidence that some complex methods are capable to outperform all  
52  
53 statistical methods of the M3, showing a substantial improvement in accuracy. To provide further  
54  
55 evidence, with the submissions of Wildi, Theta, ForecastPro, Autobox for statistics, and Illies and  
56  
57 Flores representing CI, 4 of the top 5 (80%) and 6 of the top 10 methods (60%) could be classified  
58  
59  
60

1  
2  
3  
4  
5  
6 as complex methods. As such, we have provided objective evidence that does not support H6. Short  
7  
8 of refuting H6 on the basis of a few algorithms, we seek to reverse it to challenge established  
9  
10 wisdom: (H6.b) *Simple methods are not better than sophisticated methods*. Despite identical  
11  
12 content, the prior connotation of H6 all too easily suggested that no benefit arose from  
13  
14 sophistication, and allowed misinterpretation that *Simpler is better*'. We conclude that complex  
15  
16 methods of CI/NN and statistics have caught up, and overall simple statistical methods can no  
17  
18 longer claim to outperform CI methods without proper empirical evaluation.  
19

20  
21  
22  
23 As every empirical study, findings hold only for the properties of the empirical dataset provided,  
24  
25 and as such the NN3 competition did not aim to be representative of all data properties in operational  
26  
27 forecasting. Still our competition was prone to certain limitations and biases that must be critically  
28  
29 reviewed. These include the obvious shortcomings that are endogenous to most competitions: no rolling  
30  
31 origin design (due to the challenge of organising such as setup), limited representativeness of datasets in  
32  
33 size, structure and heterogeneity, and the exclusion of certain error metrics that assess the final impact on  
34  
35 decision making, i.e. inventory costs arising from operational forecasting (Timmermann and Granger,  
36  
37 2004). As in prior M competitions, our assessment considered only empirical accuracy and neglected  
38  
39 computational resources required, an important aspect in forecasting for operations. As Expert software  
40  
41 systems such as Autobox and ForecastPro contain much faster forecasting engines than CI (i.e. we  
42  
43 received the submission of Autobox almost instantaneous after the release of the data), algorithms and  
44  
45 systems employing efficient statistical methods may still remain the first choice in operations.  
46  
47

48  
49 Despite our efforts, biases in the representativeness of algorithms may exist. In tailoring the NN3  
50  
51 to algorithms of NN and CI we may have biased the sample of contestants by attracting more CI  
52  
53 contestants than those from statistics. Furthermore, the majority of submissions came from researchers in  
54  
55 CI, while professionals and (possibly advanced) software companies in NN, CI or AI (e.g. Siemens,  
56  
57 Alyuda, NeuroDimensions, and SAS) chose not to participate despite personal invitations. Also, more  
58  
59 participation from econometrics and forecasting software vendors active in forecasting for operations (e.g.  
60

1  
2  
3  
4  
5  
6 SAP, Oracle, John Galt, Smart etc.) would have increased the validity of results. However, we tried to be  
7  
8 as objective and inclusive as we could, taking into consideration the design suggestions of prior  
9  
10 competitions and reaching out to the communities omitted before. Therefore we are confident that NN3  
11  
12 provides a more comprehensive and up-to-date assessment of the performance of CI-methods in  
13  
14 predicting monthly time series than M3, as well as more valid and reliable evidence than that of prior  
15  
16 competitions in CI. One fundamental flaw - grounded in the nature of a replication - exists in the prior  
17  
18 availability of the data, although its origin was undisclosed and masked in a sample. Although we are  
19  
20 convinced of the integrity of all contestants, it reminds us of the importance of true ex-ante evaluations on  
21  
22 unknown data to avoid any data snooping for future competitions.  
23  
24

## 25 **6. Conclusions**

26  
27  
28 Replicating and extending the prominent M3 competition, NN3 aspired to challenge prior  
29  
30 evidence on the inferior forecasting accuracy of NN approaches in operational forecasting. The final  
31  
32 results assess the accuracy of over 60 forecasting algorithms, the largest assessment of methods on time  
33  
34 series data to date. Ex ante accuracies were evaluated on 111 or 11 empirical time series using multiple  
35  
36 established error metrics and following a rigorous competition design; conditions examined include the  
37  
38 presence of seasonality, the length of the series, and the forecasting horizon.  
39

40  
41 The NN3 objective, extending the M3 competition towards NN and CI algorithms, was  
42  
43 successfully achieved in attracting 46 CI contestants and novel statistical benchmarks, making it the  
44  
45 largest empirical evaluation in NN, CI and forecasting to date. The main findings confirm prior  
46  
47 hypothesis, but also initiate new research discussions. New algorithms are feasible, in CI, NN and  
48  
49 statistics alike. The competition assessed a novel statistical - and complex - method by Wildi (B9), which  
50  
51 showed exceptional performance on both datasets. Illies' et al. (C27) introduced a NN-methodology that  
52  
53 outperformed Dampen Trend ES, but still lacked the performance of the Theta across all series. This  
54  
55 algorithm also outperformed all other algorithms on 25 short and seasonal time series, the most difficult  
56  
57 subset of the competition, and Yan (C44) outperformed all others on a subset of complex/difficult series.  
58  
59 These achievements are surprising considering prior beliefs on the data properties required to use NN  
60

1  
2  
3  
4  
5  
6 methods on empirical data, and demand further attention. Overall, we hope the success of complex  
7 algorithms on such an established dataset will at least rekindle the discussion of innovative, sophisticated  
8 algorithms for time series extrapolation.  
9

10  
11  
12 Results of NN3 suggest that methods of NN and CI can perform competitively to established  
13 statistical methods in time series prediction, but still cannot outperform them. However, in the absence of  
14 any (statistically significant) differences between algorithms we can no longer assume that they are  
15 inferior either. Considering the results of the M3, we have consciously included the top-performers of  
16 ForecastPro, Theta, and Comb S-H-D as hard benchmarks for NN to compete against. As such, we  
17 expected that the methods of ES, the workhorses of operational forecasting in practice for over 30 years,  
18 would be serious contenders that would prove challenging to outperform - after all they had outperformed  
19 most others in the original M3. It should however be noted that the other 20 statistical methods in M3  
20 performed less admirably and would not be expected to do better in comparison to many CI contestants.  
21 We feel that CI has closed in on established benchmarks, showing a range of different algorithms capable  
22 to predict both datasets as accurate as AutomatANN, the only contestant entering the M3 some 10 years  
23 ago, thus indicating improvements in feasibility and empirical accuracy to forecast with NN, and hence a  
24 motivation for further research.  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39

40 Disappointingly, it seems impossible to provide more focussed guidance as to promising routes of  
41 future CI research, as no common 'best practises' can be identified for the top NN or CI contenders. Each  
42 submission was unique conceptually and methodologically, combining freely (and often seemingly  
43 arbitrarily) from the repository of algorithms and techniques available to machine learning today, without  
44 any evaluation of the contribution each fragment of the methodology made to increasing accuracy. As an  
45 example, for Illies' et al. it remains unclear if the accuracy stems from pooling time series for training,  
46 combining predictions in ensembles, or the algorithm of Echo State Networks itself. In an attempt to  
47 generalise, only the paradigm of forecast combinations seemed to drive accuracy, an observation made  
48 before. Ensembles of CI and statistical algorithms performed very well, but again no consensus on its  
49 meta-parameters of ensemble size or combination metric could be determined, although heterogeneity of  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 its base learners seemed to positively affect accuracy. As no two algorithms are alike, it becomes  
7  
8 impossible to attribute positive performance to a particular modelling choice, allowing only an evaluation  
9  
10 of composite yet distinct algorithms but not to derive guidance for promising areas of future research.  
11  
12 Without insight, progress in CI may be slow and undirected. If this heterogeneity cannot be overcome,  
13  
14 only a meta-learning analysis could yield insights in partial contributions, linking properties of algorithms  
15  
16 and data conditions to guide future research effort.  
17

18  
19 The NN3 competition has proven a stimulating exercise that has attracted, engaged and unified  
20  
21 researchers from forecasting, informatics, machine learning, data mining and engineering. We therefore  
22  
23 hope that the NN3 will not only provide a means to disseminate best-practices on CI-methods, but more  
24  
25 importantly on competition design outside the forecasting community. We conclude that the findings of  
26  
27 the NN3 competition provide encouraging evidence for the potential of NN and CI-methods in time series  
28  
29 prediction, even for a well established domain as monthly time series prediction. The promising results of  
30  
31 NN3 motivate us to run future competitions to add knowledge in modelling neural networks for time  
32  
33 series prediction. Already, it has sparked a resurgence of interest in competitions in CI, with regular  
34  
35 competitions tracks held at ESTSP, IJCNN and WCCI conferences since. For future competitions, we see  
36  
37 the need to evaluate novel application domains that are empirically important yet previously omitted, in  
38  
39 particular those of high frequency data where NN are regularly employed in practice. Still, no method will  
40  
41 be a panacea. Yet only in extending competition designs to novel data conditions, beyond that of the M-  
42  
43 style competitions, will we be able to determine on what data the application of neural networks indeed  
44  
45 holds a breakthrough or a passing fad.  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60
- Adya, M., Armstrong, J. S., Collopy, F., & Kennedy, M. (2000). An application of rule-based forecasting to a situation lacking domain knowledge. *International Journal Of Forecasting*, 16, 477-484.
- Adya, M., & Collopy, F. (1998). How effective are neural networks at forecasting and prediction? A review and evaluation. *Journal Of Forecasting*, 17, 481-495.
- Armstrong, J. S. (2006). Findings from evidence-based forecasting: Methods for reducing forecast error. *International Journal Of Forecasting*, 22, 583-598.
- Armstrong, J. S. (2007a). Significance tests harm progress in forecasting. *International Journal Of Forecasting*, 23, 321-327.
- Armstrong, J. S. (2007b). Statistical significance tests are unnecessary even when properly done and properly interpreted: Reply to commentaries. *International Journal Of Forecasting*, 23, 335-336.
- 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., Song, H., & Wu, D. C. (2009). The tourism forecasting competition. In *Monash University Working Paper: Monash University*.
- Balkin, S. D., & Ord, J. K. (2000). Automatic neural network modeling for univariate time series. *International Journal Of Forecasting*, 16, 509-515.
- Breiman, L. (1984). *Classification and regression trees*. Belmont, Calif.: Wadsworth International Group.
- Cawley, G. C., Janacek, G. J., Haylock, M. R., & Dorling, S. R. (2007). Predictive uncertainty in environmental modelling. *Neural Networks*, 20, 537-549.
- Chatfield, C. (1993). Neural Networks - Forecasting Breakthrough or Passing Fad. *International Journal Of Forecasting*, 9, 1-3.
- Chen, B. J., Chang, M. W., & Lin, C. J. (2004). Load forecasting using support vector machines: A study on EUNITE competition 2001. *Ieee Transactions on Power Systems*, 19, 1821-1830.
- Crone, S. F., & Pietsch, S. (2007). A naive support vector regression benchmark for the NN3 forecasting competition. *2007 Ieee International Joint Conference on Neural Networks, Vols 1-6*, 2453-2458.
- Crone, S. F., & Preßmar, D. B. (2006). An extended evaluation framework for publications on artificial neural networks in sales forecasting. In IASTED (Ed.), *AIA'06 (Vol. 1)*. Innsbruck: IASTED Press.
- Curry, B. (2007). Neural networks and seasonality: Some technical considerations. *European Journal Of Operational Research*, 179, 267-274.
- Dawson, C. W., See, L. M., Abrahart, R. J., Wilby, R. L., Shamseldin, A. Y., Anctil, F., Belbachir, A. N., Bowden, G., Dandy, G., Lauzon, N., Maier, H., & Mason, G. (2005). A comparative study of artificial neural network techniques for river stage forecasting. *Proceedings of the International Joint Conference on Neural Networks (IJCNN), Vols 1-5*, 2666-2670.
- de Menezes, L. M., & Nikolaev, N. Y. (2006). Forecasting with genetically programmed polynomial neural networks. *International Journal Of Forecasting*, 22, 249-265.
- Faraway, J., & Chatfield, C. (1998). Time series forecasting with neural networks: A comparative study using the airline data. *Applied statistics*, 47, 20.
- Fildes, R. (1992). The Evaluation Of Extrapolative Forecasting Methods. *International Journal Of Forecasting*, 8, 81-98.
- Fildes, R. (2006). The forecasting journals and their contribution to forecasting research: Citation analysis and expert opinion. *International Journal Of Forecasting*, 22, 415-432.

- 1  
2  
3  
4  
5  
6 Fildes, R., Hibon, M., Makridakis, S., & Meade, N. (1998). Generalising about univariate  
7 forecasting methods: further empirical evidence. *International Journal Of Forecasting*,  
8 *14*, 339-358.  
9  
10 Fildes, R., & Makridakis, S. (1995). The Impact Of Empirical Accuracy Studies On Time-Series  
11 Analysis And Forecasting. *International Statistical Review*, *63*, 289-308.  
12 Fildes, R., Nikolopoulos, K., Crone, S. F., & Syntetos, A. A. (2008). Forecasting and operational  
13 research: a review. *Journal of the Operational Research Society*, *59*, 1150-1172.  
14 Fildes, R., & Ord, K. (2002). Forecasting competitions: their role in improving forecasting  
15 practice and research. In *A companion to economic forecasting*. Malden, Mass. [u.a.]:  
16 Blackwell.  
17 Fogel, D. B., & Fogel, L. J. (1994). EVOLUTIONARY COMPUTATION. *Ieee Transactions on*  
18 *Neural Networks*, *5*, 1-1.  
19 Foster, W. R., Collopy, F., & Ungar, L. H. (1992). Neural Network Forecasting of Short, Noisy  
20 Time-Series. *Computers & Chemical Engineering*, *16*, 293-297.  
21 Goodrich, R. L. (2000). The Forecast Pro methodology. *International Journal Of Forecasting*,  
22 *16*, 533-535.  
23  
24 Goodwin, P. (2007). Should we be using significance tests in forecasting research? *International*  
25 *Journal Of Forecasting*, *23*, 333-334.  
26 Goodwin, P., & Lawton, R. (1999). On the asymmetry of the symmetric MAPE. *International*  
27 *Journal Of Forecasting*, *15*, 405-408.  
28 Groff, G. K. (1973). Empirical Comparison of Models for Short Range Forecasting.  
29 *Management Science Series a-Theory*, *20*, 22-31.  
30 Guyon, I., Saffari, A., Dror, G., & Cawley, G. (2008). Analysis of the IJCNN 2007 agnostic  
31 learning vs. prior knowledge challenge. *Neural Networks*, *21*, 544-550.  
32 Haykin, S. S. (1999). *Neural networks : a comprehensive foundation* (2nd ed.). Upper Saddle  
33 River, NJ: Prentice Hall.  
34 Hill, T., O'Connor, M., & Remus, W. (1996). Neural network models for time series forecasts.  
35 *Management Science*, *42*, 1082-1092.  
36 Hyndman, R. J., & Koehler, A. B. (2006). Another look at measures of forecast accuracy.  
37 *International Journal Of Forecasting*, *22*, 679-688.  
38 Kolarik, T., & Rudorfer, G. (1994). Time Series Forecasting Using Neural Networks. In  
39 *Proceedings of the international conference on APL* (pp. 86-94). Antwerp, Belgium.  
40 Koning, A. J., Franses, P. H., Hibon, M., & Stekler, H. O. (2005). The M3 competition:  
41 Statistical tests of the results. *International Journal Of Forecasting*, *21*, 397-409.  
42 Lendasse, A., Oja, E., Simula, O., & Verleysen, M. (2007). Time series prediction competition:  
43 The CATS benchmark. *Neurocomputing*, *70*, 2325-2329.  
44 Liao, K. P., & Fildes, R. (2005). The accuracy of a procedural approach to specifying  
45 feedforward neural networks for forecasting. *Computers & Operations Research*, *32*,  
46 2151-2169.  
47 Makridakis, S. (1994). Book Review: Time Series Prediction - Forecasting the Future and  
48 Understanding the Past - Weigend, A.S., Gershenfeld, N.A. . *International Journal Of*  
49 *Forecasting*, *10*, 463-466.  
50 Makridakis, S., Andersen, A., Carbone, R., Fildes, R., Hibon, M., Lewandowski, R., Newton, J.,  
51 Parzen, E., & Winkler, R. (1982). The Accuracy Of Extrapolation (Time-Series) Methods  
52 - Results Of A Forecasting Competition. *Journal Of Forecasting*, *1*, 111-153.  
53 Makridakis, S., Chatfield, C., Hibon, M., Lawrence, M., Mills, T., Ord, K., & Simmons, L. F.  
54 (1993). The M2-competition: A real-time judgmentally based forecasting study.  
55 *International Journal Of Forecasting*, *9*, 5-22.  
56  
57  
58  
59  
60

- 1  
2  
3  
4  
5  
6 Makridakis, S., & Hibon, M. (1979). Accuracy Of Forecasting - Empirical-Investigation. *Journal*  
7 *Of The Royal Statistical Society Series A-Statistics In Society*, 142, 97-145.
- 8 Makridakis, S., & Hibon, M. (2000). The M3-Competition: results, conclusions and implications.  
9 *International Journal Of Forecasting*, 16, 451-476.
- 10 McNames, J., Suykens, J. A. K., & Vandewalle, J. (1999). Winning entry of the K. U. Leuven  
11 time-series prediction competition. *International Journal of Bifurcation and Chaos*, 9,  
12 1485-1500.
- 13 Nelson, M., Hill, T., Remus, W., & O'Connor, M. (1999). Time series forecasting using neural  
14 networks: Should the data be deseasonalized first? *Journal Of Forecasting*, 18, 359-367.
- 15 Newbold, P., & Granger, C. W. J. (1974). Experience With Forecasting Univariate Time Series  
16 And Combination Of Forecasts. *Journal Of The Royal Statistical Society Series A-*  
17 *Statistics In Society*, 137, 131-165.
- 18 Ord, K., Hibon, M., & Makridakis, S. (2000). The M3-Competition. *International Journal of*  
19 *Forecasting*, 16, 433-436.
- 20 Preminger, A., & Franck, R. (2007). Forecasting exchange rates: A robust regression approach.  
21 *International Journal Of Forecasting*, 23, 71-84.
- 22 Reid, D. J. (1969). *A comparative study of time series prediction techniques on economic data.*  
23 Unpublished PhD thesis, University of Nottingham, Nottingham, UK.
- 24 Reid, D. J. (1972). A comparison of forecasting techniques on economic time series. In M. J.  
25 Bramson, I. G. Helps & J. A. C. C. Watson-Grady (Eds.), *Forecasting in Action.*  
26 Birmingham, UK: Operations research Society.
- 27 Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1994). Learning representations by back-  
28 propagating errors (from Nature 1986). *Spie Milestone Series Ms*, 96, 138.
- 29 Sharda, R., & Patil, R. B. (1992). Connectionist Approach To Time-Series Prediction - An  
30 Empirical Test. *Journal Of Intelligent Manufacturing*, 3, 317-323.
- 31 Sincák, P., Strackeljanc, J., Kolcun, M., Novotný, D., & Szathmáry, P. (2002). Electricity Load  
32 Forecast Using Intelligent Technologies. . In: EUNITE European Network of Intelligent  
33 Technologies.
- 34 Suykens, J. A. K., & Vandewalle, J. (1998a). *The K.U.Leuven time series prediction competition.*  
35 Suykens, J. A. K., & Vandewalle, J. (1998b). Nonlinear Modeling: advanced black-box  
36 techniques. In. Boston: Kluwer Academic Publishers.
- 37 Syntetos, A. A., Nikolopoulos, K., & Boylan, J. E. (2010). Judging the judges through accuracy-  
38 implication metrics: The case of inventory forecasting. *International Journal Of*  
39 *Forecasting*, 26, 134-143.
- 40 Tang, Z. Y., & Fishwick, P. A. (1993). Feed-forward Neural Nets as Models for Time Series  
41 Forecasting. *ORSA Journal on Computing*, 5, 374-386.
- 42 Tashman, L. J. (2000). Out-of-sample tests of forecasting accuracy: an analysis and review.  
43 *International Journal Of Forecasting*, 16, 437-450.
- 44 Terasvirta, T., van Dijk, D., & Medeiros, M. C. (2005). Linear models, smooth transition  
45 autoregressions, and neural networks for forecasting macroeconomic time series: A re-  
46 examination. *International Journal Of Forecasting*, 21, 755-774.
- 47 Weigend, A. S. (1994). *Time series prediction: forecasting the future and understanding the*  
48 *past. proceedings of the NATO Advanced Research Workshop on Comparative Time*  
49 *Series Analysis held in Santa Fe, New Mexico, May 14 - 17,1992* (1. printing ed.).  
50 Reading: Addison-Wesley.
- 51 Zadeh, L. A. (1965). FUZZY SETS. *Information and Control*, 8, 338-&.
- 52 Zellner, A. (1986). Bayesian-Estimation and Prediction Using Asymmetric Loss Functions.  
53 *Journal of the American Statistical Association*, 81, 446-451.
- 54  
55  
56  
57  
58  
59  
60



1  
2  
3  
4  
5  
6 Zhang, G. P., & Qi, M. (2005). Neural network forecasting for seasonal and trend time series.  
7 *European Journal Of Operational Research*, 160, 501-514.  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

For Review Only