

DIY forecasting

judgment, models & judgmental model selection



Fotios Petropoulos

Cardiff University, UK

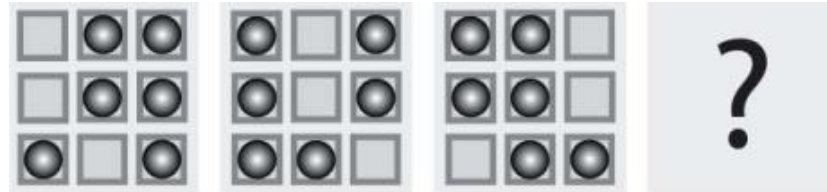
joint work with Nikolaos Kourentzes (Lancaster University, UK)
and Konstantinos Nikolopoulos (Bangor University, UK)

Judgment & forecasting

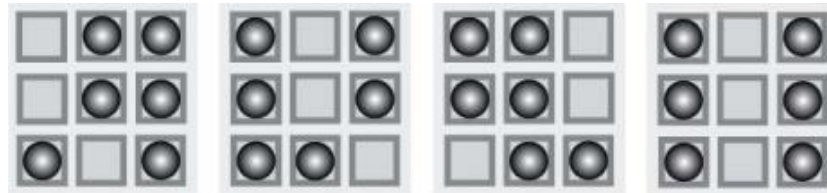
- Judgmental (point) forecasting
Lawrence et al., 1985; 1986; Lawrence & Makridakis, 1989; Sanders, 1992; Makridakis et al., 1993; Goodwin & Wright, 1993; 1994; ...
- Judgmental adjustments of a statistical baseline
Willemain, 1989; 1991; Mathews & Diamantopoulos, 1990; Goodwin & Fildes, 1999; Goodwin, 2000; Fildes et al., 2009; Syntetos et al., 2009; Franses & Legerstee, 2009; ...
- Judgmental probability forecasts and prediction intervals
Weinstein, 1982; Wright & Ayton, 1989; 1992; Onkal & Muradoglu, 1994; Eggleton, 1982; O'Connor & Lawrence, 1992; ...
- Improving judgmental forecasts: feedback, decomposition, combining, ...
Remus et al., 1996; Sanders, 1997; Goodwin & Fildes, 1999; Edmunson, 1990; Armstrong & Collopy, 1993; Lawrence et al., 1986; Blattberg & Hoch 1990; ...
- Judgmental model selection
Bunn & Wright, 1991

Forecasting with judgment: an IQ test analogy

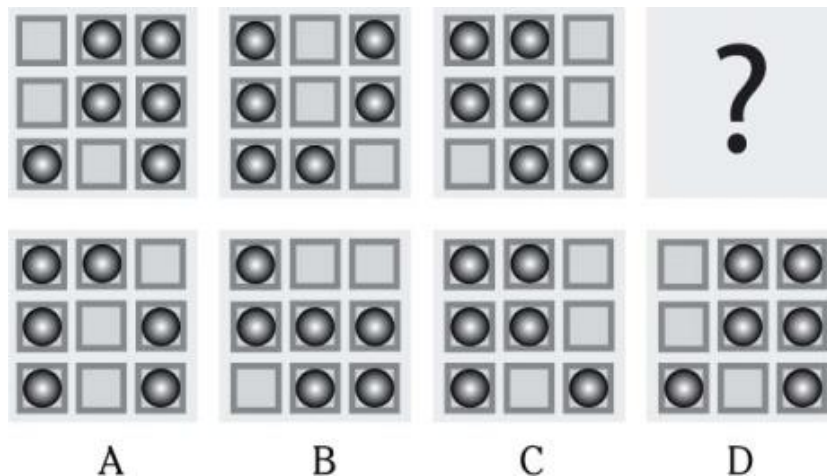
- Judgmental (point) forecasting



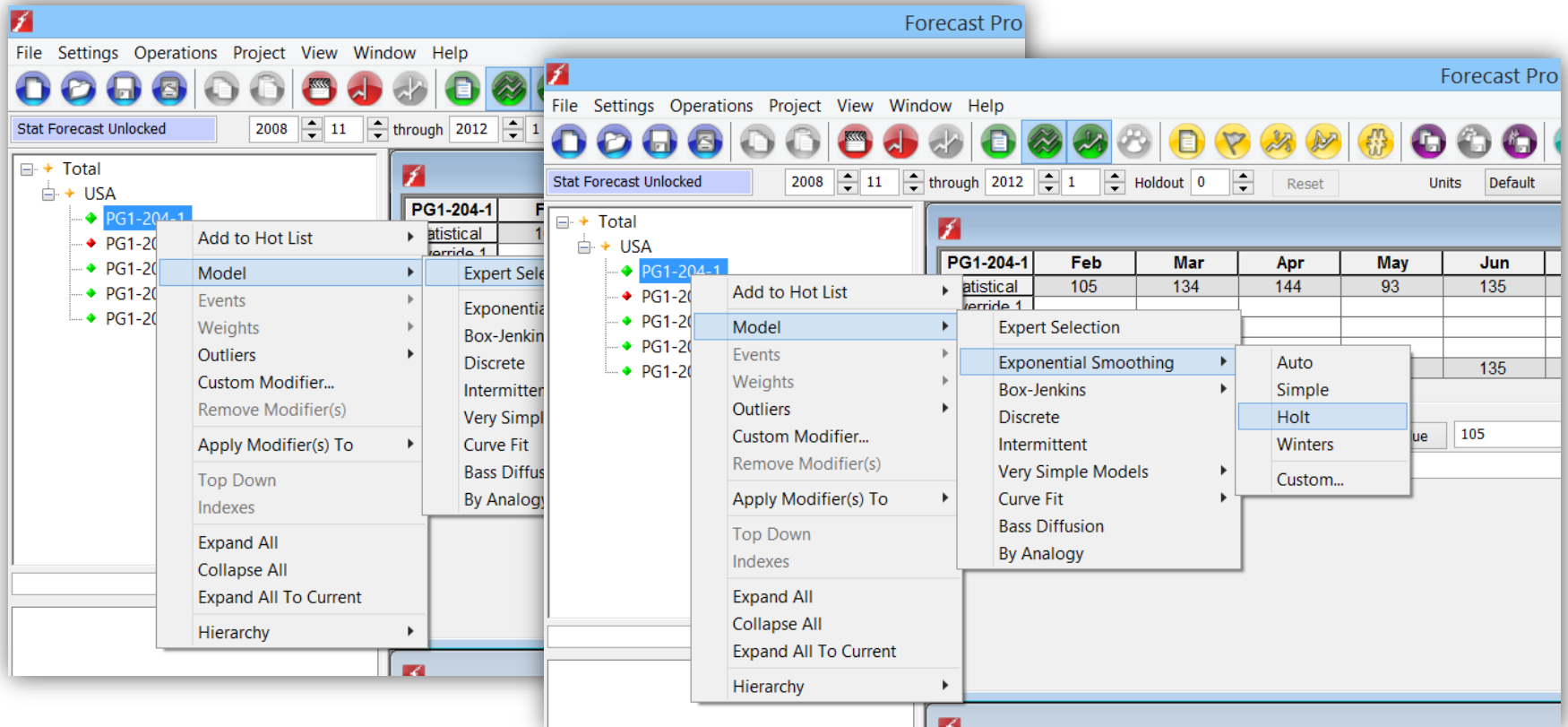
- Judgmental adjustments of a statistical baseline



- Judgmental model selection

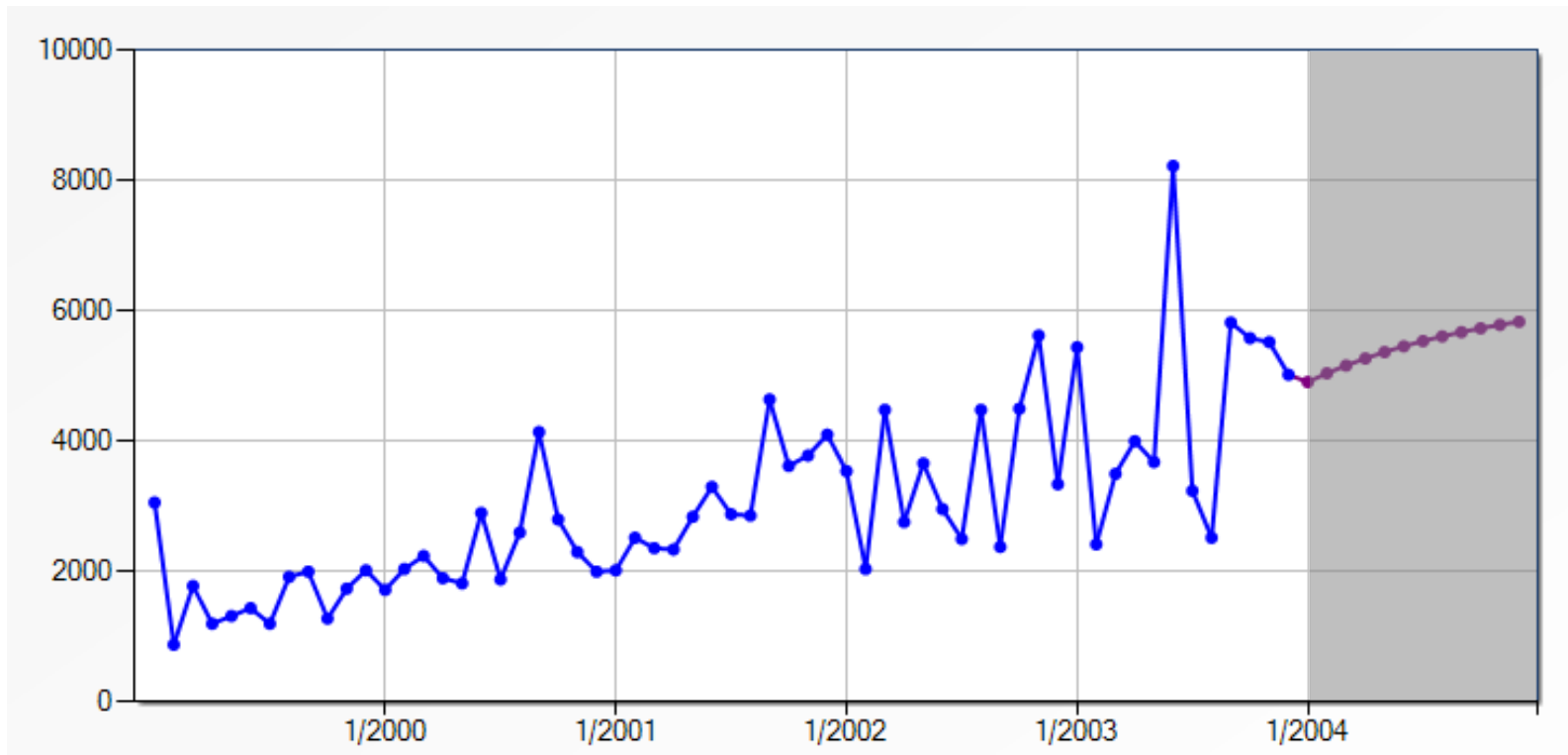


Model selection in a FSS



- What about judgment?
 - This strategy is implied by the majority of the world-leading FSSs.
 - However, an empirical investigation of how subjects perform in such tasks is a research gap.

Why do we expect to work?



- Statistical approaches cannot ex-ante assess the out-of-sample forecasts.
- Forecasters can select a method based on the quality of the out-of-sample forecasts.

Hypotheses

The **Brain**: Human Judgment = Judgmental Selection

The **Computer**: Forecasting System = Automatic Selection
based on Information Criteria

H1: Brain performs model selection differently than Computer.

H2: Brain is better in building models than selecting ones.

H3: Combination and aggregation will outperform both Brain
and Computer.

Laboratory experiment

Model Selection

Model Build

Model	Trend	Seasonality
Simple exponential smoothing (SES)	✗	✗
SES with seasonality	✗	✓
Damped trend	✓	✗
Damped trend with seasonality	✓	✓

Select the most appropriate model

Model A
 Model B
 Model C
 Model D

Submit
Press only once!

Select patterns that are applicable (if any):

Trend
 Seasonality

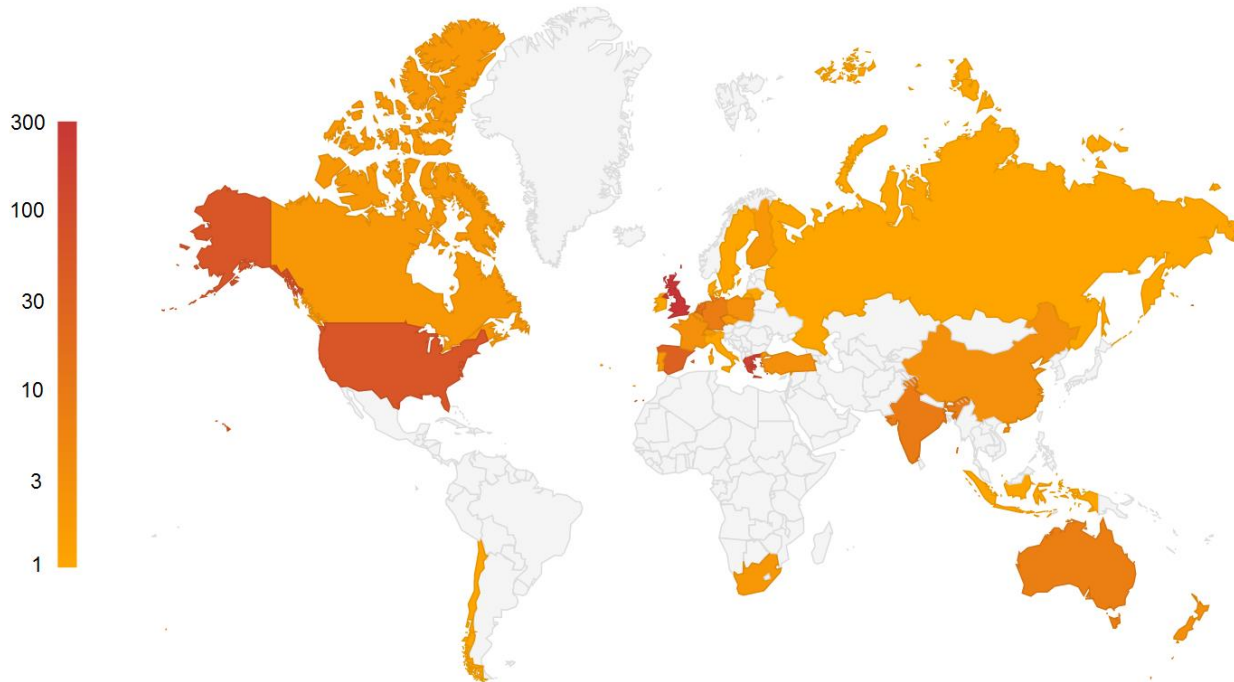
Submit
Press only once!

Each participant was randomly assigned in one of the two approaches and was asked to provide selections for 32 time series, based on different types of information.

Participants

Role	Model Selection	Model Build	Total
UG students	139	137	276
PG students	103	108	211
Researchers	13	31	44
Practitioners	46	44	90
Other	40	32	72

693
participants



Individual judgmental selections

Model A 0

Model B 0

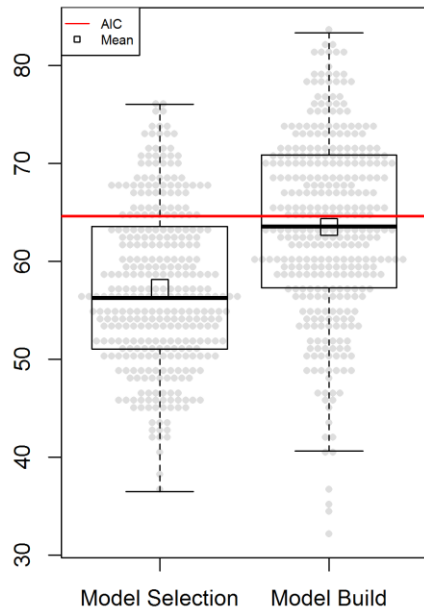


Model C 1

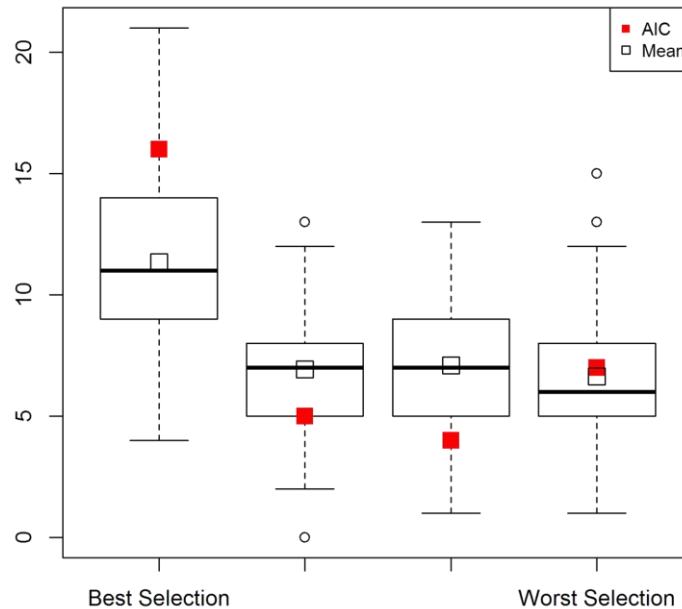
Model D 0

Selecting models judgmentally

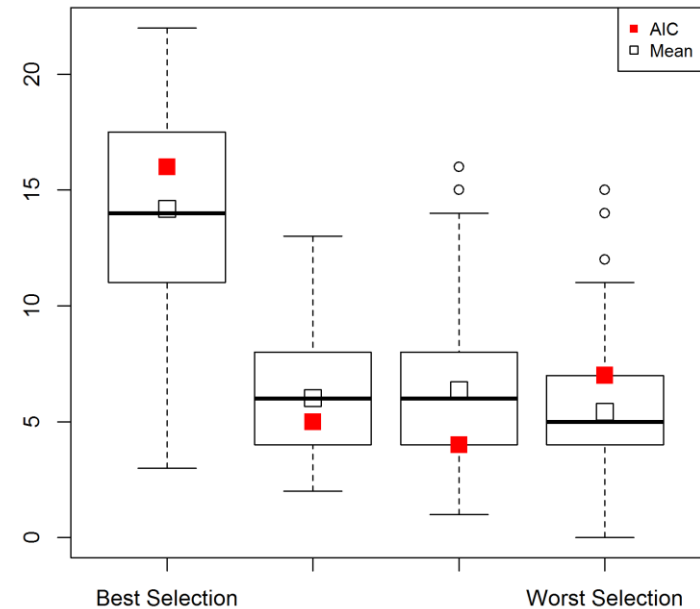
Scores



Model Selection



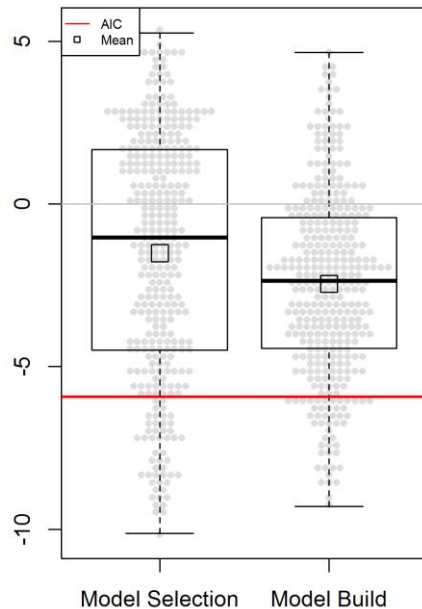
Model Build



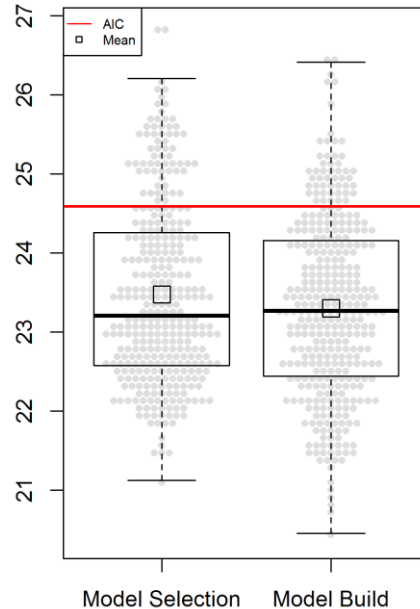
- Overall, humans' score is lower than statistics...
...while they select the ex-post best model less frequently.
- However, they do succeed in avoiding the worst model.
- How does this translate to forecasting performance?

Forecasting performance overall

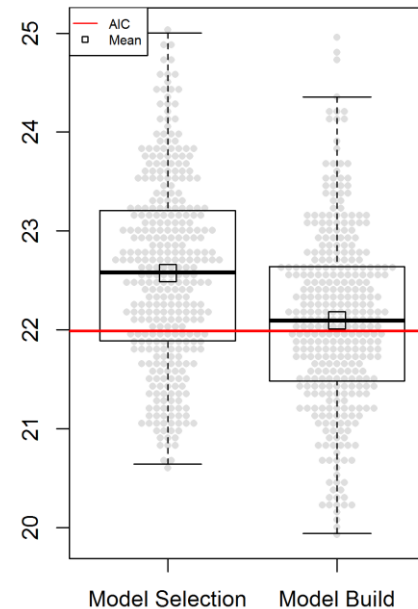
MPE (%)



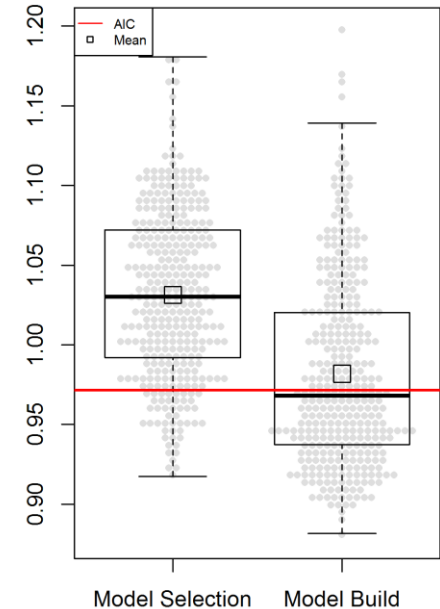
MAPE (%)



sMAPE (%)



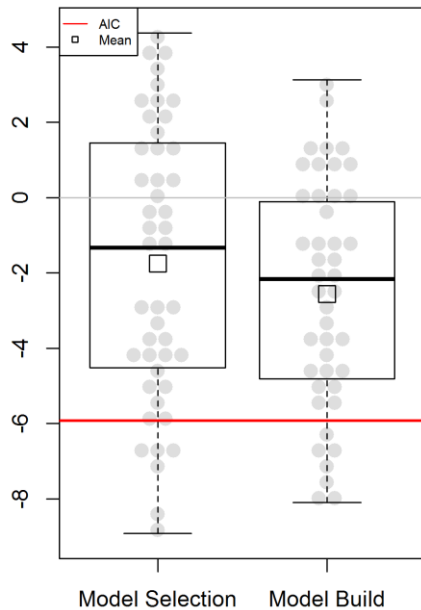
MASE



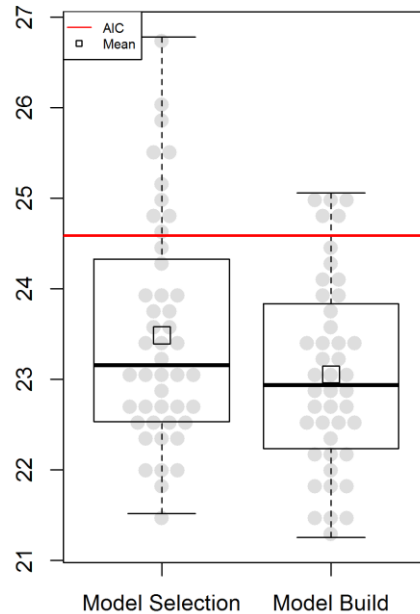
- In terms of bias and MAPE, humans perform significantly better than AIC.
- Participants in the Model Build experiment are as good as statistics, in terms of sMAPE or MASE.

Forecasting performance of practitioners

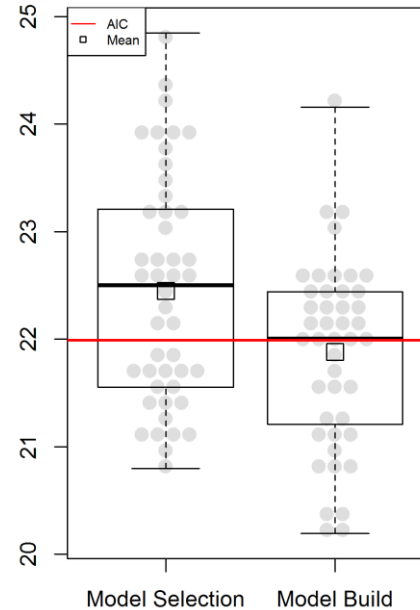
MPE (%)



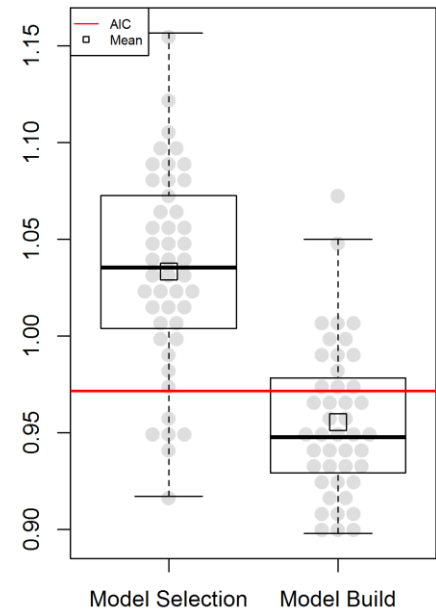
MAPE (%)



sMAPE (%)



MASE



- Practitioners on “model build” approach generally outperform the statistical model selection.

50% statistics + 50% manager [Blattberg & Hoch, 1990]



Model A

$$\frac{1}{2} = 0.5$$

Model B

$$\frac{0}{2} = 0$$



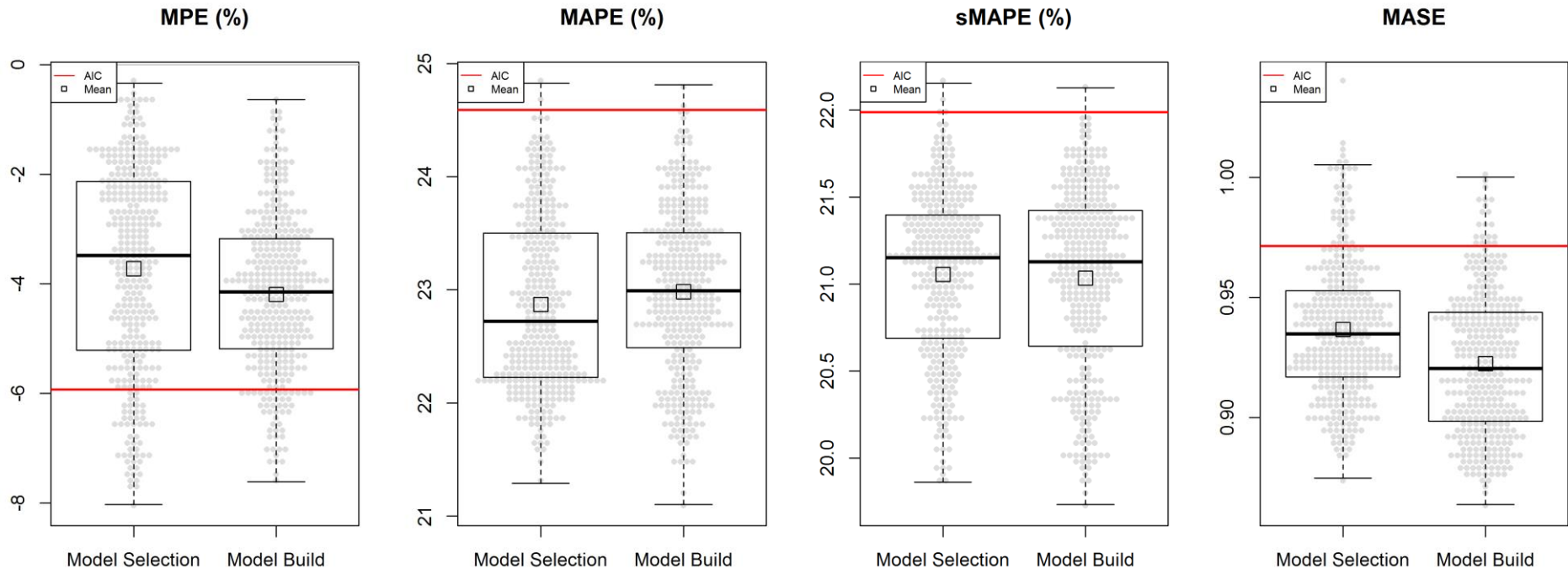
Model C

$$\frac{1}{2} = 0.5$$

Model D

$$\frac{0}{2} = 0$$

50% statistics + 50% manager: results



- The Blattberg-Hoch approach works for 86% of the cases for bias, 99% of the cases for MAPE and sMAPE and for 90% of the cases for MASE.
- The differences in the performance between the two approaches (model selection and model build) are also minimised.

Wisdom of crowds



Model A

$$\frac{6}{20} = 0.3$$



Model B

$$\frac{8}{20} = 0.4$$



Model C

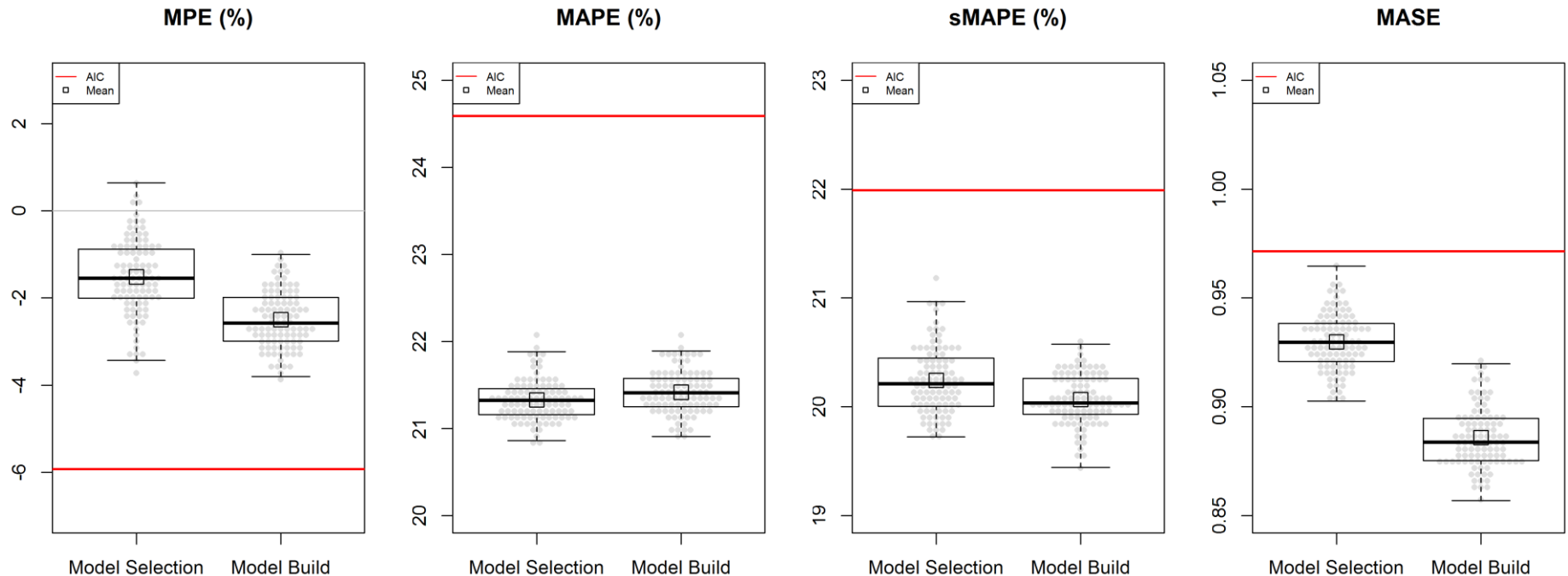
$$\frac{2}{20} = 0.1$$



Model D

$$\frac{4}{20} = 0.2$$

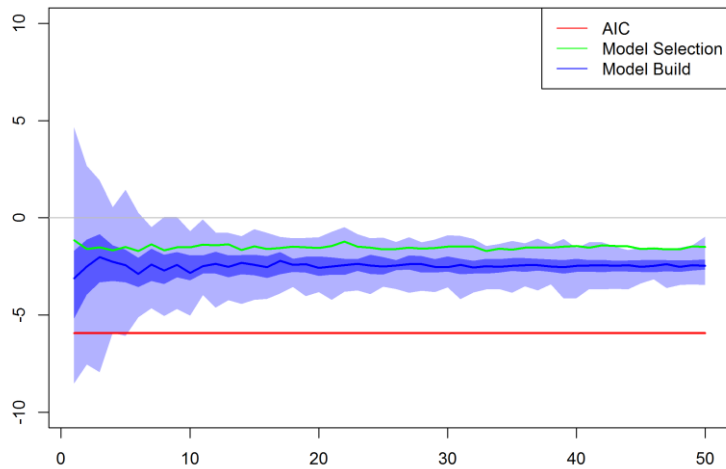
Wisdom of crowds: results



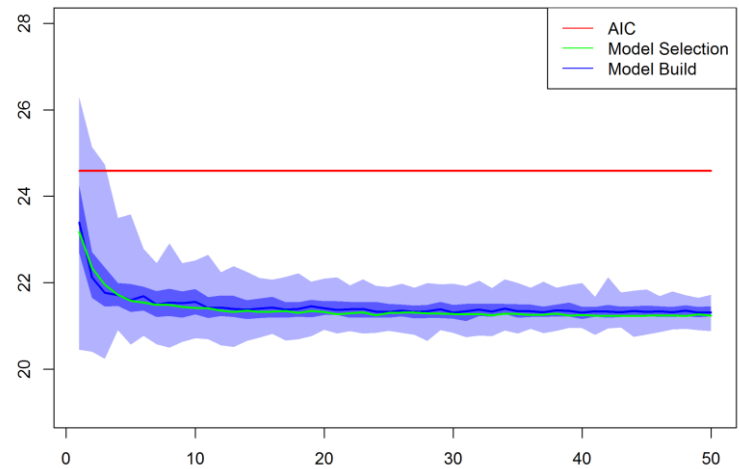
- 20 experts: the forecasting performance of a grouped judgmental model selection approach is significantly better than statistical model selection.
- How many experts are enough?

Wisdom of crowds: results (cont'd)

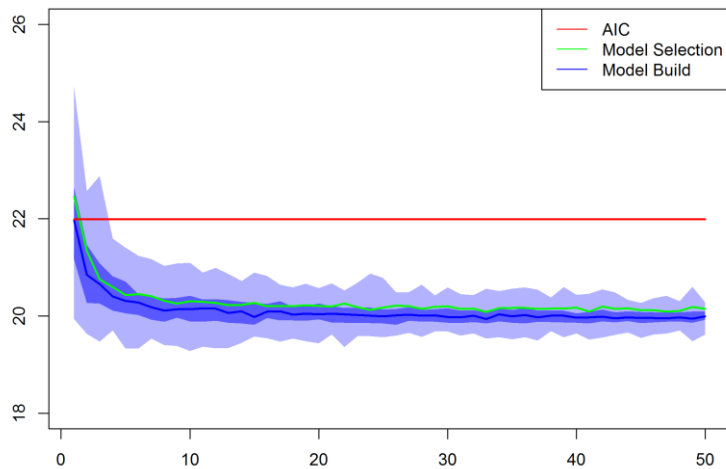
MPE (%)



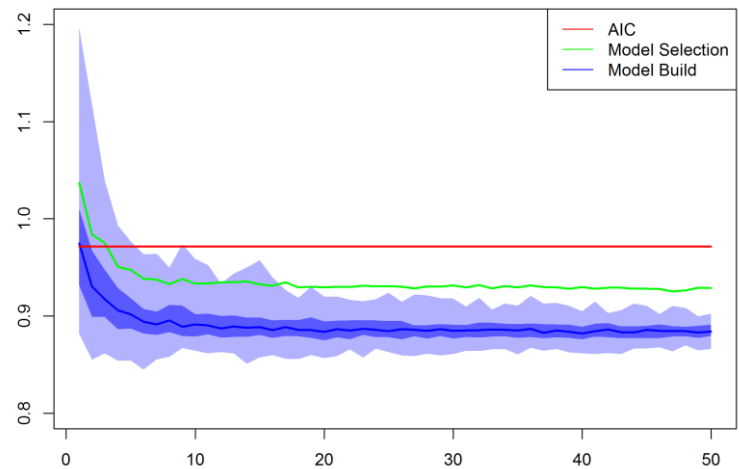
MAPE (%)



sMAPE (%)



MASE

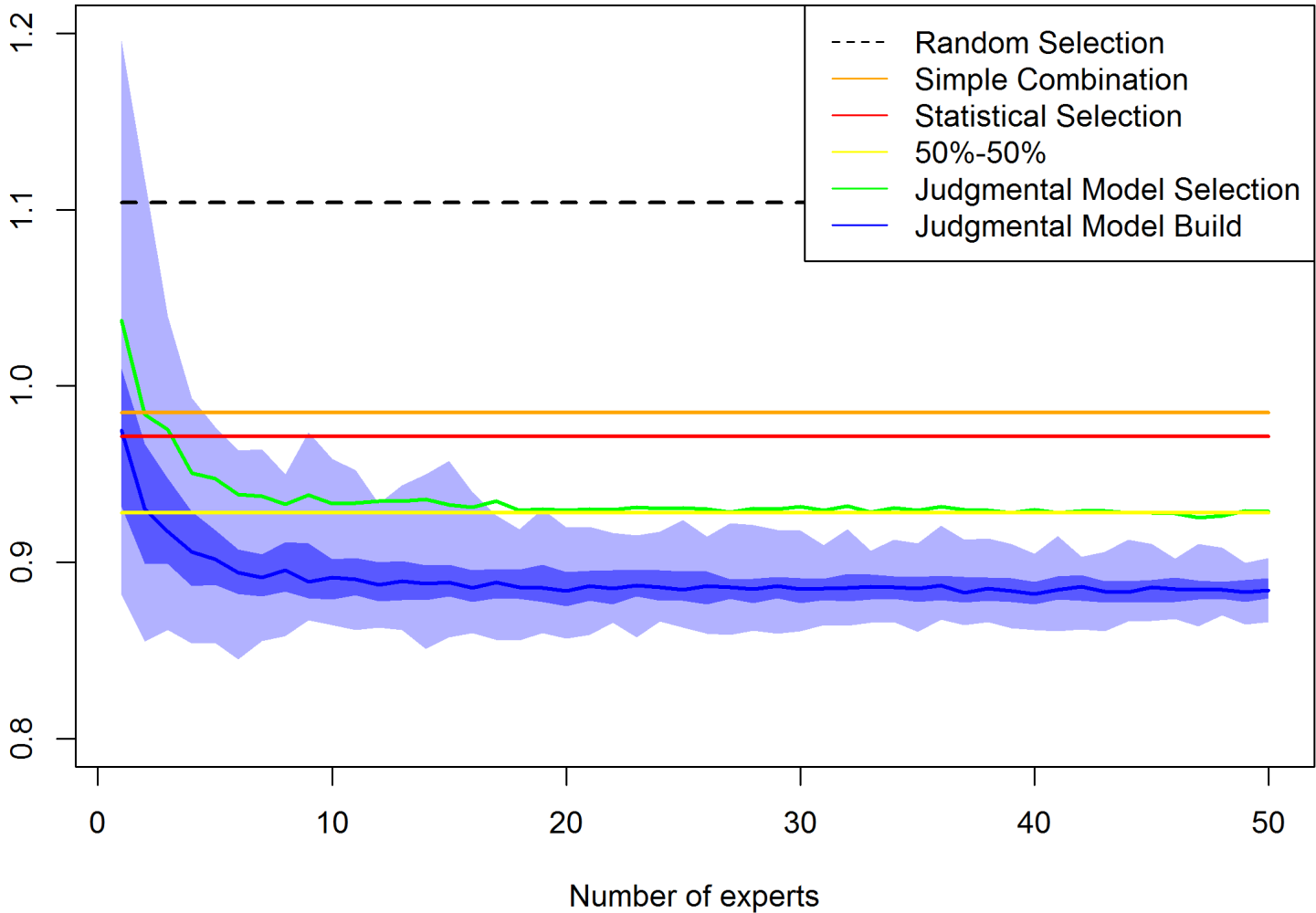


Number of experts

Number of experts

Summary of results

MASE

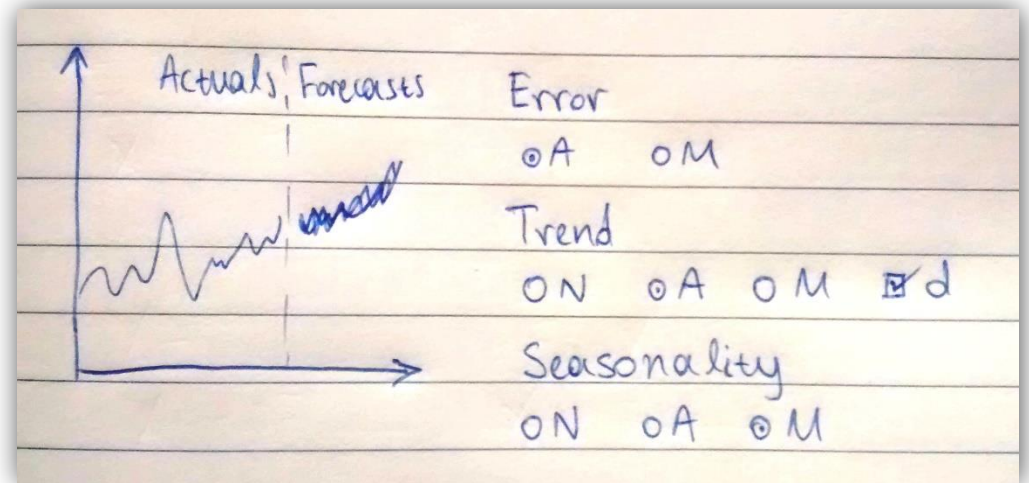


Conclusions

- Judgmental model selection is offered by every FSS, but its performance has never been empirically evaluated before.
- Judgmental model selection and, especially, model build may offer improvements over a statistical selection strategy.
- The improvements are more apparent when we focus on the participants self-described as practitioners.
- Wisdom of crowd (grouped judgmental model selection) or a 50%-50% combination approach appear to be very promising.

To do... and extensions

- Wisdom of crowds & decision trees: automatically derive optimal weights to emulate humans' selection strategy.
- Next experiment: judgmental build of ETS



- Next experiment: judgmental selection of model's parameters

Questions?

PetropoulosF@cardiff.ac.uk

<http://fpetropoulos.eu>

Acknowledgments

- Sylwia Macinska (Summer Intern at LUMS) for performing an initial literature review on judgmental forecasting studies.
- Vassilios Assimakopoulos, Philip Hans Franses, Sinan Gonul, Kesten Green, Rob Hyndman, Qinyun Li, Pamela Stroud, Aris Syntetos, and Juan Ramon Trapero Arenas for helping in recruiting participants.
- Robert Fildes and Paul Goodwin for their useful comments.
- FP and NK are supported by a Lancaster University Early Career Research Grant (MTA7690).