



Contents lists available at ScienceDirect

International Journal of Forecasting

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

A neglected dimension of good forecasting judgment: The questions we choose also matter



Edgar C. Merkle^{a,*}, Mark Steyvers^b, Barbara Mellers^c, Philip E. Tetlock^c

^a University of Missouri, United States

^b University of California, Irvine, United States

^c University of Pennsylvania, United States

ARTICLE INFO

Keywords:

Forecast evaluation
Brier score
Question selection
Forecast nonresponse

ABSTRACT

Forecasters are typically evaluated via proper scoring rules such as the Brier score. These scoring rules use only the reported forecasts for assessment, neglecting related variables such as the specific questions that a person chose to forecast. In this paper, we study whether information related to question selection influences our estimates of forecaster ability. In other words, do good and bad forecasters tend to select questions in different ways? If so, can we capitalize on these selections when estimating forecaster ability? We address these questions by extending a recently-developed psychometric model of forecasts to include question selection data. We compare the extended psychometric model to a simpler model, studying its unidimensionality assumption and highlighting the unique information that it can provide. We find that the model can make use of the fact that good forecasters tend to select more questions than bad forecasters, and we conclude that question selection data can be beneficial above and beyond reported forecasts. As a side benefit, the resulting model can potentially provide unique incentives for forecaster participation.

© 2017 International Institute of Forecasters. Published by Elsevier B.V. All rights reserved.

1. Introduction

The issue of question selection is of considerable importance in many areas of forecasting. Does a forecaster look good because he/she chooses to forecast only easy questions? Should we reward forecasters for attempting difficult questions, even if their forecasts for those questions are poor? Is forecasting ability related to question choice; that is, are good forecasters better at selecting questions on which they will excel? These questions cannot be answered via classical metrics such as proper scoring rules

(e.g., [Gneiting & Raftery, 2007](#)), which generally assume that all forecasters report on all questions.

Instead of proper scoring rules, model-based approaches to forecast evaluation make it feasible to study issues that relate to question selection. One prime candidate is a recently-proposed psychometric model of probabilistic forecasts ([Merkle, Steyvers, Mellers, & Tetlock, 2016](#)), which is related to the previously-proposed item response models for doubly-bounded variables ([Bejar, 1977](#); [Ferrando, 2001](#); [Müller, 1987](#); [Muthén, 1989](#); [Noel & Davier, 2007](#); [Samejima, 1973](#)). This model simultaneously provides estimates of forecaster ability and of question difficulty and discrimination. For example, if a particular question is worded ambiguously, good forecasters' judgments may be indiscernible from bad forecasters' judgments. The model can recognize this, discounting forecasters' judgments on ambiguous questions as we estimate their general abilities across questions. Conversely, certain ques-

* Correspondence to: Department of Psychological Sciences, University of Missouri, Columbia, MO 65211, United States.

E-mail address: merkle@missouri.edu (E.C. Merkle).

tions may be particularly suitable for distinguishing between forecasters of different abilities. Forecasters' judgments on these good questions would then be weighted more heavily than those on other questions.

In addition to addressing novel substantive issues, model-based forecaster assessment allows us to make explicit our assumptions related to missing forecasts. That is, excluding (or including) question selection data from a model implicitly makes assumptions as to why forecasters do not respond to some questions. For example, the *missing completely at random* (MCAR; e.g., Little & Rubin, 2002) assumption says that missingness is independent of the data (both observed and unobserved). This assumption, which is unlikely to be fulfilled in practice, generally implies that we can ignore missing data.

The models of Merkle et al. (2016) instead employed the *missing at random* assumption, where all observed forecasts (even those from forecasters with incomplete data) are used for model estimation. The MAR assumption states that the probability of missingness can be predicted exclusively from the observed data; our predictions would not improve if we could observe the missing data. This assumption excludes the possibility that forecasters of greater/lesser abilities will differ in their frequency of responding or in the types of questions that they choose. When forecaster ability is related to question selection, models that employ the MAR assumption may lead to suboptimal substantive conclusions regarding forecaster ability or question attributes.

To study question selection issues in this paper, we develop a psychometric model of forecasts that accommodates question selection and reported forecasts jointly. This model draws on the psychometric literature on explicitly modeling (as opposed to ignoring) missing data (e.g., Chang, Tsai, & Hsu, 2014; Holman & Glas, 2005; O'Muircheartaigh & Moustaki, 1999; Rose, von Davier, & Xu, 2010; Wang, Jin, Qiu, & Wang, 2012), which explores the idea that information can be gained from missing data in standardized testing contexts. Following model development, we apply the model to data from a recent forecasting tournament. This allows us to study a major model assumption in relation to the unidimensionality of forecasting ability, and it also allows us to compare the proposed model with a previous model that employs the MAR assumption. In addition, we compare the model-based estimates with other forecaster ability estimates that are based on the Brier score, and we illustrate the general use of question selection data for forecaster assessment.

In what follows, we begin by providing the technical detail of the models, starting with previous developments and continuing on to novel developments. Next, we apply the model to data from a recent forecasting tournament. The application includes an examination of model assumptions, a small example that provides readers with an intuition of the model's estimates, and a larger example involving the full dataset. Finally, we report on a simulation that illustrates further the benefits of modeling question selection data.

2. Models

Assume that each of I forecasters responds to some subset of J questions, with each forecaster's subset possibly being unique. Let y_{ij}^* be forecaster i 's probit-transformed forecast for the realized outcome of question j (with the possibility that it is missing), and let d_{ij} be a 0/1 variable indicating whether or not y_{ij}^* is missing (0 for missing, 1 otherwise). We begin by briefly reviewing the MAR model proposed by Merkle et al. (2016), and we then introduce a new model that handles d_{ij} in addition to y_{ij}^* .

2.1. MAR model

The models described by Merkle et al. (2016) focus on the observed y_{ij}^* , providing estimates of forecasters' abilities and questions' difficulties and discriminations. Because they focus exclusively on the observed y_{ij}^* , they employ the MAR assumption.

Most of the concepts underlying the Merkle et al. (2016) model are derived from the classical item response literature (e.g., Embretson & Reise, 2000; Lord & Novick, 1968; McDonald, 1999), with the application to probabilistic forecasts being relatively novel. That is, instead of being applied to binary data reflecting whether or not a student answers a test question correctly (say), the models are applied to probability judgments. The model can be written as

$$y_{ij}^* | t_{ij}, \theta_{a,i}, d_{ij} = 1 \sim N(\mu_{ij}, \sigma_j^2) \quad (1)$$

$$\mu_{ij} = \beta_{0j} + (\beta_{1j} - \beta_{0j}) \exp(-\beta_2 t_{ij}) + \lambda_j \theta_{a,i} \quad (2)$$

$$\theta_{a,i} \sim N(0, 1), \quad (3)$$

where t_{ij} is the time at which person i forecasts question j (measured as days until the question expires), $\theta_{a,i}$ is person i 's forecasting ability (the a subscript stands for "ability"), and the β_j and λ_j parameters are related to item j 's difficulty and discrimination, respectively.

This model is related to a factor analysis model, but with extra parameters (the β s) that allow the question difficulty to change over time. This is necessary because forecasters often report on a question at different points in time, and available, relevant information changes over time. For example, imagine two forecasters predicting the chance of rain for February 1. A forecaster who responds on January 31 will have a natural advantage over a forecaster who responds on January 28 because the question is easier on January 31. The model can take this into account by allowing the difficulty to change over time, based on the way in which the full group of forecasters responds over time.

Merkle et al. (2016) used Bayesian methods to fit the above model to data from a forecasting tournament (data from the same source as is used in this paper, described further below), and found that (i) the model could predict out-of-sample forecasts successfully; (ii) the forecaster ability estimates were the forecaster ability estimates were more highly related to a forecaster's future ability, as compared to the Brier score; and (iii) the item parameter estimates were related to external covariates in the way

that was theoretically expected. The next section extends this model to handle question selection data, resulting in a model that allows for missing not at random (MNAR) data.

2.2. MNAR model

The MNAR model allows for the possibility that missing data provide information about forecaster abilities (and about item attributes), over and above the observed data. It is a generalization of the above model that accounts for the missingness indicators d_{ij} and the reported forecasts y_{ij}^* simultaneously. In developing the model, we adopted an approach similar to those of Holman and Glas (2005) and O’Muircheartaigh and Moustaki (1999), both of whom studied methods for handling missing data in traditional item response contexts. For each person i , we model $2 \times J$ variables simultaneously: the probit-transformed forecasts for the J items (y_{ij}^* , some of which are missing), and the missingness indicators for the J items (d_{ij}).

The J forecast variables are all modeled in a manner similar to that of Merkle et al. (2016):

$$y_{ij}^* | t_{ij}, \theta_{a,i}, d_{ij} = 1 \sim N(\mu_{ij}, \sigma_j^2) \tag{4}$$

$$\mu_{ij} = \beta_0 j + \beta_1 j t_{ij} + \lambda_{j,1} \theta_{a,i}. \tag{5}$$

This is a simplification of the Merkle et al. (2016) model, where the “time” covariate has a linear influence on μ_{ij} instead of an exponential curve. This function is simpler than the exponential function while still allowing for curvilinear for potential curvilinear influences of time on the reported forecasts (because we are modeling the probit-transformed forecasts, as opposed to the original forecasts). We identify this part of the model by fixing a single $\lambda_{j,1}$ parameter (in $j = 1, \dots, J$) to 1.

In addition to the forecast variables, the J missingness indicators are handled via a two-factor model

$$d_{ij} | \theta_{a,i}, \theta_{r,i} \sim \text{Bernoulli}(p_{ij}) \tag{6}$$

$$\text{probit}(p_{ij}) = \beta_{0,(j+j)} + \lambda_{(j+j),1} \theta_{a,i} + \lambda_{(j+j),2} \theta_{r,i}, \tag{7}$$

where $\theta_{r,i}$ is person i ’s response propensity. This equation implies that a person’s forecasting ability can play a role in both the questions that he/she selects and the forecasts that he/she reports. There is also a response propensity factor that accounts for a person’s general level of activity in making forecasts. The subscripts above are based on the fact that the missingness variables can be treated as new questions within the model. That is, for person i , questions 1 to J include the reported forecasts, while questions $(J + 1)$ to $2J$ include the binary missingness indicators (for completeness, we define the parameters $\lambda_{1,2}$ to $\lambda_{J,2}$ to all equal zero). We identify this part of the model by fixing a single $\lambda_{(j+j),2}$ parameter (where j is in $1, \dots, J$) to 1.

In addition to the above constraints, parameter identification is completed by assuming that

$$\theta_i = (\theta_{a,i} \theta_{r,i})' \sim N(\mathbf{0}, \mathbf{D}_\psi), \tag{8}$$

where \mathbf{D}_ψ is a diagonal covariance matrix with unique entries ψ_1 and ψ_2 . The assumption of diagonality here could potentially be relaxed, but parameter constraints would be required elsewhere in the model to trade off with this relaxation. Preliminary testing indicated that

the model without diagonality was slow to converge, so we did not pursue it further in this paper. Holman and Glas (2005) show that parameter estimates under the above constraints can be transformed linearly to parameter estimates under the alternative constraints (with the diagonality assumption relaxed), meaning that the parameter estimates from the two approaches are related to each another.

2.3. Estimation

The model can be represented as a path diagram, illustrated in Fig. 1. Each forecaster potentially contributes $2J$ observed variables: a forecast and a selection indicator for each of the J questions. These observed variables are shown in the boxes labeled forecast₁ to forecast _{J} and select₁ to select _{J} . The former consist of probit-transformed forecasts, with each forecast variable being observed only if the corresponding select variable is equal to 1. For example, a forecaster only supplies forecast₁ if select₁ equals 1.

The path diagram further shows the two latent variables labeled “forecast ability” and “response propensity”, with “forecast ability” influencing both the reported forecasts and the question selections. In terms of notation, the λ parameters represent the paths from the latent variables to the observed variables, the θ parameters represent the latent variables, and the β parameters (corresponding to the time covariate) are excluded for the sake of simplicity (we would have a unique β parameter for each observed variable, cluttering the diagram).

To incorporate the time covariate in the model and to easily obtain θ estimates, we rely on Bayesian methods of model estimation. Specifically, we employ Markov chain Monte Carlo (MCMC) methods, adopting an approach that is similar to existing MCMC methods for estimating psychometric models (e.g., Ghosh & Dunson, 2009). We use the following prior distributions on classes of model parameters (subscripts are absent because the same prior was used on each free parameter):

$$\beta_0 \sim N(0, 2) \tag{9}$$

$$\beta_1 \sim N(0, 2) \tag{10}$$

$$\lambda \sim N(0, 1) \tag{11}$$

$$\psi \sim \text{Gamma}^{-1}(0.01, 0.01) \tag{12}$$

$$\sigma^2 \sim \text{Gamma}^{-1}(0.01, 0.01), \tag{13}$$

where the second parameter of each normal distribution is a variance, as opposed to a precision.

These priors were intended to place a high density in sensible parameter ranges, which can improve model convergence and sampling efficiency. The parameter ranges are sensible because the model parameters are generally used to make predictions on the probit scale, meaning that the predictions are akin to z-scores. Thus, we would be surprised to observe values of β_0 or β_1 that were drastically outside $(-2, 2)$, because these values would correspond to extreme probabilities near 0.025 and 0.975, respectively. We would also be surprised to observe values of λ that were much larger than 1, given the diverse questions in our dataset (further discussion below). Finally,

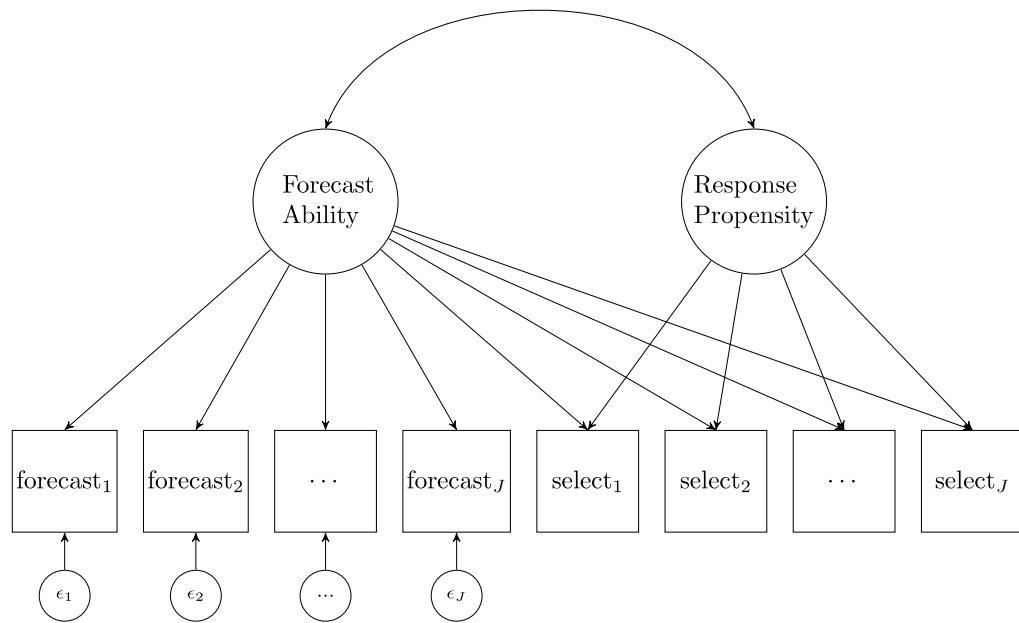


Fig. 1. Path diagram of the proposed model. The time covariate t_{ij} is excluded for simplicity.

the priors on ψ and σ^2 are traditional, noninformative priors on variance parameters.

Unless otherwise mentioned, the models were burned in for three chains of 2000 iterations each, after which parameters were sampled for an additional 2000 iterations each. Chain convergence was fast and was monitored using time series plots and the Gelman–Rubin potential scale reduction statistic (Gelman & Rubin, 1992).

2.4. Parameter interpretation

The model parameters supply many pieces of information about relationships between forecasting problems, forecaster abilities, and forecaster selection. In particular, the model allows us to address the following questions (the relevant parameters are given in parentheses):

- Which questions are more popular than others? ($\beta_{0,J+1}$ to $\beta_{0,2J}$).
- Who are the frequent/infrequent responders? ($\theta_{r,i}$).
- Do good forecasters tend to select/avoid certain questions? ($\lambda_{J+1,1}$ to $\lambda_{2J,1}$).
- Do frequent forecasters tend to select/avoid certain questions? ($\lambda_{J+1,2}$ to $\lambda_{2J,2}$).

The first two issues can be addressed easily by examining the raw data (i.e., response proportions), but the last two issues are more difficult to address via simple, data-based metrics. This is an advantage of the model-based approach described here.

Along with the above topics, the new model can also address the same issues that the Merkle et al. (2016) model addressed. These include:

- Which questions are easier/harder than others? ($\beta_{0,1}$ to $\beta_{0,J}$).
- Who are the better/worse forecasters? ($\theta_{a,i}$).

- Are some questions better than others for distinguishing between forecasters of varying abilities? ($\lambda_{1,1}$ to $\lambda_{J,1}$).

The applications below focus on the $\theta_{a,i}$ parameters, examining the ways in which the estimated forecaster abilities change from the MAR model to the MNAR model. While we eventually fit the model to a large dataset, we initially fit the model to data from only four questions, because this makes it easier to illustrate the model's behavior. First, however, we describe the data source and study the extent to which model assumptions are fulfilled.

3. Application: geopolitical forecasting

The forecasts used in this paper arise from a four-year geopolitical forecasting tournament sponsored by IARPA. The tournament involved five research teams, each of which was required to forecast hundreds of diverse questions related to world events. Example questions include:

- Will Australia formally transfer uranium to India by 1 June 2012?
- Will Mario Monti resign, lose re-election/confidence vote, or vacate the office of Prime Minister of Italy before 1 January 2013?
- Will there be a significant outbreak of H5N1 in China in 2012?
- Will the Yuan to Dollar exchange rate on 31 December 2012 be more than 5% different than the 31 August 2012 exchange rate?

For each question, the research teams elicited forecasts from large groups of individuals. The teams then aggregated the forecasts via statistical methods and reported them to the funder on a daily basis.

We focus here on assessing individual forecasters who were part of the winning team in the tournament (the

Good Judgment Project). This team collected forecasts from thousands of individuals, each of whom was active for at least one of the four tournament years. We begin by providing some background details on the dataset (see also Mellers, Stone, Atanasov et al., 2015; Mellers, Stone, Murray et al., 2015; Mellers et al., 2014), and we then discuss issues of dimensionality in relation to the dataset. The dimensionality issues are important because the proposed model makes specific assumptions here.

3.1. Data

Adult forecasters of all ages were recruited from across the United States via email lists, professional societies, university organizations, and social media. The forecasters voluntarily logged on to a website and selected the questions that they wished to forecast. Forecasters were motivated to participate in various manners, including monetary payments for participation and leaderboards of the best forecasters.

For each question, forecasters read the question details and reported a probability of event occurrence (from 0 to 1 inclusive; forecasts of exactly 0 and 1 were transformed to 0.001 and 0.999, respectively, for modeling). Forecasters were randomly assigned to experimental conditions that differed in whether, e.g., the forecasters worked individually (vs on teams) and the types of training the forecasters received. For the purposes of this paper, we ignore experimental conditions and model only individuals' reported forecasts and question selections. This is facilitated by the fact that even forecasters who worked in teams still reported their own individual forecasts.

Below, we use a data set consisting of 775 forecasters who each report on a subset of 157 binary (event occurs/does not occur) questions. To speed model estimation, forecasters were initially included if they responded to 70 or more questions; we later apply the model to forecasters with sparser data. While the forecasters were free to respond to the same question multiple times (i.e., to update their forecasts), we used only the first forecast supplied on a given question for simplicity.

3.2. Unidimensionality

The model studied in this paper assumes a single "forecaster ability" dimension and a single "response propensity" dimension, with the reported forecasts being influenced only by the "ability" dimension and the question selections being influenced by both dimensions. The assumption of a single "forecaster ability" dimension is almost certainly violated for the application considered here, which involves forecasts of diverse world events. For example, we could imagine a forecaster having expertise on a specific topic like European politics, so that his/her forecasts are better on questions related to that topic than on other, unrelated questions. The proposed MNAR model would assign this forecaster a single ability estimate, representing some combination of his/her ability at forecasting European politics and his/her ability at forecasting other questions. However, this single estimate will

not be an optimal assessment of the forecaster's true ability, which requires two dimensions to describe it fully (one for European politics and one for other questions). Likewise, if a forecaster's ability improves over time, his/her single model estimate will not reflect this. However, use of a single model estimate does mimic applied forecaster assessments where the Brier score is averaged indiscriminately across all available questions (e.g., Carvalho, 2016).

In addition to mimicking practical assessments, we can draw on the psychometric literature to explicitly assess dimensionality. Researchers here have pointed out that strict unidimensionality will not hold in practice, even in the case of standardized educational tests (e.g., Reise, Scheines, Widaman, & Haviland, 2013; Thissen, 2016; Zhang, 2007). Thus, considerable effort has been devoted to assessing the magnitude of the unidimensionality violation, as opposed to assessing whether or not such a violation occurs (e.g., Bonifay, Reise, Scheines, & Meijer, 2015; Stout et al., 1996; van Abswoude, van der Ark, & Sijtsma, 2004; Zhang, 2007). This effort has provided metrics that can tell us whether or not a set of questions is "unidimensional enough" to be useful. The metrics are nonparametric in nature, because model-based assessments tend to be overly sensitive to minor violations of unidimensionality.

One of the most popular metrics in this body of literature, which we adapt to forecasting data here, is called DETECT (Zhang, 2007; Zhang & Stout, 1999). This metric makes use of the fact that, for unidimensional tests, all nonzero covariances/correlations between questions should be due to the single, underlying ability dimension. Thus, partial covariances/correlations between questions (conditioning on the single underlying dimension) should all equal zero. While this is the idea underlying DETECT, the specific algorithm is more complex than mere covariance calculation. Further computational details are provided in Appendix A.

The DETECT index is useful for our purposes because previous researchers have provided rules of thumb for its interpretation. Roussos and Ozbeck (2006) state that values below 0.2 are often taken to represent approximate unidimensionality, whereas values greater than 1.0 are taken to represent strong multidimensionality. As we move from 0.2 to 1.0, multidimensionality increases in strength. Thus, for the unidimensionality assumption to be useful, we should look for D values below 1.0, with values closer to 0 being better.

We computed this statistic separately for the reported forecasts y and the question selections d . For the question selections, the DETECT index indicated strong multidimensionality, achieving a maximum value of 2.3 at two clusters (subgroups) of questions. These subgroups had a strong temporal component: when we re-computed the index using only data from a single year, the maximum DETECT value was 0.58 (indicating only moderate multidimensionality). For the reported forecasts, we obtained a maximum DETECT statistic of 0.72 at three subgroups.

These results provide some evidence that, for this particular dataset, multidimensionality is moderate and results from changes in the forecasters over time, as opposed to forecasters having specific expertise or interest in particular question topics. To address these findings, we later fit the model to subsets of the data arising from a single year of the tournament and compare it to a model fitted to the full dataset.

Table 1
Brier scores and response rates for four questions.

Question	Mean Brier score	Response rate
1067	0.08	0.87
1106	0.01	0.77
1147	0.09	0.59
1177	0.07	0.59

4. A simple example

For an initial investigation of the proposed model's behavior, we use data from only four questions. The four questions used here (with identification numbers in parentheses) were all open during 2012–2013; they are:

- Will Traian Basescu resign, lose referendum vote, or vacate the office of President of Romania before 1 April 2012? (1067)
- Will Kim Jong-un resign or otherwise vacate the office of Supreme Leader of North Korea before 1 April 2013? (1106)
- Before 1 April 2013, will the Egyptian government officially announce it has started construction of a nuclear power plant at Dabaa? (1147)
- Will Mohammed Morsi cease to be President of Egypt before 1 April 2013? (1177)

In the tournament, 771 of the 775 forecasters in our dataset responded to at least one of the four questions. We use all of the data supplied by these 771 forecasters, including missing observations. Below, we further describe the questions and the model before examining the results.

4.1. Data summary

The questions' mean Brier scores and response rates (out of the number of people who responded to any of the four questions) are displayed in Table 1; scatterplots and the distributions of reported forecasts are displayed in Fig. 2; and response pattern frequencies and mean Brier scores are displayed in Table 2. Questions 1067 and 1147 had the worst Brier scores, and Questions 1147 and 1177 were less popular than the other two. Fig. 2 (and most notably the panels for Question 1067) also shows that there is some overuse of "nice" numbers like 0.5, which indicates that some participants might be reporting 0.5 to reflect complete uncertainty, as opposed to the actual probability of the event occurring. Our model does not account for this phenomenon, and it is unclear whether accounting for it is worth the additional model complexity that would be required.

Finally, Table 2 shows that 345 forecasters responded to all four questions, while 173 forecasters only responded to the first question (1067). The forecasters who responded only to question 1067 appear to have worse Brier scores than the other forecasters, although this result is clouded by differences in question difficulty and in response pattern frequencies. The estimated model, described below, can help to provide a clearer assessment of these issues.

Table 2
Response pattern frequencies and mean Brier scores for the simple example.

Response pattern	Frequency	Mean Brier
0001	1	0.023
0010	1	0.000
0011	1	0.006
0100	15	0.052
0101	11	0.063
0110	8	0.027
0111	64	0.056
1000	173	0.102
1010	1	0.061
1100	87	0.043
1101	31	0.052
1110	34	0.066
1111	345	0.060

Note: the four numbers in the "Response pattern" column correspond to questions 1067, 1106, 1147, and 1177, respectively, and equal 0 for question nonresponse and 1 otherwise.

4.2. Results

Several results are notable when we examine the estimated IRT model of forecasts and question selections. We start with the λ parameters that describe the influence of forecaster ability on the reported forecasts and on question selection. We then move to the forecaster ability estimates.

The λ parameters that relate to question discrimination ($\lambda_{1,1}$ to $\lambda_{4,1}$) are all close to 1, which (unsurprisingly) means that better forecasters tended to do better on all four questions. Perhaps more surprisingly, better forecasters were more likely to select certain questions (based on $\lambda_{5,1}$ to $\lambda_{8,1}$). This was particularly true for the two questions that had lower response rates and worse Brier scores, 1147 and 1177. Question 1106 showed forecaster ability to have a smaller influence on question selection, while question 1067 showed virtually no influence.

Fig. 3 compares the ability estimates from the MAR model of Merkle et al. (2016) (x -axis; note that this model includes a linear effect of time that is similar to Eq. (5)) with those from the new model of forecasts and question selection (y -axis). Each point represents a single forecaster, with the point's color and shape representing the total number of questions answered (out of a possible four). We can roughly see three diagonal lines going from the bottom left to the top right: one line of red circles and green triangles, one line of blue squares, and one line of purple pluses. The red circles and green triangles tend to be closest to the top, indicating that the forecasters who answered three or four questions generally received the highest ability estimates under the MNAR model, followed by those who answered two questions, followed by those who answered only one question. The MNAR model penalizes non-responders automatically, because non-responders tend to supply worse forecasts than frequent responders.

The figure also includes a small number of forecasters who stand out; one such forecaster in the middle of the plot is circled. The circled forecaster answered only one question, but obtained a higher ability estimate under the new model than similar people who responded to all four questions. The forecaster responded only to the question

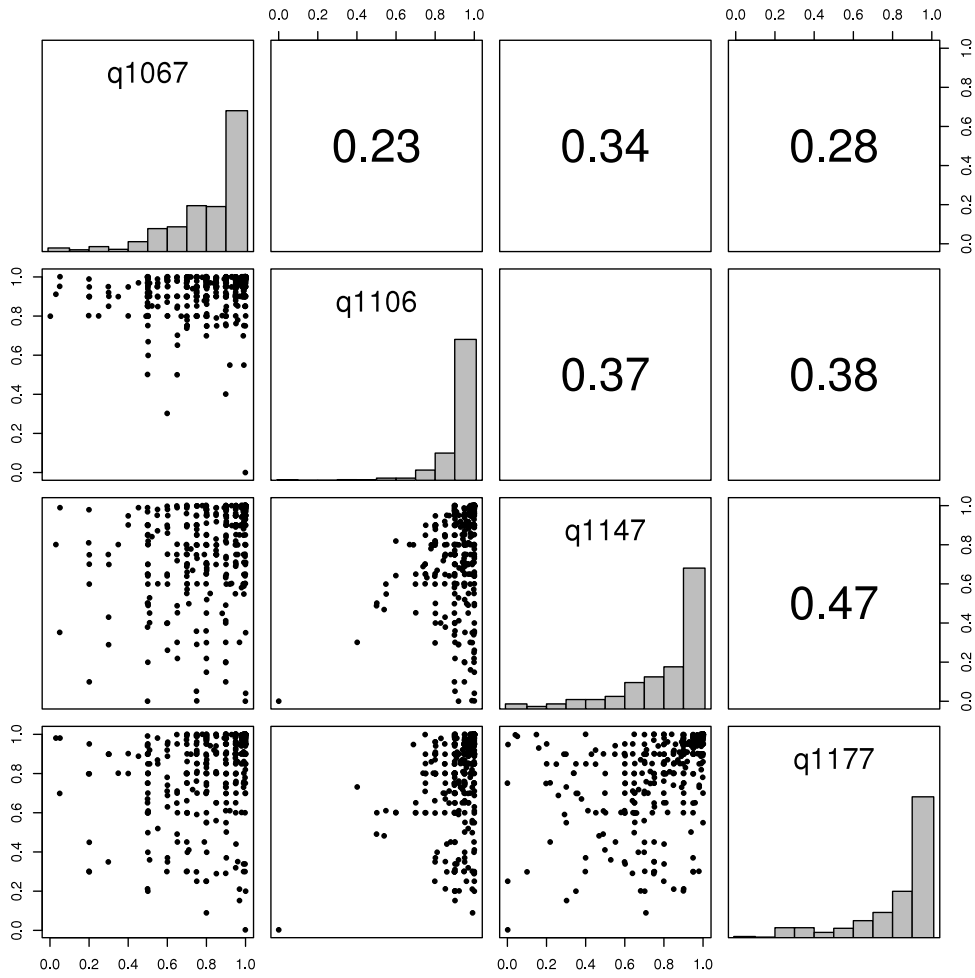


Fig. 2. Visual summaries of the forecasts for each question's realized outcome for the simple example. The upper triangle displays the Pearson correlations associated with the scatterplots in the lower triangle.

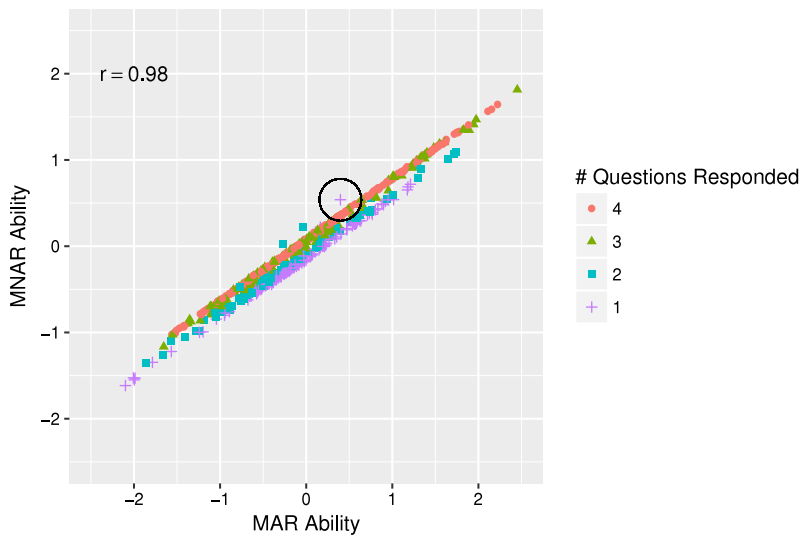


Fig. 3. Comparison of MAR ability estimates with MNAR ability estimates that incorporate question selection, for the simple example. The Spearman correlation appears in the upper left.

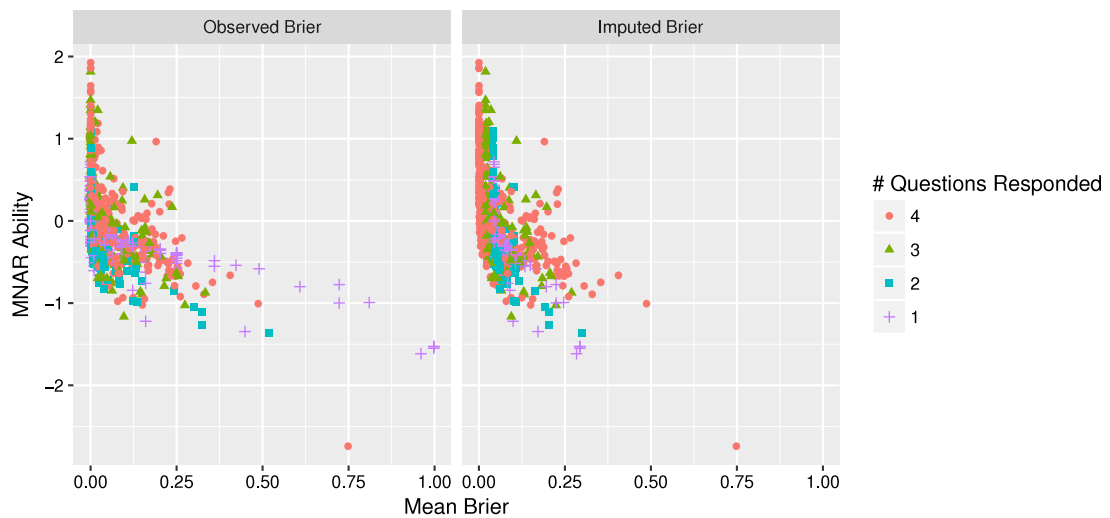


Fig. 4. Mean observed Brier score (x -axis, left panel) and mean imputed Brier scores (x -axis, right panel) versus MNAR ability estimates.

that was most highly associated with forecaster ability (question 1147), which is a very uncommon response pattern: this is the only person who responded to question 1147 and no others. In addition, the forecaster made a near-perfect forecast of 0.99 in favor of the realized outcome on that question. Thus, the model rewards the forecaster for making a good forecast on the question that was most highly associated with good forecasting. This reward is relative to the person's MAR ability estimate; that is, the person's new ability estimate is still in the middle of the pack, relative to the full set of forecasters. In order to obtain the highest ability estimate, a forecaster must report exceptional forecasts on most or all of the questions. This is because the shrinkage of each forecaster's ability estimate is related to the amount of data available on a forecaster: a forecaster's ability estimate can become more extreme as he/she responds to more questions.

Fig. 4 compares the new model's ability estimates to two types of Brier scores: a mean observed Brier score, and a mean imputed Brier score. These reflect heuristic methods for handling missing data while still using a proper scoring rule. The mean observed Brier score for each forecaster averages Brier scores only across the questions to which the forecaster responded (similar to treating missing forecasts as "not reached", for example). The mean imputed Brier score, on the other hand, fills in the missing observations with the corresponding question's mean Brier score based on the observed forecasts for that question (similar to treating missing forecasts as "incorrect", for example).

In Fig. 4, the x -axis reflects the Brier scores and the y -axis reflects the model estimates. For reference, the red circles in the two figures are exactly the same, as Brier score imputing has no impact on people who forecast all four questions. The figure shows that the ability estimates from the model are generally related to the Brier scores, with correlations in the range -0.6 to -0.7 . Comparing the two panels, we see that the Brier score imputing helped many people with bad Brier scores. In the left panel, these

people are generally closer to the right side of the x -axis, with points that are triangles, squares, or pluses. In the right panel, they have all moved further left on the x -axis (improved), while those who responded to all four questions have stayed in the same locations. Perhaps the most striking result of this figure involves the fact that we observe multiple vertical "lines" of points. This shows that the model assigns different ability estimates to people who receive nearly the same Brier scores. The specific questions that were selected, along with the times when the forecasts were reported, are responsible for these differences.

5. Full dataset

Now that we have illustrated the model's application to a small number of questions, we fit the model to the larger data set of 775 forecasters responding to 157 questions (again maintaining only the first forecast reported by each person on each question). We focus on comparing the MNAR model with the Merkle et al. (2016) model that does not handle question selection. This comparison provides information about the impact of the "missing at random" assumption on model estimates.

A comparison of the Merkle et al. (2016) ability estimates (missing at random) and the new ability estimates (missing not at random) is displayed in Fig. 5. The points are now displayed in various shades of blue, depending on the response rate: blue points represent forecasters who responded to nearly all of the questions, while black points represent those who responded to fewer questions. The figure clearly shows that the response rate influences the ability estimates in the new model: forecasters who received similar ability estimates under the old model can now receive very different estimates from the new model. The extent to which the new ability estimates change depends on response rate: light blue points are always closest to the top of the graph, and darker points are lower. Just like in the simple example, the extent to

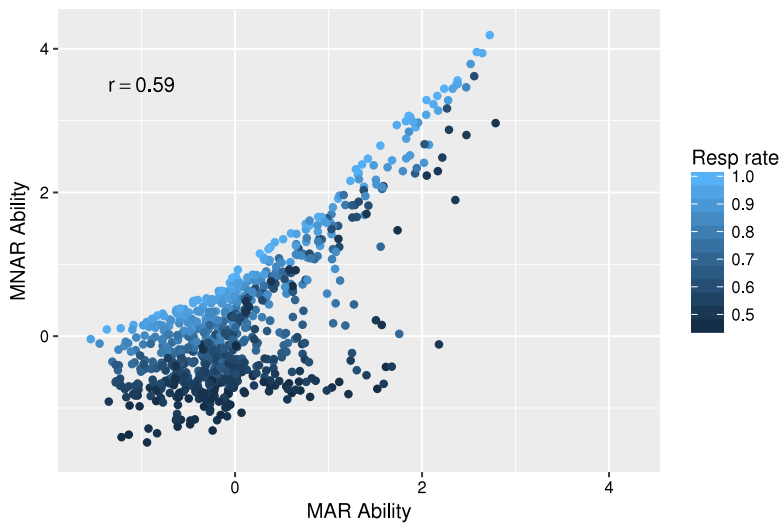


Fig. 5. Comparison of the MAR ability estimates with the MNAR ability estimates obtained using data across all four years of the tournament. The Spearman correlation appears in the top left.

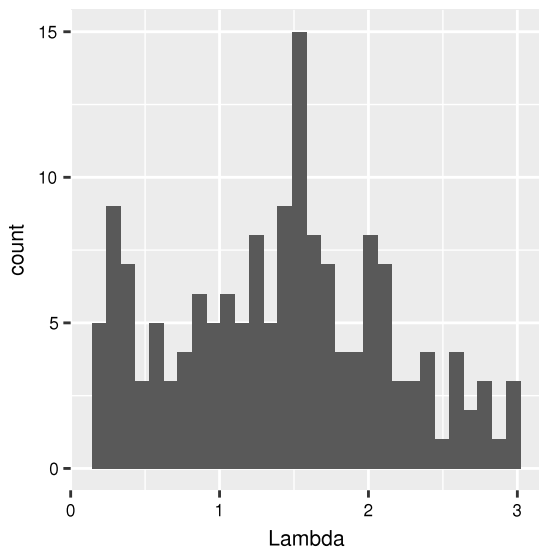


Fig. 6. Histogram of λ estimates that correspond to paths from the “forecaster ability” latent variable to the “question selection” variables.

which the darker points are penalized depends on the specific questions to which forecasters responded: if a “low response rate” forecaster responds to many questions that good forecasters select, then that forecaster is penalized less. On the other hand, a “low response rate” forecaster who responds in other ways will have a larger penalty.

Fig. 6 displays a histogram of λ estimates that correspond to paths from “Forecaster ability” to the question selection variables (see Fig. 1). These estimates provide information as to whether good forecasters are more/less likely to select certain questions. The histogram indicates that the chance of responding to any given question increases with the forecaster’s ability, regardless of that

question’s difficulty. This result has at least two further implications. First, there is a deviation from the MAR assumption, because the MAR model is obtained when all of these λ parameters are equal to zero. Second, a forecaster can improve his/her ability in two ways: by reporting good forecasts and by responding to a large number of questions. This may be especially useful for forecast consumers, in that the model developed here can give forecasters an incentive to increase their response rates. We return to this issue in the general discussion.

Finally, the histogram in Fig. 6 indicates question variability: some of the estimates are close to zero, indicating that forecaster ability is nearly unrelated to the selection of certain questions, whereas some of the other estimates are far from zero. Table 3 shows a few specific questions that fell at each extreme. The bottom section contains questions whose λ estimates were near zero, indicating that their selection was “unrelated to ability”. These questions were all open near the start of the tournament, when people were first getting accustomed to forecasting. Of these people, some became good forecasters and others dropped out, likely explaining the model results. Conversely, the top section contains questions whose selections were “related to high ability”. These questions were open later in the tournament, and they comprise less-popular topics that beginning forecasters may have avoided.

6. Impact of dropouts

As was shown in an earlier section, the multidimensionality in forecasts and question selection is related in part to the fact that the forecasting tournament was divided into four separate years. Many existing forecasters dropped out at the end of each year, and many new forecasters entered for the subsequent year. Thus, the results in the previous section (see Fig. 5) were influenced by two types of missing

Table 3
Notable questions illuminated by the model estimates.

Questions related to high ability	Question text
1174	Will the Turkish government release imprisoned Kurdish rebel leader Abdullah Ocalan before 1 April 2013?
1177	Will Mohammed Morsi cease to be President of Egypt before 1 April 2013?
1183	Will the United Nations Security Council pass a new resolution directly concerning Iran between 17 December 2012 and 31 March 2013?
Questions unrelated to ability	Question text
1004	Will the United Nations General Assembly recognize a Palestinian state by 30 September 2011?
1010	Will the 30 Sept 2011 “last” PPB for Nov 2011 Brent Crude oil futures exceed \$115?
1022	Will the South African government grant the Dalai Lama a visa before 7 October 2011?

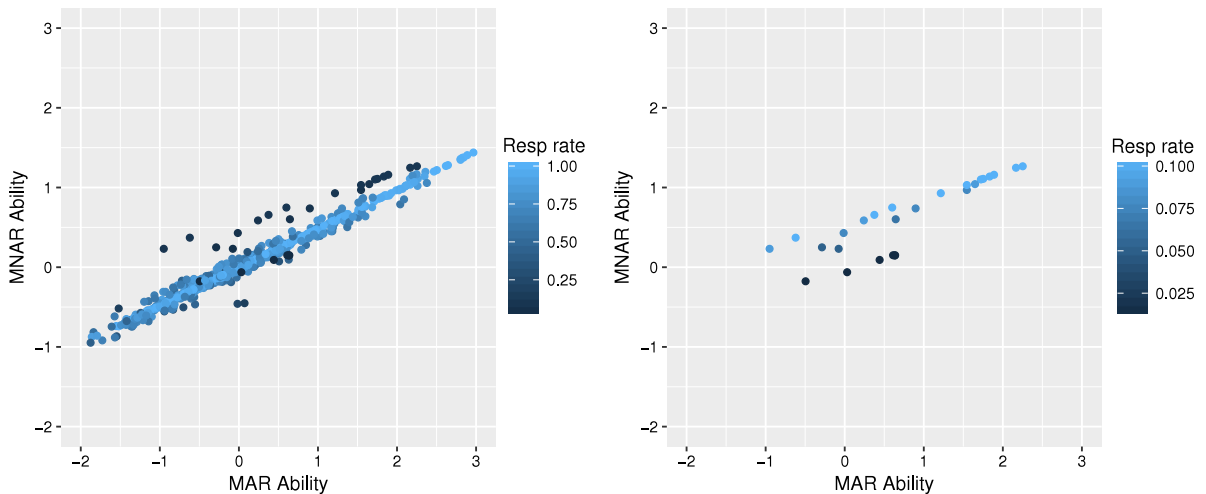


Fig. 7. Comparison of MAR and MNAR ability estimates during Year 1 (left panel). The right panel contains a subset of points in the left panel from infrequent responders.

data: dropouts who only reported forecasts during a subset of the tournament, and selective responders who forecasted a subset of questions across the entire tournament.

These two groups may influence the model differently, given that the best forecasters (the “superforecasters;” see [Mellers, Stone, Murray et al., 2015](#)) tended to continue reporting forecasts throughout the tournament, while worse forecasters were more likely to drop out. The fact that bad forecasters dropped out more often implies that bad forecasters had more missing data, so that the model learned to penalize forecasters with low response rates. However, if we can avoid the bad forecasters who dropped out after each year, the model may penalize/reward forecasters differently. Thus, this section fits the model to Year 1 forecasts only, which eliminates year-to-year dropout effects from our analysis.

6.1. Method

We fit the model to 771 forecasters who made at least one forecast during Year 1. This is a subset of the original data and includes some forecasters with very sparse data

(who reported infrequently during Year 1 but more frequently during subsequent years). We restricted ourselves to 78 questions that both opened and closed during Year 1.

6.2. Results

The left panel of [Fig. 7](#) contains the main results, with the MAR ability estimates on the x-axis and the MNAR ability estimates on the y-axis. The light blue points form a diagonal line, showing that forecasters who responded to nearly all of the questions receive similar ability estimates across models (except for some rescaling). There are also a few darker points that cluster around the main diagonal line, showing that some people who responded to fewer questions received small rewards or penalties depending on their response patterns. Aside from this, we see a small number of dark points that are farther from the diagonal line, with many of these points receiving higher ability estimates under the MNAR model.

The dark points above the line represent people who responded to a small number of questions that tended to be selected by good forecasters. The right panel of [Fig. 7](#) provides a closer look at the dark points from the left

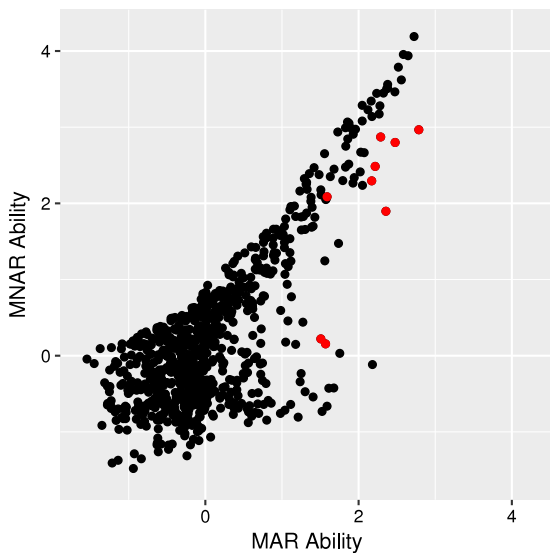


Fig. 8. Display of the Year 1 infrequent responders in the full dataset (red points), relative to other forecasters. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

panel. The right panel contains forecasters who responded to no more than 10% (seven) of the questions, so that the shading reflects response rates that go from 0 to 0.1 instead of from 0 to 1. It can be seen that the nonresponders who received the largest boosts generally answered six or seven questions (close to 10% of the questions). In general, these questions were ones that good forecasters tended to answer, and the nonresponders reported good forecasts on them. The model deemed this sufficient evidence to give the forecasters a boost.

Do these forecasters deserve the boost? We answer this question by looking at the performance of the forecasters on the full dataset (including data from other years). We focus on the nine nonresponders in the right panel of Fig. 7 whose ability estimates from the new models were greater than 0.8. We then re-create Fig. 5 in Fig. 8, except that the nine nonresponders are now highlighted in red. It can be seen that, when we compare the forecasters on abilities across years, those who originally received a boost now receive a penalty. This is likely because the nonresponders had larger amounts of missing data across years. Despite this finding, the nonresponders who originally received boosts during Year 1 all remain in the top half of forecasters, with seven of the nine being above the 90th percentile in terms of ability. This suggests that the new model can help us to identify good forecasters who have responded to only a small number of questions. We explore this suggestion further in the next section.

7. Improvements in ability estimates

While the previous sections have illustrated that the new IRT model rewards/penalizes (non)response in an intuitive fashion, ultimately we wish to know whether the resulting ability estimates are better than those of the

model that employs the missing at random assumption. This issue is more complex than it appears initially because it requires us to explicitly define what we mean by “ability”. For example, imagine that we estimate forecaster ability using three metrics: the mean Brier score, the MAR model, and the MNAR model. It is likely that, if we compute each of these metrics on a training sample, they will be correlated most highly with the analogous metric on a test sample: the training Brier score will be most correlated with the test Brier score, the training MAR estimates will be most correlated with the test MAR estimates, and the training MNAR estimates will be most correlated with the test MNAR estimates. In order to say which model is best, we need to decide explicitly which metric counts as the “official” measure of ability. This amounts to dealing with the validity of each of the ability metrics (e.g., Borsboom, Meltenbergh, & van Heerden, 2004), which is a difficult topic to address in the current context.

We sidestep validity issues here, showing that, if we provide the MNAR model only with the questions that some forecasters selected (and not with their actual forecasts), those forecasters’ ability estimates are related to the estimates that would be obtained if we used the full dataset. This implies that information can be obtained from the item selections, independently of the reported forecasts. This, in turn, illustrates the utility of the proposed model in practice.

7.1. Method

We conducted a simulation study of the MNAR model, using only data from Year 1 of the forecasting tournament. Similarly to the previous section, this was done to prevent the model from capitalizing on year-to-year dropout effects. For each of 100 replications, we randomly selected 25% of the 775 forecasters in the data and deleted all of their forecasts. We maintained the questions that these forecasters selected (i.e., the d_{ij}), however, fitting the model to these selections along with all of the data provided by the remaining 75% of forecasters. Following model estimation (2000 burn-in samples followed by 2000 posterior draws), we computed the posterior mean ability estimates of the forecasters whose forecasts were deleted. Finally, we examined the relationships between these ability estimates and those associated with the forecasters’ full data from Years 1 to 4. We included the data from Years 2 to 4 in our comparison because it served as a more stringent generalizability measure. That is, because data from Years 2–4 were completely held out of the initial model estimation, it is more impressive if the resulting ability estimates are correlated with the estimates that include data from Years 2–4.

7.2. Results and discussion

Fig. 9 contains a histogram of correlations between (i) the ability estimates resulting from the Year 1 deleted dataset (where 25% of the forecasters had only question selection data), and (ii) the ability estimates resulting from the model developed in this paper (utilizing reported forecasts and question selections from all four years). The

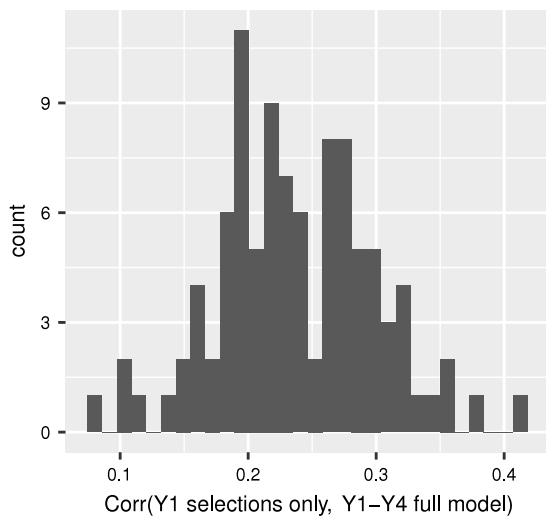


Fig. 9. Simulated correlations between ability estimates under two models: a model that only uses question selections from Year 1, and a model that uses both reported forecasts and question selections from Years 1–4.

histogram depicts only 98 correlations, as the model failed to converge for two of the 100 simulation replications. This is likely due to bad, randomly-generated initial values in these replications.

The histogram shows that the ability estimates from the two models are correlated positively across all replications, with a mean correlation of 0.24 and an interquartile range of (0.2, 0.28). This result provides evidence that the question selections contain information that is related to the full ability estimates (which would be obtained if we included the reported forecasts in the model).

This result is weakened by the fact that the question selection data were included in both models; we might expect a positive correlation between the models' estimates because they were partially based on the same data. To explore this criticism, we also examined the relationship between the “Year 1, question selection” ability estimates and the MAR ability estimates (based on the model of Merkle et al., 2016). The latter model utilizes only the reported forecasts from Years 1–4 (without question selection data), so that the forecasters with deleted data contribute unique data points to each model. These correlations are nearly always positive (in 97 of 98 replications), with a mean of 0.12 and an interquartile range of (0.08, 0.16). This mean (and range) is lower than those of the correlations from Fig. 9, potentially illustrating the impact of repeating the data across models. However, given that the correlations remain positive, we conclude that the question selection data contains useful information. This information may not always lead to major, practical improvements in ability estimates, but it is worth considering in scenarios where forecasters are free to select their own questions.

8. Discussion

In this paper, we first developed a psychometric model that allows us to assess forecasters' abilities while si-

multaneously handling data on question selection. This is potentially useful in situations where forecasters are free to select the questions that they wish to forecast, so that the selected questions provide information about forecasting ability above and beyond the forecasts reported on those questions. After developing the model and checking its assumptions, we illustrated the extent to which the proposed model differed from a previous model that did not account for question selection. The results from the new model implied that good forecasters tended to select more questions, regardless of the question difficulty, and that some specific question selections had an influence on forecaster ability estimates. We also studied the extent to which we can estimate forecaster ability based on question selections alone (not forecasts), finding that these ability estimates exhibited correlations of 0.24 (on average) with the full data ability estimates. This implies that there is information in the question selections that can be capitalized upon, a result that has also been studied in other contexts (e.g., Rubin & Steyvers, 2009). In the remainder of this paper, we provide further ideas on missingness mechanisms, relationships to traditional scoring rules, model assumptions, and methods of model estimation.

8.1. Missingness mechanisms

One appeal of the proposed MNAR model is the fact that it handles missingness in a manner that agrees with intuition: good forecasters select questions differently from bad forecasters (in the specific context of the current data, good forecasters generally selected more questions than bad forecasters), and we should be able to use these differences in a forecaster assessment. On the other hand, the statistical literature on missing data (e.g., Little & Rubin, 2002) clearly states that (i) there are an infinite number of missingness mechanisms that qualify as “missing not at random”, and (ii) if the mechanism in the model does not match the truth, then the parameter estimates may exhibit more bias than the corresponding “missing at random” estimates. This implies that the extra complexity of the proposed model may hurt us.

At least for the model proposed in this paper, there appears to be little danger in employing the MNAR model instead of the MAR model. This is because the MAR model is a special case of the MNAR model, being obtained by fixing a subset of the λ parameters to zero. Thus, if the MAR assumption is approximately fulfilled, the model should account for this automatically during estimation.

8.2. Relationship to scoring rules

The model described here could also be used to develop new types of model-based scoring rules (for related ideas, see Budescu & Bo, 2015). Existing scoring rules (such as the Brier score or logarithmic score; see, e.g., Gneiting & Raftery, 2007) work only on the forecasts themselves, requiring every forecaster to respond to exactly the same questions. Such is seldom the case in practice, and it is

awkward to tailor these scoring rules to missing data. For example, for each question that a forecaster fails to answer, we could substitute the mean observed Brier score on the corresponding question. This substitution is somewhat similar to IRT procedures that code unanswered questions as incorrect (though, in a forecasting context, the notion of “incorrect” is unclear).

Beyond substitution for missing observations, we can consider new scoring rules in which question selections and missing data play a role. A rough definition, corresponding to the models estimated in this paper, is as follows. A forecaster will receive the highest expected score if he/she:

- consistently makes better forecasts than the crowd, and
- responds to more questions.

This definition requires the best forecasters to be the best on both question selection and forecast reporting. Furthermore, middling forecasters could receive the same ability estimates through different routes. For example, say that Forecasters A and B receive the same ability estimates from the model. It would be possible for Forecaster A to obtain this estimate by selecting many questions but providing relatively bad forecasts on those questions, while Forecaster B obtains this estimate by selecting few questions but providing relatively good forecasts on those questions. Further work could examine the extent to which these two criteria simultaneously provide incentives for honest forecasting and frequent responding. A game-theoretic framework similar to that of [Prelec \(2004\)](#) might be useful here, because we can depict each forecaster as striving to do the minimal amount of forecasting required to be the best. Under these conditions, forecasters might be motivated to respond to all of the questions when they do not know other forecasters’ response patterns.

8.3. Model assumptions

As has been mentioned throughout, the model proposed here assumes a single dimension of forecaster ability; that is, each forecaster’s ability is summarized via a single number. While the analyses in this paper suggest that this assumption is not grossly violated in our dataset, there remains the possibility that it is grossly violated on other datasets. As an extreme example, we could imagine a forecaster who follows only local occurrences and knows nothing about broader world events. If this forecaster responds only to questions related to her locale, she may receive a good ability estimate in spite of the fact that her forecasts on the other, unanswered questions would be awful.

Despite this violation, the model’s handling of this extreme forecaster could still be reasonable. First, if other high-ability forecasters tend to respond to questions that do not involve this particular locale, then the model will temper the extreme forecaster’s ability estimate so that it is not as high as others. Second, if the extreme forecaster

does not respond to many questions (i.e., there are few questions about the forecaster’s locale), then the model will again temper her ability estimate: the model requires large amounts of data from the forecaster before it is “willing” to assign an extremely good ability estimate. While these results do not guarantee that the model will be robust to all dimensionality violations, they seem applicable to many situations in which evaluators wish to rank order forecasters across all questions.

On a related note, the model also assumes that each forecaster has a static level of forecasting ability and a static response propensity. In contrast, forecasters tend to change over time, gaining (losing) interest in the forecasting tournament and reporting improved (diminished) judgments. The model as proposed here cannot accommodate forecaster attributes that change over time, though it may be possible to directly model changes in forecaster ability over time via new parameters and/or increased dimensions of forecaster ability. It would also be of interest to relax the distributional assumptions by employing, say, t distributions instead of normal distributions or mixture models that accommodate subclasses of homogeneous forecasters. As is described further in the next section, traditional psychometric modeling frameworks can be helpful for including these model extensions.

8.4. Model estimation

The estimation of traditional item response models with multiple ability dimensions is generally difficult (e.g., [Cai, 2010](#)), and the same is true for the two-dimensional model proposed here. The Bayesian approach that we adopted introduces an additional complication, in that we must employ Markov chain Monte Carlo, sampling the forecaster ability parameters instead of integrating them out (e.g., [Lee, 2007](#)). This means that we must be careful to ensure that the model parameters are identified and that the model converges (e.g., [Ghosh & Dunson, 2009](#); [Merkle & Wang, in press](#); [Peeters, 2012](#)), which may introduce an undesirable practical complication.

Depending on the data, some simplifications are possible. In particular, we adopted the Bayesian approach in this paper so that we could easily include the “time of reported forecast” covariate in the model easily. However, this covariate is not necessary when all forecasters report their judgments at approximately the same time. If this covariate is not necessary, then the model proposed here could be estimated via maximum likelihood, using popular SEM software such as Mplus ([Muthén & Muthén, 1998–2012](#)) or lavaan ([Rosseel, 2012](#)). These approaches would make use of ideas related to the path diagrams from [Fig. 1](#). However, when the data are very sparse (i.e., each forecaster reports on a small proportion of questions), these programs may fail in situations where the Bayesian approach would succeed. This failure is again related to the fact that ML estimation methods integrate the forecaster latent variables out of the likelihood, whereas Bayesian estimation methods sample the forecaster latent variables directly (and are “smoothed” by the prior distributions). The integration of the latent variables requires us to work with the covariance matrix of a multivariate normal likelihood, which can

often become non-positive definite during model estimation (resulting in a failed estimation).

Sample size is an additional consideration for all of the models discussed here. Because the proposed model is related to traditional psychometric models (including factor analysis and item response models), we can draw on the psychometric literature for sample size recommendations. In that literature, it is customary to observe hundreds or thousands of participants reporting on small numbers of items. Other researchers proposing models similar to ours have tended to follow this trend: [Holman and Glas \(2005\)](#) applied their model to 171 participants responding to 32 items, whereas [O'Muirheartaigh and Moustaki \(1999\)](#) applied their model to two datasets, of which one had 2691 participants responding to five items and the other had 1270 participants responding to four items. While our application had many more items than the others, we generally recommend large numbers of participants and suggest artificial data simulation as a way of determining whether one's particular sample size is appropriate for estimating the parameters of interest. The Bayesian approach of sampling forecaster latent variables directly can be helpful again here, allowing us to bypass non-positive definite covariance matrices.

8.5. Summary

In situations where the respondents are free to select their own questions or stimuli, the selections themselves can provide valuable information about the latent respondent attributes that we wish to measure. While these selections are often viewed as nuisance characteristics of the data that cause difficulties for analysis, this paper has illustrated a model-based approach for capturing the information inherent in the selections. The ability to incorporate multiple types of variables (forecasts, question selections) in forecaster assessment is a major advantage of model-based approaches over data-based metrics (i.e., scoring rules), which rely exclusively on the reported forecasts. In forecasting scenarios and beyond, a detailed consideration of selection/missingness mechanisms could lead to improved estimation of the latent traits of interest.

Computational details

All results were obtained using the R system for statistical computing ([R Core Team, 2016](#)) version 3.3.3 and the JAGS software for Bayesian computation ([Plummer, 2003](#)) version 4.2.0, employing the add-on package `runjags` 2.0.4-2 ([Denwood, 2016](#)). R and the package `runjags` are available freely from the Comprehensive R Archive Network at <http://CRAN.R-project.org/> under the General Public License 2. JAGS is available freely from Sourceforge at <http://mcmc-jags.sourceforge.net/> under the General Public License 2.

Acknowledgments

For data access, please contact Phillip Rescober at rescober@goodjudgment.com. This research was supported by the Intelligence Advanced Research Projects Activity (IARPA) via the Department of Interior National Business Center (DoI/NBC) Contract No. D11PC20061. The US government is authorized to reproduce and distribute reprints for government purposes notwithstanding any copyright annotation thereon. The views and conclusions expressed herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of IARPA, DoI/NBC, or the US government.

Appendix A. DETECT technical details

We computed the DETECT statistic separately for the reported forecasts y and for the question selections d . The statistics were calculated via the `expl.detect()` function in the R package `sirt` ([Robitzsch, 2016](#)).

The calculation of the DETECT statistic for question selection was straightforward, because all forecasters had complete data corresponding to standard item response data. That is, each forecaster's data consisted of a series of zeros and ones, with a zero indicating that he/she did not respond to a particular question and a one indicating the opposite. In addition to the observed data, the DETECT statistic also requires unidimensional estimates of a person's ability. For this, we used the weighted likelihood estimates arising from a Rasch model.

We computed DETECT statistics for the reported forecasts by first restricting ourselves to a subset of 241 forecasters who responded to at least 136 of the 176 questions (with most of the forecasters responding to at least 160 of the questions). We did this so that we could ignore missing data mechanisms when examining forecast dimensionality. Next, we attempted to maximize the DETECT statistic by transforming the data to account for the fact that the forecasts were reported at different points in time. In particular, for each question j , we regressed the y^* s associated with question j on the time at which the forecast was reported (i.e., the t_{ij}). We then used the fitted model to push each person's reported forecast to the question's "halfway" point (i.e., the time at which the question is halfway between introduction and resolution). Finally, we computed the statistic by creating binary variables from the aligned forecasts (equal to 0 if the forecast was less than 0.5, 1 otherwise). We used the average forecast reported for each question's realized outcome as the estimate of a person's ability.

Appendix B. JAGS model estimation

JAGS code for estimating the model is displayed below. The probit-transformed forecasts y_{star} are given in long format, while the missingness indicators d are in a data matrix where the rows are forecasters and the columns are questions. Following the JAGS code, we provide R code to illustrate its usage.

```

model{
  for (i in 1:nr){ ## Rows of forecast data
    ystar[i] ~ dnorm(mu[i], invsig2[qidx[i]])

    mu[i] <- b0[qidx[i]] + b1[qidx[i]]*nd[i] +
      lambda[qidx[i], 1] *
      theta[pidx[i], 1]
  }

  for (i in 1:n){ ## Forecasters
    for (j in 1:J){ ## Questions
      d[i, j] ~ dbern(pd[i, j])

      probit(pd[i, j]) <- b0[(J + j)] +
        lambda[(J + j), 1]*theta[i, 1] +
        lambda[(J + j), 2]*theta[i, 2]
    }

    ## Person parameters
    theta[i, 1] ~ dnorm(0, invpsi[1])
    theta[i, 2] ~ dnorm(0, invpsi[2])
  }
  invpsi[1] ~ dgamma(0.01, 0.01)
  invpsi[2] ~ dgamma(0.01, 0.01)

  ## Equality constraints + priors for question parameters
  lambda[1,1] <- 1
  lambda[1,2] <- 0
  lambda[(J+1), 1] ~ dnorm(0, 1)
  lambda[(J+1), 2] <- 1

  b0[1] ~ dnorm(0, 0.5)
  b0[(J + 1)] ~ dnorm(0, 0.5)
  b1[1] ~ dnorm(0, 0.5)
  invsig2[1] ~ dgamma(0.01, 0.01)

  for (j in 2:J){
    ## loadings for forecasts
    lambda[j, 1] ~ dnorm(0, 1)
    lambda[j, 2] <- 0

    ## loadings for d parameters
    lambda[(J + j), 1] ~ dnorm(0, 1)
    lambda[(J + j), 2] ~ dnorm(0, 1)

    ## Intercept priors
    b0[j] ~ dnorm(0, 0.5)
    b0[(J + j)] ~ dnorm(0, 0.5)
    b1[j] ~ dnorm(0, 0.5)

    ## Error precision prior
    invsig2[j] ~ dgamma(0.01, 0.01)
  }
}

```

The R code below gives an example with artificial data, showing how the JAGS code can be run from within R using the runjags package.

```

library("runjags")
set.seed(1080)

## Generate data
n <- 500
K <- 100

## Probability judgments
b0 <- runif(K, -1, 2)
b1 <- runif(K, 0, 3)

lambda <- runif(K, -0.5, 3.5)
theta1 <- rnorm(n, 0, 1)
nd <- runif(n*K, -0.5, 0)

dat <- expand.grid(uidx=1:n, ifpidx=1:K)
dat$nd <- nd
mny <- b0[dat$ifpidx] + b1[dat$ifpidx]*nd +
  lambda[dat$ifpidx]*theta1[dat$uidx]
dat$ystar <- rnorm(n*K, mny, 0.4)
dat$ystar[dat$ystar < -3.5] <- -3.5
dat$ystar[dat$ystar > 3.5] <- 3.5
dat$ystar[dat$ystar > -0.1 & dat$ystar < 0.1] <- 0

dat$fcast1 <- pnorm(dat$ystar)

## Missingness indicators
b0 <- runif(K, 0.5, 2)
lambda <- matrix(runif(K*2, -0.5, 2.5), K, 2)
theta2 <- rnorm(n, 0, 1)

ppd <- lambda[,1] %>% matrix(theta1, 1, n) +
  lambda[,2] %>% matrix(theta2, 1, n)
ppd <- apply(ppd, 2, function(x) x + b0)
d <- apply(ppd, 2, function(x)
  rbinom(length(x), 1, pnorm(x)))
for(i in 1:K){
  subs <- which(d[i,] == 0)
  dat$ystar[dat$ifpidx == i & dat$uidx %in% subs] <- NA
}

rmrows <- which(is.na(dat$ystar))
dat <- dat[~rmrows,]

## Data formatted for JAGS
data <- list(nr = nrow(dat), n = length(unique(dat$uidx)),
  J = length(unique(dat$ifpidx)), ystar = dat$ystar,
  qidx = dat$ifpidx, pidx = dat$uidx, nd = dat$nd,
  d = t(d))

## Starting values
inits <- list(b0 = rep(0, 2*data$J), invpsi = rep(1, 2),
  theta = matrix(0, data$n, 2),
  b1 = rep(0.1, data$J),
  invsig2 = rep(1, data$J))

## MCMC run, will take some time
runjags.options(force.summary = TRUE)
mdraws <- run.jags("paper_model.jag", data=data,
  inits=inits, monitor=c("theta", "b0", "b1", "lambda"),
  n.chains=3, burnin=5000, sample=1000)

## Parameter summaries, posterior means
mdraws$summaries
mdraws$summaries[, "Mean"]

```

References

- Bejar, I. (1977). An application of the continuous response level model to personality measurement. *Applied Psychological Measurement, 1*, 509–521.
- Bonifay, W. E., Reise, S. P., Scheines, R., & Meijer, R. R. (2015). When are multidimensional data unidimensional enough for structural equation modeling? An evaluation of the DETECT multidimensionality index. *Structural Equation Modeling, 22*, 504–516.
- Borsboom, D., Mellenbergh, G. J., & van Heerden, J. (2004). The concept of validity. *Psychological Review, 111*, 1061–1071.

- Budescu, D. V., & Bo, Y. (2015). Analyzing test-taking behavior: Decision theory meets psychometric theory. *Psychometrika*, 80, 1105–1122.
- Cai, L. (2010). High-dimensional exploratory item factor analysis by a Metropolis–Hastings Robbins–Monro algorithm. *Psychometrika*, 75, 33–57.
- Carvalho, A. (2016). An overview of applications of proper scoring rules. *Decision Analysis*, 13, 223–242.
- Chang, Y. W., Tsai, R. C., & Hsu, N. J. (2014). A speeded item response model: Leave the harder till later. *Psychometrika*, 79, 255–274.
- Denwood, M. J. (2016). runjags: An R package providing interface utilities, model templates, parallel computing methods and additional distributions for MCMC models in JAGS. *Journal of Statistical Software*, 71(9), 1–25. Retrieved from <http://www.jstatsoft.org/v71/i09/>.
- Embretson, S. E., & Reise, S. P. (2000). *Item response theory for psychologists*. Mahwah, NJ: Erlbaum Associates.
- Ferrando, P. J. (2001). A nonlinear congenic model for continuous item responses. *The British Journal of Mathematical and Statistical Psychology*, 54, 293–313.
- Gelman, A., & Rubin, D. B. (1992). Inference from iterative simulation using multiple sequences (with discussion). *Statistical Science*, 7, 457–511.
- Ghosh, J., & Dunson, D. B. (2009). Default prior distributions and efficient posterior computation in Bayesian factor analysis. *Journal of Computational and Graphical Statistics*, 18, 306–320.
- Gneiting, T., & Raftery, A. E. (2007). Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102, 359–378.
- Holman, R., & Glas, C. A. W. (2005). Modelling non-ignorable missing-data mechanisms with item response theory models. *The British Journal of Mathematical and Statistical Psychology*, 58, 1–17.
- Lee, S. Y. (2007). *Structural equation modeling: A Bayesian approach*. Chichester: Wiley.
- Little, R. J. A., & Rubin, D. B. (2002). *Statistical analysis with missing data* (2nd ed.). New York: Wiley.
- Lord, F. M., & Novick, M. R. (1968). *Statistical theories of mental test scores*. Reading, MA: Addison-Wesley.
- McDonald, R. P. (1999). *Test theory: A unified treatment*. Mahwah, NJ: Erlbaum.
- Mellers, B., Stone, E., Atanasov, P., Rohrbaugh, N., Metz, S. E., Ungar, L., et al. (2015). The psychology of intelligence analysis: Drivers of prediction accuracy in world politics. *Journal of Experimental Psychology: Applied*, 21, 1–14.
- Mellers, B., Stone, E., Murray, T., Minster, A., Rohrbaugh, N., Bishop, M., et al. (2015). Identifying and cultivating superforecasters as a method of improving probabilistic predictions. *Perspectives on Psychological Science*, 10, 267–281.
- Mellers, B., Ungar, L., Baron, J., Ramos, J., Gurcay, B., Fincher, K., et al. (2014). Psychological strategies for winning a geopolitical forecasting tournament. *Psychological Science*.
- Merkle, E. C., Steyvers, M., Mellers, B., & Tetlock, P. E. (2016). Item response models of probability judgments: Application to a geopolitical forecasting tournament. *Decision*, 3, 1–19.
- Merkle, E. C., & Wang, T. (2017). Bayesian latent variable models for the analysis of experimental psychology data. *Psychonomic Bulletin and Review*, (in press).
- Müller, H. (1987). A Rasch model for continuous ratings. *Psychometrika*, 52, 165–181.
- Muthén, B. (1989). Tobit factor analysis. *The British Journal of Mathematical and Statistical Psychology*, 42, 241–250.
- Muthén, L. K., & Muthén, B. (1998–2012). *Mplus user's guide* (7th ed.). Los Angeles, CA: Muthén & Muthén.
- Noel, Y., & Davier, B. (2007). A beta item response model for continuous bounded responses. *Applied Psychological Measurement*, 31, 47–73.
- O'Muirheartaigh, C., & Moustaki, I. (1999). Symmetric pattern models: A latent variable approach to item non-response in attitude scales. *Journal of the Royal Statistical Society. Series A. Statistics in Society*, 162, 177–194.
- Peeters, C. F. W. (2012). Rotational uniqueness conditions under oblique factor correlation metric. *Psychometrika*, 77, 288–292.
- Plummer, M. (2003). JAGS: A program for analysis of Bayesian graphical models using Gibbs sampling. In K. Hornik, F. Leisch, & A. Zeileis (Eds.), Proceedings of the 3rd international workshop on distributed statistical computing.
- Prelec, D. (2004). A Bayesian truth serum for subjective data. *Science*, 306, 462–466.
- R Core Team (2016). R: A language and environment for statistical computing [Computer Software Manual]. Vienna, Austria. Retrieved from <https://www.R-project.org/>.
- Reise, S. P., Scheines, R., Widaman, K. F., & Haviland, M. G. (2013). Multi-dimensionality and structural coefficient bias in structural equation modeling: A bifactor perspective. *Educational and Psychological Measurement*, 73, 5–26.
- Robitzsch, A. (2016). sirt: Supplementary item response theory models [Computer software manual]. Retrieved from <https://CRAN.R-project.org/package=sirt> (R package version 1.12-2).
- Rose, N., von Davier, M., & Xu, X. (2010). *Modeling nonignorable missing data with item response theory (IRT)* (Tech. Rep.). ETS Research Report.
- Rossee, Y. (2012). lavaan: An R package for structural equation modeling. *Journal of Statistical Software*, 48(2), 1–36. Retrieved from <http://www.jstatsoft.org/v48/i02/>.
- Roussos, L. A., & Ozbeck, O. Y. (2006). Formulation of the DETECT population parameter and evaluation of DETECT estimator bias. *Journal of Educational Measurement*, 43, 215–243.
- Rubin, T.N., & Steyvers, M. (2009). A topic model for movie choices and ratings. In *Proceedings of the 9th international conference on cognitive modeling – ICCM2009*. Manchester, UK.
- Samejima, F. (1973). Homogeneous case of the continuous response model. *Psychometrika*, 38, 203–219.
- Stout, W., Habing, B., Douglas, J., Kim, H. R., Roussos, L. A., & Zhang, J. (1996). Conditional covariance-based nonparametric multidimensionality assessment. *Applied Psychological Measurement*, 20, 331–354.
- Thissen, D. (2016). Bad questions: An essay involving item response theory. *Journal of Educational and Behavioral Statistics*, 41, 81–89.
- van Abswoude, A. A. H., van der Ark, L. A., & Sijtsma, K. (2004). A comparative study of test dimensionality assessment procedures under nonparametric IRT models. *Applied Psychological Measurement*, 28, 3–24.
- Wang, W. C., Jin, K. Y., Qiu, X. L., & Wang, L. (2012). Item response models for examinee-selected items. *Journal of Educational Measurement*, 49, 419–445.
- Zhang, J. (2007). Conditional covariance theory and DETECT for polytomous items. *Psychometrika*, 72, 69–91.
- Zhang, J., & Stout, W. (1999). The theoretical DETECT index of dimensionality and its application to approximate simple structure. *Psychometrika*, 64, 213–249.

Edgar C. Merkle is associate professor in the Department of Psychological Sciences at the University of Missouri. He received a Ph.D. in quantitative psychology and an MS in statistics, both from The Ohio State University. His research interests include latent variable models, subjective probability and forecasts, and statistical computing. He has authored numerous journal articles within these areas.

Mark Steyvers is a Professor of Cognitive Science at U.C. Irvine and holds joint appointments with the Departments of Computer Science and of Psychology & Social Behavior. His research interests include collective intelligence, human memory and decision-making as well as statistical machine learning and information retrieval. He is a past president of the Society of Mathematical Psychology, a recipient of the Early Investigator Award from The Society for Experimental Psychologists and a new investigator award from the American Psychological Association.

Barbara A. Mellers is the George I. Heyman University Professor at the University of Pennsylvania with appointments in the Department of Psychology and the Wharton Marketing Department. She serves as Associate and Consulting Editors on numerous academic journals and has received several grants from NSF and IARPA. Her research focuses on how people make judgments and decisions and how they can do it better. She is currently investigating how people can improve their forecasts and whether forecasting tournaments are a mechanism for increasing open mindedness.

Philip E. Tetlock holds the Annenberg University Professor chair at the University of Pennsylvania, with cross-appointments in Wharton and the School of Arts and Sciences. His work addresses a wide range of topics, including cognitive biases, accountability systems, value conflict, and taboo trade-offs. He has received awards from many scientific societies, including the American Psychological Association, the American Political Science Association, the National Academy of Sciences, and the American Association for the Advancement of Science. His most recent work focuses on forecasting tournaments and their potential to improve accuracy and to depolarize unnecessarily polarized debates. He is the author of *Superforecasting*, which tells the story of how the Good Judgment Project won a series of forecasting tournaments sponsored by the U.S. intelligence community.