

# Statistical analysis of the National Institutes of Health peer review system

Valen E. Johnson\*

University of Texas M.D. Anderson Cancer Center, 1400 Pressler Street, Unit #1411, Houston, TX 77030

Communicated by James O. Berger, Duke University, Durham, NC, May 15, 2008 (received for review February 17, 2008)

**A statistical model is proposed for the analysis of peer-review ratings of R01 grant applications submitted to the National Institutes of Health. Innovations of this model include parameters that reflect differences in reviewer scoring patterns, a mechanism to account for the transfer of information from an application's preliminary ratings and group discussion to final ratings provided by all panel members and posterior estimates of the uncertainty associated with proposal ratings. Application of this model to recent R01 rating data suggests that statistical adjustments to panel rating data would lead to a 25% change in the pool of funded proposals. Viewed more broadly, the methodology proposed in this article provides a general framework for the analysis of data collected interactively from expert panels through the use of the Delphi method and related procedures.**

hierarchical model | item response model | latent variable model | ordinal data

Every year, the National Institutes of Health (NIH) spend more than \$22 billion to fund scientific research (1). Approximately 70% of these funds are awarded through a peer-review process overseen by the NIH Center for Scientific Review (CSR). Despite the vast sum of money involved, the absence of statistical methodology appropriate for the analyses of peer-review scores generated by this system has precluded the type of detailed assessment applied to other national health and educational systems (2, 3). As a consequence, statistical adjustments to account for uncertainties and biases inherent to these scores are not made before funding decisions. To address this deficiency, this article examines the properties of these ratings and proposes methodology to more efficiently extract the information contained in them.

It is useful to begin with a brief review of the NIH peer-review system. Upon submission to the NIH, most grant applications (e.g., R01, R03, R21, etc.) are assigned to a study section within an Integrated Review Group (IRG) for review, and to an NIH Institute and Center (IC) for eventual funding. IRG study sections typically contain  $\approx 30$  members and review  $\approx 50$  grant applications (proposals) during each of three annual meetings. Because it is impractical for every member of a study section to review every application, between two and five reviewers are typically assigned to read and score each application before the study section convenes. In the sequel, these individuals are called the proposal's "readers," and the scores they assign before a study section convenes are called "pre-scores." Proposals are scored on a 1.0–5.0 scale in increments of 0.1 units, with 1.0 representing the best score. When the study section convenes, the scientific review officer (SRO) and the study section chair suggest a list of proposals that might be "streamlined." Based on their pre-scores, proposals on this list are viewed as unlikely to receive fundable priority scores and, if no one in the study section objects, are not considered further. The remaining proposals are discussed and scored by all members of the study section.

Readers of a grant application begin the discussion by announcing their pre-scores and summarizing the proposal for other members of the study section, most of whom will not have read it. After these summaries, there is an open discussion of the

application. Proposal readers then state their "post-scores" for the application, and all other members of the study section (i.e., the proposal's nonreaders) also score the proposal. Nonreaders are required to either score the proposal within 0.5 units of the range of scores established by reader post-scores or provide a written statement to the SRO explaining why they scored the proposal outside of that range. Scores received from all study section members are then averaged to obtain the proposal's priority score. In "established" study sections, priority scores are converted to a percentile ranking through a comparison with recent priority scores from other grant applications scored within that study section. In newer study sections or special emphasis panels (i.e., panels that are convened to rate a limited number of proposals), percentile scores are calculated by comparing the proposal's priority score to established norms. Finally, proposal percentile ratings are used by ICs to determine which applications will be funded. Although the exact criteria by which ICs use these percentiles to make funding decisions vary by the IC, funding decisions are thought to be highly correlated with percentile scores.

In this article, I propose statistical methodology to account for the effect of the selection of readers on a proposal's final percentile score, quantify the uncertainty associated with the percentile scores, and demonstrate how such uncertainties can be incorporated into a decision-theoretic framework to improve the probability that the greatest proportion of top proposals are funded. Viewed more generally, methods developed in this article extend existing statistical methodology for the analysis of multirater ordinal data (4–7) and item response data (8–12) to provide a framework for the analysis of panel rating data collected by using the Delphi method and related interactive rating schemes (13).

The data that form the basis for this study were collected as part of a contract awarded to the author by the CSR in November 2004. As part of that study, all preliminary and final reader scores and nonreader scores for all R01 grant proposals submitted to the NIH and reviewed under the auspices of the CSR over two review cycles (June and October 2005) were collected and redacted.

## Description of Data

Ratings for 18,959 R01 proposals rated by 14,041 reviewers in 744 study sections (including special emphasis panels) were available for analysis. Fig. 1 displays a histogram of all scores, including reader pre-scores and post-scores, and nonreader

Author contributions: V.E.J. designed research, performed research, analyzed data, and wrote the paper.

The author declares no conflict of interest.

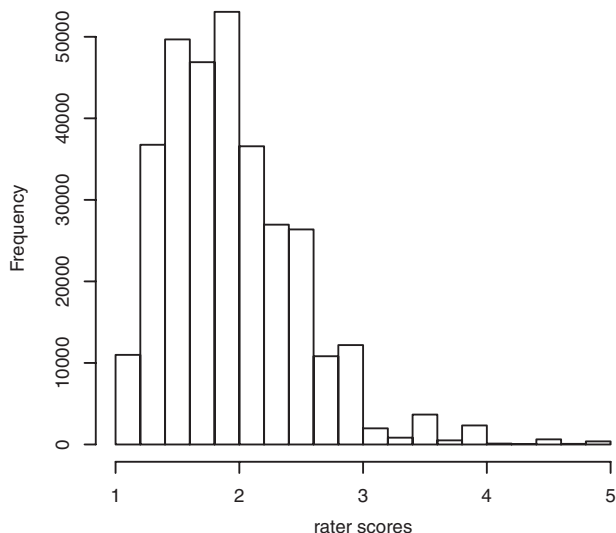
Freely available online through the PNAS open access option.

Data deposition: Dr. Johnson will provide the data in ASCII format upon request.

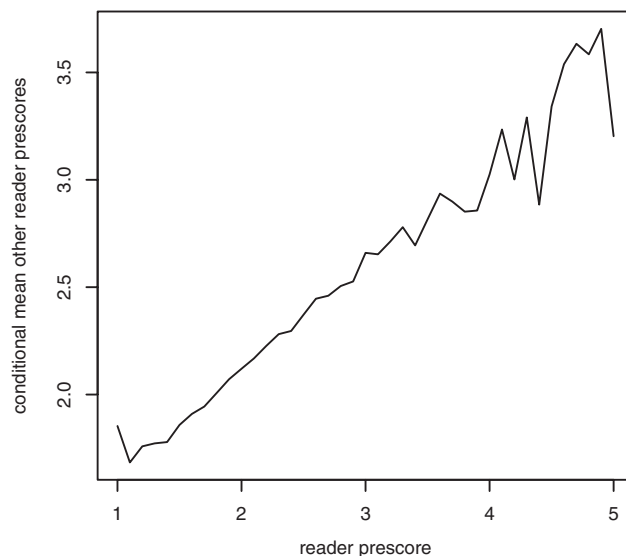
\*E-mail: vejohanson@mdanderson.org.

This article contains supporting information online at [www.pnas.org/cgi/content/full/0804538105/DCSupplemental](http://www.pnas.org/cgi/content/full/0804538105/DCSupplemental).

© 2008 by The National Academy of Sciences of the USA



**Fig. 1.** Histogram of rater scores (including reader pre-scores and post-scores, and nonreader scores) assigned to R01 proposals.



**Fig. 2.** Plot of the conditional mean pre-scores assigned by other readers versus single reader pre-scores.

scores. Table 1 provides a summary of the mean and standard deviation of the rater scores.

Several interesting features of the data are apparent from Fig. 1. Among these is a tendency for reviewers to use two distinct scales to score proposals. The first scale, nominally assumed by the CSR, runs from 1.0 to 5.0 in increments of 0.1 units. The second scale, used more frequently for less competitive proposals, runs from 1.0 to 5.0 in increments of 0.5 units. Evidence for the operation of these dual scales is provided in Fig. 2, in which the conditional means of reader pre-scores are displayed as a function of the prescore assigned to a proposal by a single reviewer. The relation between a reader prescore and the mean of other reader pre-scores for the same proposal is nearly linear between  $\approx 1.1$  and 3.0, but, outside of that range, the relationship is not monotonic. For example, among proposals that receive one prescore of 5.0, the mean of the remaining pre-scores is 3.2; for proposals receiving a prescore of 4.9, the mean of the remaining pre-scores is 3.7. Although not a central focus of this article, these observations suggest that a 20-point scale, anchored at an “average” rating of 10, might be better supported by current rating procedures. Such a scale would nominally provide a 10-point scale for nonstreamlined proposals.

## Results

I used a latent variable model (14, 15) to formally describe the relation among application merit, reader pre- and post-scores, and nonreader scores. Within this model, reader pre-scores were assumed to represent independent assessments of application merit, whereas reader post-scores and nonreader scores were assumed to represent weighted averages of information elicited during the proposal discussion and the scores of (other) proposal readers. I used a continuous-valued latent variable  $\mu_i$  to represent the merit of the  $i$ th application. The resulting model was

**Table 1.** Summary statistics for R01 proposal rater scores

	pre-scores (all)	pre-scores (not streamlined)	post-scores	non-reader scores
sample mean	2.21	1.88	1.90	1.96
std. deviation	0.77	0.51	0.49	0.50

Columns provide the mean and standard deviations of reader pre-scores for all proposals, reader pre-scores for proposals that were not streamlined, reader post-scores, and nonreader scores.

then used to estimate the effects of reader biases and to assess the uncertainty in final proposal rankings. A description of this statistical model is provided in the [supporting information \(SI\)](#).

**Adjustments for Reader Bias.** Demonstrating the benefit of corrections for reviewer bias is difficult because true proposal merits are not known. For this reason, I examined the effectiveness of bias corrections in two stages. First, I performed a cross-validation study that used only reader pre-scores. Because reader pre-scores can be considered to be conditionally independent, they can be analyzed without modeling the complex structure among their values, reader post-scores, and nonreader scores. Therefore, a comparison of the model-based prediction errors based on reader pre-scores to the NIH prediction error provides an indication of the effectiveness of corrections for reader biases and a partial model validation. Second, I applied the full statistical model to all rater scores to illustrate the impact of reader bias on the final estimates of the proposals' merits.

I implemented the cross-validation experiment by first splitting reader pre-scores into two samples, randomly assigning 90% of the scores to a training sample and assigning the remaining 10% to a test sample. I used the training data to estimate model parameters. The posterior means of merit parameters for the proposals were then converted back to the original rating scale and were used to predict pre-scores in the test sample. The mean squared error for these predictions was 0.373.

In the NIH scoring system, proposal merit is estimated by the sample mean of the raters' scores. Thus, the estimate of a proposal's merit based on the training sample is the sample mean of the training sample pre-scores. The mean squared error of the corresponding prediction of pre-scores in the test sample was 0.413. Use of the statistical model to predict reader pre-scores in the test sample thus reduced the mean squared error of prediction by  $\approx 10\%$ .

The improvement in mean squared error enjoyed by the model-based estimate can be attributed primarily to the estimation of parameters that represent rater biases, or the tendency of some raters to score proposals more stringently than others.

When propagated through the full statistical model for reader post-scores and nonreader scores, these effects can be quite dramatic. For example, consider the posterior estimates of the proposal merits listed in Table 2. These proposals represent the







The mean of rater biases  $\zeta$  was included in the model to account for the fact that reader pre-scores have a lower mean value than either the rater post-scores or nonreader post-scores.

**Second-Stage Model.** In the next stage of the data generating process, I assumed that readers modified their pre-scores by using both the reported values of other reader pre-scores and the group's discussion of the proposal. The resulting reader post-scores were thereby represented as a weighted average of these three information sources.

The latent value  $x_{i,j}^{\text{post}}$  assumed to be responsible for the generation of reader  $j$ 's postscore of proposal  $i$ ,  $y_{ij}^{\text{post}}$ , can be written as

$$x_{i,j}^{\text{post}} = u_{i,j} x_{i,j}^{\text{pre}} + v_i \mu_i + \sum_{k \in A_i; k \neq i} w_{i,j,k} x_{i,j,k}^{\text{pre}} + \varepsilon_{ij}^{\text{post}}, \quad [2]$$

where

$$y_{i,j}^{\text{post}} = s \Leftrightarrow \gamma_{m,s-1} \leq x_{i,j}^{\text{post}} < \gamma_{m,s}, \quad [3]$$

and

$$u_{i,j} + v_i + \sum_{k \in A_i; k \neq i} w_{i,j,k} = 1. \quad [4]$$

The error terms  $\varepsilon_{ij}^{\text{post}}$  were assumed to be independently distributed as  $N(0, \sigma_1^2)$  random variables. Here,  $A_i$  denotes the set of reviewers who provided pre-scores for proposal  $i$ .

On the latent scale of measurement, the model specification described so far resembles a standard hierarchical model with a Gaussian error structure. Unfortunately, the usual Gaussian model does not provide an accurate representation of reader post-scores and nonreader scores at higher levels in the model hierarchy. This difficulty stems from the high proportion of reader post-scores that fall within the range defined by the reader pre-scores, and the even higher proportion of nonreader scores that fall within the range defined by the reader post-scores. There also is a tendency for nonreaders to assign scores that are identical to a reader postscore.

To account for these tendencies, the weights  $u_{ij}$ ,  $v_i$ , and  $w_{ijk}$  were assumed to be generated from a Dirichlet model with a parameter vector containing a component  $a$  for each  $u_{ij}$ , a component  $b$  for each  $v_i$ , and a component  $c$  for

each  $w_{ijk}$ . The distribution of hyperparameters estimated at higher levels in the model hierarchy make it likely that these weights are assigned values that are either close to 0 or 1; this permits the model to mimic the tendency of nonreaders to concentrate their scores around and between the scores recorded by the proposal's readers.

Another innovation of the statistical model involves the inclusion of the term  $v_i \mu_i$  in the weighted average defining the latent variable  $x_{i,j}^{\text{post}}$  (Eq. 2). The purpose of this term is to model systematic shifts between reader pre-scores and reader post-scores that result from a proposal's discussion. In the construction of this term,  $v_i$  weights  $\mu_i$ , the parameter that represents the true merit of the proposal. That is, the model implicitly assumes what might be regarded as the ideal situation from the NIH's standpoint. Alternative assumptions regarding the distributions of these weights can be incorporated into the model framework, but for the purposes of this article the NIH's "ideal" was assumed. It is important to note, however, that the rating data themselves cannot be used to validate this assumption in the absence of an external "gold standard" for relative proposal merits.

The values of the hyperparameters  $a$ ,  $b$ , and  $c$  determine, respectively, the average relative weights that readers assign to their own pre-scores, the proposal discussion, and the pre-scores of other readers when determining their final postscore ratings.

**Third-Stage Model.** The model for nonreader scores  $y_{i,j}^{\text{non}}$  is similar to the model specified for reader post-scores  $y_{i,j}^{\text{post}}$ , except that nonreader scores were assumed to be based on a latent variable  $x_{i,j}^{\text{non}}$  that represents a weighted average of reader post-scores and proposal merit. That is, the model for nonreader scores was obtained by replacing Eq. 2 with

$$x_{i,j}^{\text{non}} = v_i \mu_i + \sum_{k \in B_i} w_{i,j,k} x_{i,j,k}^{\text{post}} + \varepsilon_{i,j}^{\text{non}}, \quad [5]$$

and modifying Eqs. 3 and 4 accordingly. The weights appearing in Eq. 5 were defined similarly to those used to model reader post-scores.

Further description of higher-level model structures [including the prior distributions imposed on model hyperparameters ( $\gamma_m$ ,  $a$ ,  $b$ ,  $c$ ,  $\sigma_0^2$ ,  $\sigma_1^2$ ,  $\sigma_2^2$ ,  $\tau^2$ )], along with model diagnostics and a brief description of the numerical algorithm used to fit this model to the peer-review data, is provided in the [SI](#).

**ACKNOWLEDGMENTS.** I thank James Berger and two referees for constructive comments and suggestions that significantly improved the manuscript.

1. Office of Budget, National Institutes of Health (2007) Summary of the FY 2008 President's Budget. Available at <http://officeofbudget.od.nih.gov/PDF/Press%20info-2008.pdf>.
2. Goldstein H, Spiegelhalter DJ (1996) League tables and their limitations: Statistical issues in comparisons of institutional performance. *J Roy Stat Soc* 159:385–443.
3. Bird SM, et al. (2005) Performance indicators: Good, bad, and ugly. *J Roy Stat Soc* 168:1–27.
4. Albert JA, Chib S (1993) Bayesian analysis of binary and polychotomous response data. *J Am Stat Assoc* 88:669–679.
5. Johnson VE (1996) On Bayesian analysis of multirater ordinal data: An application to automated essay grading. *J Am Stat Assoc* 91:42–51.
6. Johnson VE, Albert JA (1999) *Ordinal Data Modeling* (Springer, New York).
7. Ishwaran H (2000) Univariate and multirater ordinal cumulative link regression with covariate specific cutpoints. *Can J Stat* 28:715–730.
8. Verhelst N, Verstralen H (2001) IRT models for multiple raters. *Essays on Item Response Theory*, eds Boosma A, van Duijn M, Snijders T (Springer, New York), pp 89–108.
9. Wilson M, Hoskens M (2001) The rater bundle. *J Educ Behav Stat* 26:283–306.
10. Patz RJ, Junker BW, Johnson MS, Mariano LT (2002) The hierarchical rater model for rated test items and its application to large-scale educational assessment data. *J Educ Behav Stat* 27:341–384.
11. Skrondal A, Rabe-Hesketh S (2004) *Generalized Latent Variable Modeling: Multilevel, Longitudinal, and Structural Equation Models* (CRC, Boca-Raton, FL).
12. Mariano LT, Junker BW (2007) Covariates of the rating process in hierarchical models for multiple ratings of test items. *J Educ Behav Stat* 32:287–314.
13. Rowe G, Wright G (1999) The Delphi technique as a forecasting tool: Issues and analysis. *Int J Forecast* 15:353–375.
14. Bollen KA (2002) Latent variables in psychology and the social sciences. *Annu Rev Psychol* 53:605–634.
15. Borsboom D, Mellenbergh GJ, van Heerden J (2003) The theoretical status of latent variables. *Psychol Rev* 110:203–209.
16. National Research Council (2005) *Strengthening Peer Review in Federal Agencies That Support Education Research* (National Academies Press, Washington, DC).
17. Kolen, MJ, Brennan RL (1995), *Test Equating: Methods and Practices* (Springer, New York).
18. McCullagh P (1980). Regression models for ordinal data. *J Roy Stat Soc Ser B Method* 42:109–142.