1        Reliability of Instructor Evaluations

2        Erin M. Buchanan[1], Jacob Miranda[2], & Christian Stephens[2]

3        [1] Harrisburg University of Science and Technology

4        [2] University of Alabama

5        Author Note

Abstract

TBA

*Keywords:* keywords

Word count: X

<p>17</p>

<center>Reliability of Instructor Evaluations</center>

<p>18</p>

The following was pre-registered: https://osf.io/czb4f

<p>19</p>

Exploratory Research Questions:

<p>20</p>

1) What is the reliability of instructor evaluations?

<p>21</p>

2) Are instructor evaluations reliable across time?

<p>22</p>

3) Is the average level of perceived fairness of the grading in the course a moderator of

<p>23</p>

reliability in instructor evaluations?

<p>24</p>

4) Does the average variability in instructor fairness rating moderate reliability of

<p>25</p>

instructor evaluations?

<p>26</p>

## Method

<p>27</p>

### Data Source

<p>28</p>

The archival study was conducted using data from the psychology department at a

<p>29</p>

large Midwestern public university. We used data from 2898 undergraduate, 274

<p>30</p>

mixed-level undergraduate, and 42 graduate psychology classes taught from 1987 to 2018

<p>31</p>

that were evaluated by students using the same 15-item instrument. Faculty followed set

<p>32</p>

procedures in distributing scan forms no more than two weeks before the conclusion of the

<p>33</p>

semester. A student was assigned to collect the forms and deliver them to the

<p>34</p>

departmental secretary. The instructor was required to leave the room while students

<p>35</p>

completed the forms. In the last several years of evaluations, online versions of these forms

<p>36</p>

were used with faculty encouraged to give students time to complete them in class while

<p>37</p>

they were outside the classroom.

<p>38</p>

The questionnaire given to students can be found at https://osf.io/4sphx. These

<p>39</p>

items were presented with a five-point scale from 1 (*strongly disagree*) to 5 (*strongly agree*).

<p>40</p>

For this study, the overall instructor evaluation question was "The overall quality of this

41    course was among the top 20% of those I have taken.". For fairness, we used the question of

42    "The instructor used fair and appropriate methods in the determination of grades.". The

43    ratings were averaged for each course, and the sample size for each rating was included.

**Planned Analyses**

45    The evaluations were filtered for those with at least fifteen student ratings for the

46    course (Rantanen, 2012). We performed a robustness check for the first research question

47    on the data when the sample size is at least $n = 10$ up to $n = 14$ (i.e., on all evaluations

48    with at least 10 ratings, then at least 11 ratings, etc.) to determine if the reliability

49    estimates are stable at lower sample sizes. We first screened the dataset (two evaluation

50    questions, sample size for course) for accuracy errors, linearity, normality, and

51    homoscedasticity. The data is assumed to not have traditional "outliers", as these

52    evaluations represent true averages from student evaluations. If the linearity assumption

53    fails, we considered potential nonparametric models to address non-linearity. Deviations

54    from normality were noted as the large sample size should provide robustness for any

55    violations of normality. If data appears to be heteroscedastic, we used bootstrapping to

56    provide estimates and confidence intervals.

57    This data was considered structured by instructor; therefore, all analyses below were

58    coded in $R$ using the *nlme* package (Pinheiro, Bates, Debroy, Sarkar, & Team, 2017) to

59    control for correlated error of instructor as a random intercept in a multilevel model.

60    Multilevel models allow for analysis of repeated measures data without collapsing by

61    participant [i.e., each instructor/semester/course combination can be kept separate without

62    averaging over these measurements; Gelman (2006)]. Random intercept models are

63    regression models on repeated data that structure the data by a specified variable, which

64    was instructor in this analysis. Therefore, each instructor's average rating score was

65    allowed to vary within the analysis, as ratings would be expected to be different from

66    instructor to instructor. In each of the analyses described below, the number of students

67 providing ratings for the course was included as a control variable to even out differences in

68 course size as an influence in the results. However, this variable was excluded if the models

69 did not converge. The dependent variable and predictors varied based on the research

70 question, and these are described with each analysis below.

71 **RQ 1.**   In this research question, we examined the reliability of instructor

72 evaluations on the overall rating and separately on the fairness rating. We calculated eight

73 types of reliability using course (same or different) by instructor (same or different) by

74 semester (same or different). The dependent variable was the first question average with a

75 predictor of the comparison question average, and both sample sizes (first sample size,

76 comparison sample size). Instructor code was used as the random intercept for both ratings

77 (i.e., two instructor random intercepts, first and comparison). The value of interest was the

78 standardized regression coefficient for the fixed effect of question from this model. Given

79 that the large sample size will likely produce "significant" $p$-values, we used the 95% CI to

80 determine which reliability values were larger than zero and to compare reliability

81 estimates to each other.

82 **RQ 2.**   We used the reliability for the same instructor and course calculated as

83 described in RQ1 at each time point difference between semesters. For example, the same

84 semester would create a time difference of 0. The next semester (Spring to Summer,

85 Summer to Fall, Fall to Spring) would create a time difference of 1. We used the time

86 difference as a fixed effect to predict reliability for the overall question only with a random

87 intercept of instructor. We used the coefficient of time difference and its confidence interval

88 to determine if there was a linear change over time. Finally, we plotted the changes over

89 time to examine if this effect was non-linear in nature and discuss implications of the graph.

90 **RQ 3.**   Using the reliability estimates from RQ 2, we then added the average rating

91 for the fairness question as the moderator with time to predict reliability. Fairness was

92 calculated as the average of the fairness question for all courses involved in the reliability

93 calculation for that instructor and time difference. Therefore, this rating represented the

94 average perceived fairness of grading at the time of ratings. If this interaction effect's

95 coefficient does not include zero, we performed a simple slopes analysis to examine the

96 effects of instructors who were rated at average fairness, one standard deviation below

97 average, and one standard deviation above average (Cohen, Cohen, West, & Aiken, 2003).

98     **RQ 4.**    Finally, we examined the average standard deviation of fairness ratings as a

99 moderator of with time to predict reliability. This variable represented the variability in

100 perceived fairness in grading from student evaluations, where small numbers indicated

101 relative agreement on the rating of fairness and larger values indicated a wide range of

102 fairness ratings. The variability in fairness ratings was calculated in the same way as the

103 mean fairness, which was only for the instructor and semester time difference evaluations

104 that were used to calculate the reliability estimate. This research question was assessed the

105 same way as research question three.

# Results

## Data Screening

108     The overall dataset was screened for normality, linearity, homogeneity, and

109 homoscedasticity using procedures from Tabachnick, Fidell, and Ullman (2019). Data

110 generally met assumptions with a slight skew and some heterogeneity. The complete

111 anonymized dataset and other information can be found online at https://osf.io/k7zh2.

112 This page also includes the manuscript written inline with the statistical analysis with the

113 *papaja* package (Aust et al., 2022) for interested researchers/reviewers who wish to recreate

114 these analyses.

## Descriptive Statistics

116     3214 evaluations included at least 15 student evaluations for analysis. Table 1

117 portrays the descriptive statistics for each course level including the total number of

118 evaluations, unique instructors, unique course numbers, and average scores for the two

119 rating items. Students additionally projected their course grade for each class ($A = 5$, $B =$

120 $4$, $C = 3$, $D = 2$, $F = 1$), and the average for this item is included for reference. Overall,

121 231 unique instructors and 70 unique courses were included in the analyses below across 94

122 semesters.

**RQ 1**

124     Each individual evaluation was compared to every other evaluation resulting in

125 5163291 total comparisons. Eight combinations of ratings were examined using instructor

126 (same, different), course (same, different), and semester (same, different) on both the

127 overall and fairness evaluation ratings separately. One of the individual ratings was used to

128 predict the comparison rating (i.e., question 1 was used to predict a comparison question 1

129 for the same instructor, different instructor, same semester, different semester, etc.), and

130 the number of ratings (i.e., rating sample size) per question were used as fixed-effects

131 covariates. The instructor(s) were used as a random intercept to control for correlated

132 error and overall average rating per instructor. The effects were then standardized using

133 the *parameters* package (Lüdecke et al., 2023). The data was sorted by year and semester

134 such that "predictor" was always an earlier semester predicting a later semester's scores,

135 except in cases of the the same semester comparisons. Therefore, positive standardized

136 scores indicate that scores tend to go up over time, while negative scores indicate that

137 scores tend to go down over time.

138     As shown in 1, reliability was highest when calculated on the same instructor in the

139 same semester and within the same course for both overall rating and fairness. This

140 reliability was followed by the same instructor, same semester, and different courses. Next,

141 the reliability for same instructor, same course, and different semesters was greater than

142 zero and usually overlapped in confidence interval with same instructor, same semester,

143 and different courses. Interestingly, the same instructor with different courses and

144 semesters showed a non-zero negative relationship, indicating that ratings generally were

145 lower for later semesters in different courses.

146     For different instructors, we found positive non-zero reliablities when they were at
147 least calculated on the same semester or course. These values were very close to zero,
148 generally in the .01 to .05 range. Last, the reliabilities that were calculated on different
149 courses, semesters, and instructors include zero in their confidence intervals. Exact values
150 can be found in the online supplemental document with the robustness analysis.

151     ROBUSTNESS REVEALED:

152 **RQ 2**

153     The paired evaluations were then filtered to only examine course and instructor
154 matches to explore the relation of reliability across time. Reliability was calculated by
155 calculating the partial correlation between the overall rating for the course first evaluation
156 and the overall rating for the course second evaluation, controlling for the number of
157 ratings within those average scores. This reliability was calculated separately for each
158 instructor and semester difference (i.e., the time between evaluations, 0 means same
159 semester, 1 means the next semester, 2 means two semesters later, etc.). The ratings were
160 filtered so that at least 10 pairs of ratings were present for each instructor and semester
161 difference combination (Weaver & Koopman, 2014). Of 36084 possible matched instructor
162 and course pairings, 30728 included at least 10 pairings, which was 1009 total instructor
163 and semester combinations.

164     The confidence interval for the effect of semester difference predicting reliability did
165 not cross zero, $b = -0.004$, 95% CI [-0.005, -0.003], $R^2 = .04$. The coefficient, while small,
166 represents a small effect of time on the reliability of instructor ratings. As shown in 2,
167 reliability appears to decrease across time.

### RQ 3

The confidence interval for the interaction of semester time difference and average fairness did cross zero, $b$ = -0.001, 95% CI [-0.007, 0.005], $R^2$ = .04. Therefore, there was no effect of the interaction of average fairness with semester differences in predicting reliability. Similarly, average fairness did not predict reliability overall, $b$ = -0.041, 95% CI [-0.226, 0.143].

### RQ 4

The confidence interval for the interaction of variability of fairness and semester time difference did cross zero, $b$ = -0.010, 95% CI [-0.022, 0.002], $R^2$ = .05. The variability of fairness also did not predict reliability overall, $b$ = 0.291, 95% CI [-0.091, 0.672].

## Discussion

- Summarize the results

### What Should I Do with This Information

- Don't expect to be reliable across other classes
- Don't expect to be reliable over long period of time, people change, students change, etc.

### Strengths

- a crap ton of data
- over a long period of time
- robust results

### Limitations

- one item versus many
- evaluations don't mean what we want them to mean

191      • one uni means maybe not generalizable

192 **Future Work**

<sup></sup>

# References

193

194 Aust, F., Barth, M., Diedenhofen, B., Stahl, C., Casillas, J. V., & Siegel, R. (2022).

195 *Papaja: Prepare american psychological association journal articles with r markdown.*

196 Retrieved from https://CRAN.R-project.org/package=papaja

197 Cohen, J., Cohen, P., West, S. G., & Aiken, L. (2003). *Applied multiple regression /*

198 *correlation analysis for the behavioral sciences* (3rd ed.). Lawrence Erlbaum Associates.

199 Gelman, A. (2006). Multilevel (hierarchical) modeling: What it can and cannot do.

200 *Technometrics*, *48*(3), 432–435. https://doi.org/10.1198/004017005000000661

201 Lüdecke, D., Makowski, D., Ben-Shachar, M. S., Patil, I., Højsgaard, S., Wiernik, B. M.,

202 . . . Luchman, J. (2023). *Parameters: Processing of model parameters.* Retrieved from

203 https://CRAN.R-project.org/package=parameters

204 Pinheiro, J., Bates, D., Debroy, S., Sarkar, D., & Team, R. C. (2017). *Nlme: Linear and*

205 *nonlinear mixed effects models.* Retrieved from

206 https://cran.r-project.org/package=nlme

207 Rantanen, P. (2012). The number of feedbacks needed for reliable evaluation. A multilevel

208 analysis of the reliability, stability and generalisability of students' evaluation of

209 teaching. *Assessment & Evaluation in Higher Education*, *38*(2), 224–239.

210 https://doi.org/10.1080/02602938.2011.625471

211 Tabachnick, B. G., Fidell, L. S., & Ullman, J. B. (2019). *Using multivariate statistics*

212 (Seventh edition). NY, NY: Pearson.

213 Weaver, B., & Koopman, R. (2014). An SPSS macro to compute confidence intervals for

214 pearson's correlation. *The Quantitative Methods for Psychology*, *10*(1), 29–39.

215 https://doi.org/10.20982/tqmp.10.1.p029

Table 1

*Descriptive Statistics of Included Courses*

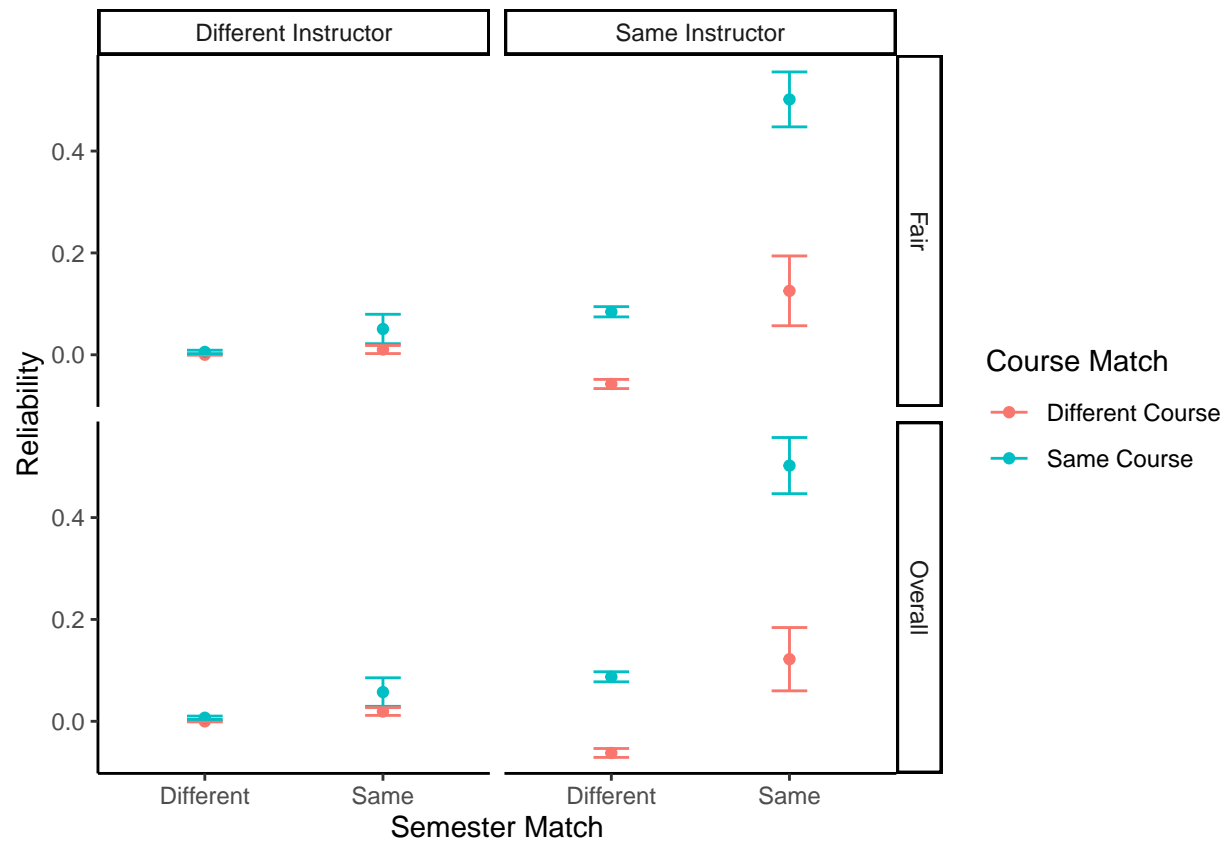| Statistic | Undergraduate | Mixed | Master's |
| --- | --- | --- | --- |
| N Total | 2898 | 274 | 42 |
| N Instructors | 223 | 40 | 10 |
| N Courses | 41 | 21 | 8 |
| Average N Ratings | 34.39 | 21.15 | 21.10 |
| Average Overall | 3.94 | 4.01 | 3.72 |
| SD Overall | 0.55 | 0.59 | 0.67 |
| Average Fairness | 4.46 | 4.50 | 4.19 |
| SD Fairness | 0.35 | 0.38 | 0.55 |
| Average Grade | 4.26 | 4.52 | 4.41 |
| SD Grade | 0.33 | 0.27 | 0.34 |

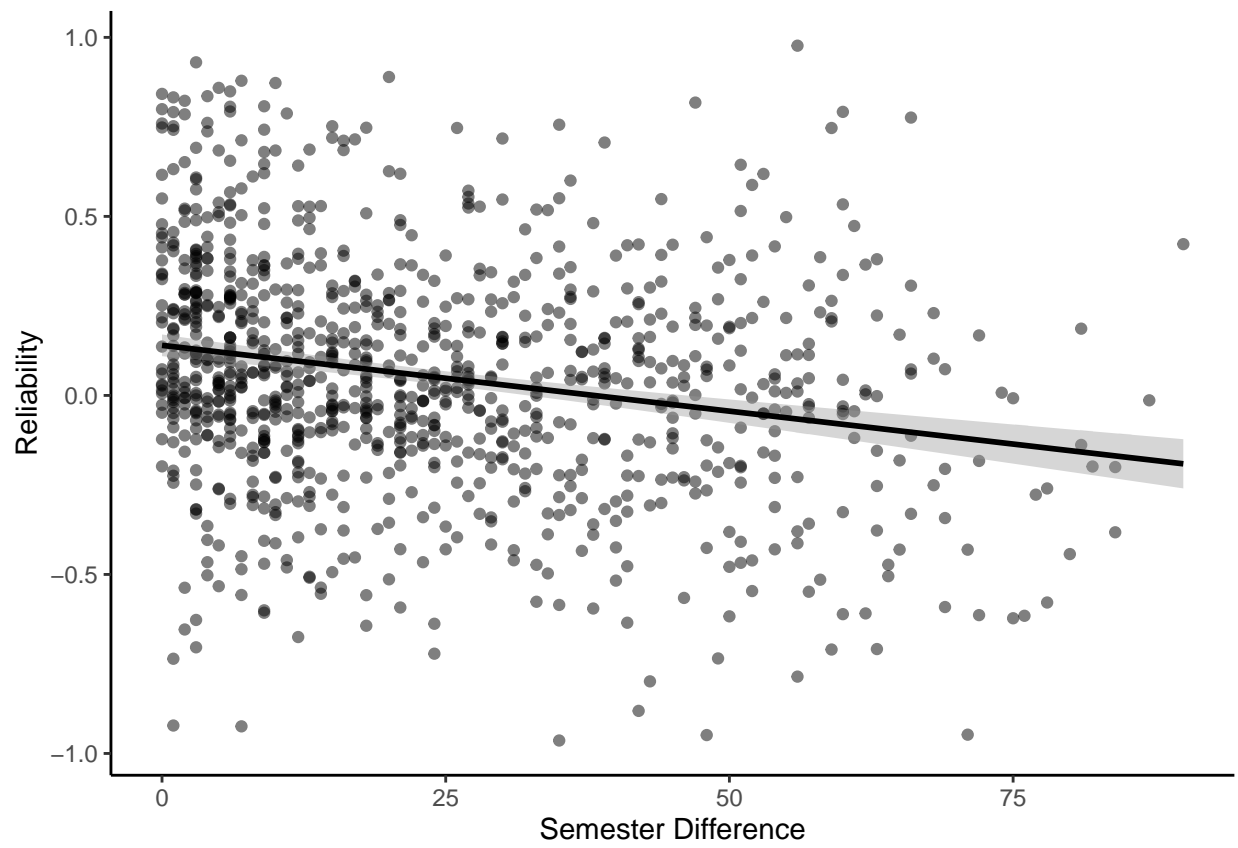*Figure 1*. Reliability estimates for instructor, course, and semester combinations.

*Figure 2*. Reliability estimates for same instructor and course across time.