ELSEVIER

# Grading leniency, grade discrepancy, and student ratings of instruction

## Bryan W. Griffin[*]

*Department of Curriculum, Foundations, and Reading, Georgia Southern University,
P.O. Box 8144, Statesboro, GA 30460, USA*

## Abstract

The purpose of this study was to examine how grading leniency and grade discrepancy (the difference between expected grades and deserved grades) were associated with various dimensions of student ratings of instruction. A sample of 754 undergraduate college students completed a student ratings of instruction instrument and provided responses to a number of other questions on topics such as course difficulty and workload. A series of multilevel regression analyses were conducted and results showed that an instructor's grading leniency, as perceived by students, was positively associated with student ratings on all dimensions of instruction examined. This finding suggests that more lenient instructors tend to receive higher student ratings. The second finding shows that grade discrepancy was negatively associated with most dimensions of instruction. This supports the self-serving bias hypothesis under attribution theory (Gigliotti & Buchtel, 1990) in that students tended to punish instructors with lower ratings when expected grades were lower than students believed they deserved, yet little evidence of a pattern of rewards existed in student ratings when students expected grades higher than they deserved.
© 2004 Published by Elsevier Inc.

*Keywords:* Student ratings of instruction; Student evaluations of instruction; Grading leniency; Grade discrepancy; Self-serving bias; Attribution theory

[*] Fax: 1-912-681-5382.
*E-mail address:* bwgriffin@gasou.edu.

25 **1. Introduction**

26    Student ratings are widespread and a common tool for evaluating faculty. When
27 asked, most faculty members approve of the use of student ratings of instruction for
28 teaching improvement (Baxter, 1991; Griffin, 1999; Moses, 1986; Schmelkin, Spen-
29 cer, & Gellman, 1997), but many are resistant to the use of student ratings for tenure,
30 promotion, and merit decisions (Feldman, 1997; McKeachie, 1997a). What many ed-
31 ucators believe is that student ratings are affected, or biased, by a number of factors
32 unrelated to teaching performance (Marsh & Overall, 1979; Wilson, 1998), and one
33 common concern is that grading standards employed by instructors could bias rat-
34 ings. As Marsh and Roche (2000) have noted, the average correlation between ex-
35 pected grades and student ratings of instruction is around .20. Typically this
36 relationship has been interpreted using one of three theoretical explanations (for re-
37 views see Greenwald & Gillmore, 1997a; Marsh & Roche, 2000; Wachtel, 1998).
38    First, the positive correlation between expected grade and student ratings of in-
39 struction may be explained as indicating a valid measurement of student ratings since
40 better instruction should result in more learning, better grades, and better ratings.
41 Second, the association between expected grades and ratings of instruction could
42 be spurious and produced by various student characteristics such as motivation.
43 For example, more motivated students who have greater interest in the subject mat-
44 ter are likely to learn more, achieve more, and rate the instructor higher. Third, an
45 association between expected grades and ratings could reflect some type of biasing
46 effect. For example, one possible biasing effect is grading leniency. Under this hy-
47 pothesis, instructors are rewarded with higher ratings for assigning higher grades
48 as a result of lenient grading practices, or conversely penalized with lower ratings
49 for assigning lower grades due to grading harshness. One important weakness of
50 studies examining the grading leniency hypothesis is that few have incorporated mea-
51 sures of student perceptions of the instructor's grading leniency (Marsh, 1987; Marsh
52 & Roche, 2000).
53    Olivares' (2001) was the only study found that incorporated a measure of grading
54 leniency. Olivares measured grading leniency by asking students to compare their
55 current instructor to others they have had and rate this instructor's grading from
56 1 ''much easier/lenient grader'' to 7 ''much harder/strict grader.'' Olivares found ze-
57 ro-order correlations of −.42 between grading leniency and an overall rating of the
58 instructor, and of −.45 between grading leniency and a composite rating of the in-
59 structor based on students' perceptions of the instructor's organization, communica-
60 tion, level of caring, and classroom atmosphere. Given the scoring system of the
61 rating scale used for grading leniency, the negative correlations indicate that more
62 lenient grading was associated with higher ratings of the instructors. Olivares
63 also found that the association between grading leniency and student ratings of
64 the instructor remained after controlling for pre-course interest, change in interest,
65 expected grade for the course, and a measure of cognitive ability.
66    In addition to the grading leniency hypothesis, another possible biasing effect in-
67 terpretation for the grades–ratings association can be found in the theories of attri-
68 bution and retribution (Feldman, 1997). Attribution theory suggests that a student

3

69  may react in one of two ways if that student receives a grade that differs from what
70  was expected. If the grade is lower than expected, then the student is likely to activate
71  a defensive mechanism commonly referred to as self-serving bias (Gigliotti & Buch-
72  tel, 1990). With self-serving bias, a student will attempt to protect his or her view of
73  self and assign blame for the lower than expected performance to an external cause.
74  The likely target will be the instructor, so the student will rate the instructor lower,
75  thus a rating penalty effect will occur. If a student receives a grade that is higher than
76  expected, then the student will assign credit to this performance to internal causes,
77  such as his or her intelligence, ability, hard work, etc. Since the better than expected
78  grade is seen as a result of the student's behavior or ability, ratings of the instructor
79  are not likely to differ from ratings given by students who receive grades as expected;
80  in essence, there is no rating reward effect. Further diminishing the possible rating
81  reward effect is the situation identified by Miller and Ross (1975) in which individ-
82  uals typically anticipate positive outcomes, so it is unlikely that many students will
83  acknowledge higher than expected grades since high grades were expected anyway.
84  In short, with attribution theory and self-serving bias, students are likely to penalize
85  instructors for lower than expected grades, but there is unlikely to be any reward ef-
86  fect for the few students who might believe they are receiving a grade higher than
87  expected. Retribution effect (Feldman, 1997) predicts simpler behavior on the part
88  of students. If, for example, a student receives lower than expected grades, this indi-
89  vidual will penalize the instructor, while a student who receives higher than expected
90  grades will reward the instructor.
91      One difficulty with student ratings research using the self-serving bias and retribu-
92  tion effect explanations has been the method for determining the
93  grade discrepancy—whether grades are higher or lower than what students expect.
94  The most direct method for assessing grade discrepancy is usually found in grade
95  manipulation experiments in which students are lead to anticipate one grade, but
96  then receive a grade inconsistent with their expectations (e.g., Abrami, Dickens, Per-
97  ry, & Leventhal, 1980; Tata, 1999; Worthington & Wong, 1979). Reviewers of these
98  studies, however, have pointed to a number of potential flaws. One important flaw is
99  that in classroom settings, often students do not know what their actual grade will be
100  before they complete instructional rating forms, so the external validity of these stud-
101  ies is limited. For correlational studies of attribution and retribution effects, re-
102  searchers often calculate grade discrepancy by considering pre-course grade point
103  average (GPA) or pre-course expected grade, and then examining how the end-of-
104  course expected grade or actual grade differs from the pre-course GPA or expected
105  grade (e.g., Gigliotti & Buchtel, 1990; Granzin & Painter, 1973; Greenwald & Gill-
106  more, 1997b; Palmer, Carliner, & Romer, 1978). A potential limitation of these de-
107  signs is that students are very likely to reassess their expectations once they are
108  exposed to the course and instructor, so pre-course grade expectation may provide
109  an inaccurate grade discrepancy baseline. Similarly, the use of GPA for determining
110  grade discrepancy could be misleading since performance, and expectation for per-
111  formance, in a given course can be independent of performance in other courses.
112  This does not mean that previous correlational studies are flawed or misleading,
113  but alternative methods for assessing grade discrepancy may prove useful.

4          *B.W. Griffin / Contemporary Educational Psychology xxx (2004) xxx–xxx*

114    The purpose of this study is twofold. First, since only one study of the grading
115  leniency hypothesis has incorporated a measure of leniency as perceived by students,
116  it is important to understand better how scores from such a measure relate to student
117  ratings, and to learn if the association between grading leniency and student ratings
118  replicates across studies. Second, the calculation of grade discrepancy for assessing
119  the self-serving bias and retribution effect hypotheses can be done in a manner that
120  is perhaps more course appropriate than previously examined. Thus, the intent of
121  this study is to examine the grading leniency explanation of student ratings by incor-
122  porating a measure of students' perceptions of leniency, and to test both self-serving
123  bias and retribution effect hypotheses by incorporating a more course specific mea-
124  sure of grade discrepancy.

125  **Method**

126  *Participants*

127    A total of 754 undergraduate students enrolled in 39 education courses at a me-
128  dium sized (14,000 students), Regional University in the southeastern United States
129  participated in this study. The classes ranged in size from 6 to 34 students. Under-
130  graduate education students at this institution are predominately White (71%) and
131  female (80%). Most respondents (76%) reported grade point averages in the range
132  of 2.5–3.5 on a 4.0 scale. Data were collected during the fall and spring semesters
133  of the 1998–1999 academic year.

134  *Instrument and variables*

135    An instrument to assess student evaluations of instruction and course character-
136  istics was developed drawing item and question wording from multiple sources (Ab-
137  rami, d'Apollonia, & Rosenfield, 1997; Feldman, 1997; Marsh, 1987; Murray, 1997).
138  To measure teaching effectiveness, 12 statements were used to assess multiple dimen-
139  sions of instruction with ratings following a five-point scale. The 12 statements fol-
140  low.
141    1.  Overall, how would you rate this course?
142    2.  Overall, how would you rate this instructor?
143    3.  The instructor was dynamic and energetic in conducting the course.
144    4.  The instructor presented the material in a clear and understandable manner.
145    5.  Course materials were well prepared and organized.
146    6.  Students were invited to share their ideas and knowledge.
147    7.  The instructor made students feel welcome in seeking help/advice in or outside of
148       class.
149    8.  The content of this course is useful, worthwhile, or relevant to you.
150    9.  Methods of evaluating student work were fair and appropriate.
151  10.  The instructor seems to have a real interest in and concern for students.
152  11.  The instructor gave students useful/helpful feedback on work.

5

153 12. The instructor is very knowledgeable in the subject of this course.
154    For the first 2 items, overall course and overall instructor, the scale ranged from 1
155 "Poor" to 5 "Excellent" and for the remaining 10 items the scale ranged from 1
156 "strongly disagree" to 5 "strongly agree."
157    The two predictors of interest in this study are grading leniency, which was as-
158 sessed by students' responses to this statement, "This instructor is a lenient/easy gra-
159 der" (1 "strongly disagree" to 5 "strongly agree"), and grade discrepancy, which was
160 calculated as the difference between the grade a student expected ("What grade do
161 you think the instructor will assign you in this course?") minus the grade a student
162 believed they deserved in the course ("What grade do you think you deserve in this
163 course?"). Both expected and deserved grades were assessed using a 12-point scale
164 (i.e., $A+ = 13$, $A = 12$, $A- = 11$, etc. through $D- = 2$, $F = 1$). The difference be-
165 tween expected minus deserved grade can be interpreted as follows: a positive differ-
166 ence indicates the expected was higher than the deserved grade (e.g., expect an $A-$
167 but deserve a $B+$), no difference shows expected and deserved are the same (e.g., ex-
168 pect a $B$ and deserve a $B$), and a negative difference shows that expected grade is low-
169 er than deserved grade (e.g., expect $B+$ and deserve $A-$).
170    In addition to these measures, students also provided information concerning: (a)
171 the instructor's reputation (1 "very bad" to 5 "very good," and 6 "didn't know about
172 the instructor"), (b) course difficulty (1 "one of easiest" to 5 "one of most difficult"),
173 (c) course workload (1 "very light" to 5 "very heavy"), (d) current GPA, and (e) pre-
174 course motivation ("You had a strong desire to take this course," with responses
175 ranging from 1 "strongly disagree" to 5 "strongly agree"). Class size and instructor's
176 sex were also included in the analysis. Three categories of instructor reputation were
177 developed for the analyses performed in this study: negative reputation, which in-
178 cluded students who selected responses 1–3 ("instructor very bad" to "about aver-
179 age") for the instructor reputation item; positive reputation, which included
180 students who choose responses 4 and 5 ("above average" to "instructor very good")
181 for the instructor reputation item; and no information, which consisted of students
182 who selected response 6 ("didn't know about the instructor") for the instructor rep-
183 utation item.
184    From these three categories of instructor reputation, two dummy variables (Pe-
185 dhazur, 1997) were created for the regression analyses performed below. The first,
186 called positive reputation, was coded 1 if student responses corresponded with the
187 positive reputation category, and 0 otherwise. The second dummy variable was la-
188 beled negative reputation and was coded 1 if student responses corresponded with
189 the negative reputation category, otherwise a 0 was used. Of the 754 respondents,
190 176 (23.3%) were classified into the positive reputation group, 420 (55.7%) into the
191 no information group, and 158 (21%) into the negative reputation group.
192    Evidence for construct validity for the scores obtained from this instrument and
193 sample can be assessed by examining correlations among scores from the dimensions
194 of instruction and various other course-related variables. Correlations and descrip-
195 tive statistics for the student-level variables are presented in Table 1. For example,
196 prior research has demonstrated a generally positive relationship between students'
197 pre-course motivation and students' ratings of instruction (Marsh, 1987), and a

YCEPS 1162
DISK / 17/2/04 / Vimala(CE)/ Jayanthi (TE)

6

B.W. Griffin / Contemporary Educational Psychology xxx (2004) xxx–xxx

ARTICLE IN PRESS

No. of pages: 16
DTD 4.3.1 / SPS

Table 1
Descriptive statistics and correlations among student-level variables

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1.000 | | | | | | | | | | | | | | | | | | | | |
| 2 | .789 | 1.000 | | | | | | | | | | | | | | | | | | | |
| 3 | .715 | .646 | 1.000 | | | | | | | | | | | | | | | | | | |
| 4 | .716 | .676 | .728 | 1.000 | | | | | | | | | | | | | | | | | |
| 5 | .654 | .618 | .700 | .759 | 1.000 | | | | | | | | | | | | | | | | |
| 6 | .479 | .391 | .499 | .492 | .481 | 1.000 | | | | | | | | | | | | | | | |
| 7 | .617 | .505 | .548 | .556 | .544 | .652 | 1.000 | | | | | | | | | | | | | | |
| 8 | .556 | .658 | .569 | .608 | .543 | .415 | .498 | 1.000 | | | | | | | | | | | | | |
| 9 | .604 | .519 | .543 | .571 | .570 | .602 | .637 | .467 | 1.000 | | | | | | | | | | | | |
| 10 | .675 | .557 | .635 | .628 | .618 | .644 | .762 | .511 | .703 | 1.000 | | | | | | | | | | | |
| 11 | .656 | .577 | .645 | .664 | .662 | .604 | .663 | .520 | .671 | .737 | 1.000 | | | | | | | | | | |
| 12 | .544 | .487 | .605 | .597 | .627 | .537 | .504 | .506 | .571 | .658 | .649 | 1. 000 | | | | | | | | | |
| 13 | .232 | .167 | .158 | .190 | .164 | .238 | .252 | .061 | .361 | .276 | .239 | .145 | 1. 000 | | | | | | | | |
| 14 | −.005 | −.051 | .033 | .007 | .003 | .049 | .062 | −.048 | .050 | .031 | .016 | −.008 | .055 | 1.000 | | | | | | | |
| 15 | −.235 | −.217 | −.147 | −.201 | −.153 | −.171 | −.249 | −.132 | −.310 | −.242 | −.205 | −. 087 | −.213 | −.114 | 1.000 | | | | | | |
| 16 | .216 | .202 | .150 | .176 | .153 | .109 | .127 | .174 | .122 | .134 | .120 | .117 | .084 | . 012 | −.058 | 1.000 | | | | | |
| 17 | −.356 | −.304 | −.213 | −.240 | −.199 | −.256 | −.294 | −.163 | −.369 | −.319 | −.284 | −. 218 | −.253 | −.053 | .229 | −.284 | 1.000 | | | | |
| 18 | .131 | .135 | .157 | .081 | .133 | .072 | .045 | .178 | .027 | .075 | .099 | .174 | −.337 | −. 101 | .199 | .012 | .112 | 1.000 | | | |
| 19 | .048 | .057 | .114 | −.021 | .083 | .015 | .023 | .039 | .021 | .039 | .089 | .086 | −.169 | −. 048 | .094 | −.047 | .067 | .478 | 1.000 | | |
| 20 | .366 | .496 | .361 | .381 | .341 | .203 | .276 | .486 | .293 | .304 | .349 | .254 | .113 | −. 034 | −.073 | .143 | −.141 | .098 | .134 | 1.000 | |
| 21 | .166 | .164 | .153 | .166 | .130 | .153 | .226 | .089 | .244 | .189 | .172 | .085 | .158 | . 094 | −.431 | .051 | −.174 | −.275 | −.110 | .116 | 1.000 |
| M | 3.86 | 3.50 | 4.06 | 3.87 | 4.12 | 4.52 | 4.26 | 4.04 | 4.19 | 4.27 | 4.13 | 4.47 | 2.94 | 0. 03 | 0.29 | 0.23 | 0.21 | 3.25 | 3.47 | 3.21 | 10.54 |
| SD | 1.16 | 1.13 | 1.11 | 1.15 | 1.02 | 0.81 | 0.99 | 1.14 | 1.02 | 0.97 | 1.01 | 0.80 | 1.16 | 0. 17 | 0.46 | 0.42 | 0.41 | 0.90 | 0.94 | 1.10 | 1.77 |

*Note.* Variables include: 1, Overall Instructor Rating; 2, Overall Course Rating; 3, Dynamic/Energetic Rating; 4, Presented Clearly Rating; 5, Materials Organized Rating; 6, Students Invited to Share Ideas Rating; 7, Students Could Seek Help Rating; 8, Course Content Worthwhile Rating; 9, Fair Evaluations Rating; 10, Instructor Show Interest in Students Rating; 11, Feedback Helpful Rating; 12, Instructor Knowledgeable Rating; 13, Grading Leniency; 14, Positive Discrepancy (coded 1 if grade higher than deserved, 0 otherwise); 15, Negative Discrepancy (coded 1 if grade lower than deserved, 0 otherwise); 16, Positive Reputation Dummy (1 if student rated instructor as having positive reputation, 0 otherwise); 17, Negative Reputation Dummy (1 if student rated instructor as having negative reputation, 0 otherwise); 18, Course Difficulty; 19, Course Workload; 20, Pre-course Motivation; 21, Expected Grade.

All correlations larger than .071 in absolute value are statistically significant at the .05 level.

*n* = 754.

7

198 similar pattern emerges for these data. Additionally, the grade students expect for a
199 course correlates positively with ratings for the course (Wachtel, 1998), and this pat-
200 tern also can be observed with these ratings. Similar findings exist for course work-
201 load and course difficulty (Greenwald & Gillmore, 1997a, 1997b; Marsh & Roche,
202 2000).

203 *Procedures*

204　Students in 39 classes were administered the evaluation instrument during the
205 last week of regular classes in the fall and spring semesters of the 1998–1999
206 academic year. Instructors were required to leave the classroom during evalua-
207 tions. Students were told that evaluations would not be made available until af-
208 ter course grades had been assigned and would only be provided to instructors
209 in aggregate form.

210 **Results**

211　Of the 754 students sampled, 67.8% ($n = 511$) believed that the grade they ex-
212 pected in the course was the grade they deserved, hence there was no difference be-
213 tween expected and deserved grade for these students. A total of 222 students
214 (29.4%) expected a grade lower than they deserved and only 23 students (3.1%) ex-
215 pected a grade higher than they deserved. Of the two competing theories, self-serving
216 bias and retribution effect, these data provide a better fit to the self-serving bias ex-
217 planation since so few students surveyed thought they were to receive a grade higher
218 than deserved. Miller and Ross (1975) predicted such behavior. It is also interesting
219 to note that the majority of students expected no discrepancy at all, so it is likely that
220 any grade discrepancy effect on student ratings of instruction may be small or limited
221 to only a minority of students overall.
222　To statistically model student ratings, it was necessary to create dummy variables
223 (Pedhazur, 1997) for grade discrepancy. The first, labeled positive discrepancy, was
224 created to represent those students who believed their expected grade would be high-
225 er than deserved. The coding for this dummy was 1 for students expecting grades
226 higher than deserved, and 0 for all other students. The second dummy variable,
227 called negative discrepancy, was created to represent those students who believed
228 their expected grade would be lower than their deserved grade, with coding of 1
229 for students expecting lower grades, and 0 for all others.
230　As the correlations in Table 1 show, grading leniency was positively corre-
231 lated with each of the 12 instructional rating items. The correlations ranged
232 from a low of .06 to a high of .36, with an average correlation of .21. The po-
233 sitive discrepancy dummy variable showed an inconsistent pattern of correla-
234 tions, with both positive and negative correlations with the 12 ratings items,
235 and with no correlation greater than .06 in absolute value. The negative discrep-
236 ancy dummy demonstrated a consistently negative pattern of correlations with
237 each of the 12 ratings items, with correlations ranging from −.08 to −.31. These

238 correlations indicate that students with lower expected than deserved grades
239 tended to rate the instructor and instruction lower on each of the 12 instruc-
240 tional rating items.
241     While the zero-order correlations are informative about the general nature of
242 the relationship among these variables, it is important to determine whether these
243 patterns of association remain once other predictors of student ratings are taken
244 into account in a regression equation. To learn whether grading leniency and
245 grade discrepancy are associated with student ratings of instruction, multilevel re-
246 gression (Bryk & Raudenbush, 1992; Goldstein, 1995; Longford, 1993) was used
247 in an effort to examine variation in student ratings both within and across classes.
248 Several researchers of student ratings of instruction (e.g., Cranton & Smith, 1990;
249 Feldman, 1998; Gigliotti & Buchtel, 1990) have noted that the level of analysis,
250 either student- or class-level, at which student ratings are examined could influ-
251 ence the nature of statistical relationships revealed. For example, the analysis
252 of class means rather than student-level data may obscure important variation
253 in ratings that result from individual student differences within the classroom.
254 Multilevel analysis allows one to combine both levels of analysis to provide a
255 more complete model of student ratings.
256     Incorporated into the multilevel analyses that follow were several covariates
257 previously identified as important predictors of student ratings of instruction.
258 At the student level, these covariates include course difficulty, course workload,
259 pre-course motivation, instructor reputation, and expected grade in the course.
260 Research on student ratings has demonstrated course difficulty and course work-
261 load, often measured together, to correlate positively with ratings of instruction
262 (Greenwald & Gillmore, 1997a, 1997b; Marsh, 1980; Marsh & Roche, 2000). In-
263 terest in the subject matter of the course before enrollment—pre-course motiva-
264 tion—has been linked to higher student ratings of instruction (Howard &
265 Maxwell, 1980; Marsh, 1980; Prave & Baril, 1993). Barké, Tollefson, and Tracy
266 (1983), Griffin (2001), and Ory (1980) found that instructor reputation was as-
267 sociated with various measures of teaching effectiveness. Finally, expected grade
268 in the course, which typically correlates positively with ratings, has been the
269 subject of much debate and research (Greenwald & Gillmore, 1997a; Marsh,
270 1987; Marsh & Roche, 1997, 2000; McKeachie, 1997b) and therefore was in-
271 cluded in the analysis.
272     At the class level, class size and instructor sex were included. Research shows that
273 class size correlates, albeit weakly, with ratings of instruction (Feldman, 1994). The
274 sex of the instructor also appears to relate to student ratings. Feldman's (1998) re-
275 views have shown that women tend to receive slightly higher ratings than men. How-
276 ever, Feldman (1998) also notes that a same-sex favorability in ratings exists;
277 students of the same sex as their instructor may provide slightly higher ratings (Cen-
278 tra & Gaubatz, 2000). Since the majority of students in the classes examined in this
279 study were women, it is likely that women instructors in this sample may have higher
280 ratings.
281     Thus, the models examined were, with variables enclosed in parentheses, as
282 follows:

283 *Student-level*

$$
\begin{aligned}
(\text{Student Rating of Instruction Item})_{ij} = {} & \beta_{0j} + \beta_1(\text{Grading Leniency})_{ij} \\
& + \beta_2(\text{Positive Discrepancy})_{ij} \\
& + \beta_3(\text{Negative Discrepancy})_{ij} \\
& + \beta_4(\text{Positive Reputation})_{ij} \\
& + \beta_5(\text{Negative Reputation})_{ij} \\
& + \beta_6(\text{Course Difficulty})_{ij} \\
& + \beta_7(\text{Course Workload})_{ij} \\
& + \beta_8(\text{Pre-course motivation})_{ij} \\
& + \beta_9(\text{Expected Grade})_{ij} + e_{ij}.
\end{aligned}
$$

285   At the class-level, mean ratings of the instructor were modeled with class size and
286 instructor sex:

287 *Class-level*

$$
\beta_{0j} = \gamma_{00} + \gamma_{01}(\text{Instructor's Sex})_j + \gamma_{02}(\text{Class Size})_j + \mu_{0j}.
$$

289   Combining the student- and class-level equations yields the following model of in-
290 structor rating:

291 *Combined*

$$
\begin{aligned}
(\text{Student Rating of Instruction Item})_{ij} = {} & \gamma_{00} + \beta_1(\text{Grading Leniency})_{ij} \\
& + \beta_2(\text{Positive Discrepancy})_{ij} \\
& + \beta_3(\text{Negative Discrepancy})_{ij} \\
& + \beta_4(\text{Positive Reputation})_{ij} \\
& + \beta_5(\text{Negative Reputation})_{ij} \\
& + \beta_6(\text{Course Difficulty})_{ij} \\
& + \beta_7(\text{Course Workload})_{ij} \\
& + \beta_8(\text{Pre-course motivation})_{ij} \\
& + \beta_9(\text{Expected Grade})_{ij} \\
& + \gamma_{01}(\text{Instructor's Sex})_j \\
& + \gamma_{02}(\text{Class Size})_j + e_{ij} + \mu_{0j}.
\end{aligned}
$$

293   This combined model was used to estimate the regression coefficients for each of
294 the 12 rating items presented above. Multilevel regression results, using full informa-
295 tion maximum likelihood to obtain estimates (Hox, 1995), are presented in Table 2.

YCEPS 1162
DISK / 17/2/04 / Vimala(CE)/ Jayanthi (TE)

10

ARTICLE IN PRESS

No. of pages: 16
DTD 4.3.1 / SPS

B.W. Griffin / Contemporary Educational Psychology xxx (2004) xxx–xxx

Table 2
Multilevel regression results for student ratings of instruction

| | Overall Instructor | | Overall Course | | Dynamic and Energetic | | Presented Clearly | | Materials Organized | | Students Shared Ideas | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | B | SE B | B | SE B | B | SE B | B | SE B | B | SE B | B | SE B |
| *Fixed Portion of Model* | | | | | | | | | | | | |
| Student Level | | | | | | | | | | | | |
| Grading Leniency | .12* | .03 | .06* | .03 | .11* | .03 | .12* | .03 | .10* | .03 | .12* | .03 |
| Grade Discrepancy | | | | | | | | | | | | |
| Positive Discrepancy | −.14 | .17 | −.32 | .17 | .19 | .18 | .06 | .18 | .01 | .18 | .20 | .15 |
| Negative Discrepancy | −.24* | .07 | −.23* | .07 | −.10 | .07 | −.21* | .08 | −.11 | .08 | −.11 | .06 |
| Instructor Reputation | | | | | | | | | | | | |
| Positive Reputation | .21* | .08 | .10 | .07 | .08 | .08 | .07 | .08 | .09 | .08 | .04 | .07 |
| Negative Reputation | −.39* | .10 | −.32* | .09 | −.19* | .10 | −.08 | .10 | −.13 | .10 | −.28* | .08 |
| Course Difficulty | .17* | .04 | .13* | .04 | .13* | .04 | .13* | .05 | .11* | .04 | .15* | .04 |
| Course Workload | .00 | .04 | .02 | .04 | .01 | .04 | −.07 | .04 | .03 | .04 | −.03 | .04 |
| Pre-course Motivation | .20* | .03 | .32* | .03 | .18* | .03 | .20* | .03 | .17* | .03 | .08* | .03 |
| Expected Grade | .08* | .02 | .07* | .02 | .07* | .02 | .10* | .02 | .06* | .02 | .06* | .02 |
| Intercept | 2.05* | .46 | 1.80* | .42 | 2.30* | .46 | 2.40* | .47 | 2.65* | .42 | 2.97* | .33 |
| Class Level | | | | | | | | | | | | |
| Class Size | −.01 | .01 | −.02 | .01 | −.01 | .01 | −.02 | .01 | −.01 | .01 | .00 | .01 |
| Instructor's Sex | −.54* | .21 | −.41* | .18 | −.46* | .20 | −.41* | .20 | −.36* | .16 | −.10 | .11 |
| *Random Portion of Model* | | | | | | | | | | | | |
| Class-level variance | .35* | | .27* | | .35* | | .33* | | .20* | | .08* | |
| Student-level variance | .62* | | .57* | | .64* | | .68* | | .64* | | .48* | |
| $R^2$ (total variance modeled) | .32 | | .36 | | .22 | | .24 | | .20 | | .17 | |

YCEPS 1162
DISK / 17/2/04 / Vimala(CE)/ Jayanthi (TE)

ARTICLE IN PRESS

No. of pages: 16
DTD 4.3.1 / SPS

B.W. Griffin / Contemporary Educational Psychology xxx (2004) xxx–xxx

11

Table 2 (*continued*)

| | Students Could Seek Help | | Course Content Worthwhile | | Fair Evaluation of Students | | Interest in Students | | Feedback Helpful | | Instructor Knowledgeable | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | B | SE B | B | SE B | B | SE B | B | SE B | B | SE B | B | SE B |
| *Fixed Portion of Model* | | | | | | | | | | | | |
| Student Level | | | | | | | | | | | | |
| Grading Leniency | .13* | .03 | .03 | .03 | .19* | .03 | .13* | .03 | .14* | .03 | .08* | .03 |
| Grade Discrepancy | | | | | | | | | | | | |
| Positive Discrepancy | .26 | .18 | −.12 | .19 | .12 | .17 | .10 | .17 | .06 | .17 | −.01 | .15 |
| Negative Discrepancy | −.25* | .08 | −.19* | .08 | −.31* | .07 | −.23* | .07 | −.17* | .07 | −.01 | .06 |
| Instructor Reputation | | | | | | | | | | | | |
| Positive Reputation | .05 | .08 | .05 | .08 | −.02 | .07 | .04 | .07 | .02 | .08 | .00 | .07 |
| Negative Reputation | −.37* | .09 | −.27* | .10 | −.49* | .09 | −.36* | .09 | −.30* | .09 | −.32* | .08 |
| Course Difficulty | .15* | .04 | .18* | .05 | .13* | .04 | .14* | .04 | .11* | .04 | .16* | .04 |
| Course Workload | −.01 | .04 | −.01 | .04 | .02 | .04 | .02 | .04 | .07 | .04 | .01 | .04 |
| Pre-course Motivation | .15* | .03 | .37* | .03 | .16* | .03 | .16* | .03 | .20* | .03 | .09* | .03 |
| Expected Grade | .07* | .02 | .04 | .02 | .08* | .02 | .06* | .02 | .08* | .02 | .04* | .02 |
| Intercept | 2.54* | .38 | 2.43* | .44 | 2.05* | .38 | 2.69* | .38 | 2.10* | .40 | 3.39* | .33 |
| Class Level | | | | | | | | | | | | |
| Class Size | −.01 | .01 | −.01 | .01 | −.01 | .01 | −.01 | .01 | −.01 | .01 | −.02 | .01 |
| Instructor's Sex | −.37* | .12 | −.43* | .17 | −.15 | .14 | −.29* | .14 | −.30* | .15 | −.20 | .11 |
| *Random Portion of Model* | | | | | | | | | | | | |
| Class-level variance | .09* | | .23* | | .13* | | .14* | | .16* | | .08* | |
| Student-level variance | .65* | | .71* | | .57* | | .56* | | .62* | | .46* | |
| $R^2$ (total variance modeled) | .26 | | .31 | | .32 | | .27 | | .25 | | .17 | |

*Note*. Positive Discrepancy coded 1 if expected grade is higher than believed deserved, 0 otherwise; Negative Discrepancy coded 1 if expected grade lower than believed deserved, 0 otherwise; Positive Reputation dummy coded 1 if student rated instructor as having positive reputation, 0 otherwise; and Negative Reputation dummy coded 1 if student rated instructor as having negative reputation, 0 otherwise. $R^2$ is calculated in the normal manner (Pedhazur, 1997), but model variance is calculated by summing both the between and within class variances (Snijders & Bosker, 1999).

$n = 754$ students in 39 courses.

* $p < .05$.

296    The regression results in Table 2 indicate that grading leniency was statistically
297 and positively related to 11 of the 12 rating items. The weakest relationship
298 ($b = .03$) was with the course content item, and this was the only partial coefficient
299 for grading leniency that was not statistically significant. The strongest relationship
300 ($b = .19$) was with the fair evaluation of students item. The latter coefficient may be
301 interpreted as showing that the more lenient the instructor's grading, the more fair
302 and appropriate was judged the instructor's evaluations of students' work. The aver-
303 age partial regression coefficient for the 12 items was .11. To put these estimates into
304 perspective, consider the situation of examining the single overall instructor rating
305 item for which the grading leniency regression estimate is $b = .12$. Assuming that
306 all other factors are held constant, two instructors who differ only on perceived grad-
307 ing leniency by one standard deviation ($SD = 1.16$, see Table 1) could expect an av-
308 erage mean difference of $1.16 \times .12 = .14$ points on their overall instructor rating
309 item. On the extremes, one instructor judged the least lenient (rating = 1) and an-
310 other judged most lenient (rating = 5) would differ by $(5–1) \times .12 = .48$ points on
311 their average overall instructor rating; for example, say 4.48 vs. 4.00 on a scale of
312 1–5.
313    The relationship between grade discrepancy and student ratings was more com-
314 plex than that found with grading leniency. The positive discrepancy dummy vari-
315 able was positively related to 8 of the 12 ratings items, and negatively related to
316 the remaining 4 ratings items. In no cases were the coefficient estimates for this dum-
317 my variable statistically significant, and in all cases the standard errors for the coef-
318 ficients were relatively large, thus indicating great variability in the estimates. Given
319 the small sample size of students who thought their expected grade was higher than
320 their deserved grade ($n = 23$), such unreliable estimates should be expected. The re-
321 gression estimates obtained for the positive discrepancy dummy show a weak and
322 inconsistent pattern of rating behavior for this group of students.
323    Unlike the positive discrepancy dummy, the dummy variable negative discrepancy
324 demonstrated a consistent and negative pattern of rating behavior for students ex-
325 pecting grades lower than they perceive they deserved. The negative discrepancy
326 dummy was found to be negatively associated with student ratings in all cases,
327 and was statistically significant for 8 of the 12 ratings items. Since negative grade dis-
328 crepancy is a dummy variable, the regression coefficient may be interpreted as the
329 mean difference in student ratings between those students who expect a grade lower
330 than they deserve and everyone else. The largest difference ($b = −.31$) was for the fair
331 evaluation of students item, and the smallest difference ($b = −.01$) was found for the
332 instructor knowledgeable item. Drawing on the example above using the overall in-
333 structor rating item, consider two instructors who differ only in the expectations held
334 by their students regarding their expected and deserved grades. The overall instruc-
335 tor rating for the instructor with students who believe their expected grades will be
336 lower than they deserve will be −.24 points lower than the instructor whose students
337 do not anticipate any difference between their expected and deserved grades, e.g.,
338 4.00 vs. 4.24.
339    For the other variables included in the models, results mirrored findings from pre-
340 vious studies. The strongest predictor of ratings was pre-course motivation. The neg-

13

341 ative instructor reputation dummy variable was negatively related to each rating
342 item except for two. Course difficulty was consistently, and positively, related to
343 all rating items. The more difficult the course, as judged by students, the more posi-
344 tive were student ratings. Course workload was not statistically related to any of the
345 rating items. Expected grade was also positively and statistically related to 11 of the
346 12 rating items. The partial regression coefficients for expected grade ranged from a
347 low of .04 to a high of .10.

348 **Discussion**

349     Recall the three possible interpretations of the positive relationship between ex-
350 pected grade and student ratings of instruction: (a) valid teaching/learning associa-
351 tion, (b) spurious association, and (c) biasing effect. Two ways of expressing the
352 biasing effect were examined in this paper, grading leniency and grade discrepancy.
353 Grading leniency was positively, and linearly, associated with 11 of the 12 rating
354 items. The positive relationship means that students tended to rate higher those in-
355 structors judged to be more lenient graders, and, conversely, instructors with harsher
356 grading practices tend to receive lower ratings. This finding replicates that reported
357 by Olivares (2001) who also found that instructors with more lenient grading prac-
358 tices tended to have higher student ratings. On the basis of results from this study
359 and Olivares' study, it appears that students rate instructors who are lenient graders
360 higher than instructors who are less lenient with their grading.
361     Also examined was the relationship between student ratings and grade discrep-
362 ancy, which is defined in this study as the difference between students' expected grade
363 and perceived deserved grade. Two theoretical explanations for such an effect were
364 listed, self-serving bias and retribution effect. As noted, self-serving bias suggests that
365 students will penalize instructors for lower than deserved grades, but will not reward
366 instructors for higher than deserved grades. Retribution effect holds that students
367 will reward instructors for higher than deserved grades, and penalize instructors
368 for lower than deserved grades. The data examined here provide a better fit to the
369 self-serving bias hypothesis. Only about 3% of the students sampled expected grades
370 higher than they deserved, and about 29% expected grades lower than they deserved.
371 There was little evidence that those who expected higher than deserved grades re-
372 warded instructors with higher ratings when compared to ratings made by other stu-
373 dents in the sample. None of the regression estimates for this group of students was
374 statistically different from zero. There is, however, evidence of a penalty effect; stu-
375 dents who expected grades lower than they deserved consistently provided ratings
376 that were lower than other students. The differences, adjusted for the modeled cova-
377 riates, ranged from low of −.01 to a high of −.31, with the overall average of −.18.
378 This penalty effect is also consistent with findings of a grading harshness effect
379 (Marsh & Roche, 2000; Worthington & Wong, 1979) in which students rate lower
380 instructors perceived to grade harshly. Note, however, that Marsh and Roche
381 (2000) point out that the self-serving bias may not be a bias under certain conditions
382 for student ratings of instruction. Perhaps, for example, if a grade discrepancy is due

383  to factors unrelated to instruction or the instructor, then students may not provide
384  lower ratings. Unfortunately, the reason for a grade discrepancy was not assessed
385  this study, so it is impossible to know further what students were thinking when they
386  identified a grade discrepancy.
387      In summary, these results suggest two things. First, there may be a grading le-
388  niency effect in student ratings, but so far only this study and Olivares' (2001)
389  study have apparently examined directly students' perceptions of grading leniency.
390  Replication studies are needed to further evaluate this finding. Second, in addi-
391  tion to a possible grading leniency effect, there appears to be an association be-
392  tween a negative grade discrepancy and student ratings. This finding supports the
393  self-serving bias hypothesis in that students appear to penalize instructors when
394  grades are lower than expected, but do not reward instructors when grades are
395  higher than expected. Since grading leniency and grade discrepancy, both possible
396  parts of the biasing effect interpretation, were statistically controlled in the mul-
397  tilevel regression models, the partial regression coefficients for expected grade may
398  represent a more pure examination of the: (a) valid teaching/learning association
399  and (b) spurious association hypotheses. Several factors that could lead to the
400  spurious association effect were included in the regression models, such as pre-
401  course motivation, course difficulty and workload. It is possible, though, that
402  other factors could contribute to the observed relationship between expected
403  grade and ratings found in this and other studies. More careful examinations tak-
404  ing into account various motivational factors such as intrinsic and extrinsic mo-
405  tivation, personal control, and autonomy may prove useful in further elimination
406  of the spurious effects hypothesis. However, since at least part of the spurious as-
407  sociation and biasing effects hypotheses have been controlled in this study, that
408  means the relationships between expected grades and student ratings of instruc-
409  tion found in the current study probably can be explained, at least in part, by
410  the valid teaching/learning hypothesis. Thus, the results provided here suggest
411  that student ratings of instruction are probably a function of both valid teaching
412  and learning and some biasing due to grading leniency and grade discrepancy.

## 413  References

414  Abrami, P. C., d'Apollonia, S., & Rosenfield, S. (1997). The dimensionality of student ratings of
415      instruction: What we know and what we do not. In R. P. Perry & J. C. Smart (Eds.), *Effective teaching
416      in higher education: Research and practice* (pp. 321–367). New York: Agathon.
417  Abrami, P. C., Dickens, W. J., Perry, R. P., & Leventhal, L. (1980). Do teacher standards for assigning
418      grades affect student evaluations of instruction? *Journal of Educational Psychology, 72*, 107–118.
419  Barké, C. R., Tollefson, N., & Tracy, D. B. (1983). Relationship between course entry attitudes and end-
420      of-course ratings. *Journal of Educational Psychology, 75*, 75–85.
421  Baxter, E. P. (1991). The TEVAL experience, 1983–88: The impact of a student evaluation of teaching
422      scheme on university teachers. *Studies in Higher Education, 16*, 151–179.
423  Bryk, A. S., & Raudenbush, S. W. (1992). *Hierarchical linear models: Applications and data analysis
424      methods.* Newbury Park, CA: Sage.
425  Centra, J. A., & Gaubatz, N. B. (2000). Is there gender bias in student evaluations of teaching? *The Journal
426      of Higher Education, 70*, 17–33.

15

Cranton, P., & Smith, R. A. (1990). Reconsidering the unit of analysis: A model of student ratings of instruction. *Journal of Educational Psychology, 82*, 207–212.

Feldman, K. A. (1994). Class size and college students' evaluations of teachers and courses: A closer look. *Research in Higher Education, 21*, 45–116.

Feldman, K. A. (1997). Identifying exemplary teachers and teaching: Evidence from student ratings. In R. P. Perry & J. C. Smart (Eds.), *Effective teaching in higher education: Research and practice* (pp. 368–395). New York: Agathon.

Feldman, K. A. (1998). Reflections on the study of effective college teaching and student ratings: One continuing question and two unresolved issues. In J. C. Smart (Ed.), *Higher education: Handbook of theory and research* (pp. 35–74). New York: Agathon.

Gigliotti, R. J., & Buchtel, F. S. (1990). Attributional bias and course evaluations. *Journal of Educational Psychology, 82*, 341–351.

Goldstein, H. (1995). *Multilevel statistical models* (2nd ed.). London: Edward Arnold.

Granzin, K. L., & Painter, J. J. (1973). A new explanation for students' course evaluation tendencies. *American Educational Research Journal, 10*, 115–124.

Greenwald, A. G., & Gillmore, G. M. (1997a). Grading leniency is a removable contaminant of student ratings. *American Psychologist, 52*, 1209–1217.

Greenwald, A. G., & Gillmore, G. M. (1997b). No pain, no gain? The importance of measuring course workload in student ratings of instruction. *Journal of Educational Psychology, 89*, 743–751.

Griffin, B. W. (1999). *Results of the faculty survey on student ratings of instruction: Preliminary report*. Statesboro, GA, USA: Georgia Southern University, Student Ratings Committee.

Griffin, B. W. (2001). Instructor reputation and student ratings of instruction. *Contemporary Educational Psychology, 26*, 534–552.

Howard, G. S., & Maxwell, S. E. (1980). Correlation between student satisfaction and grades: A case of mistaken causation? *Journal of Educational Psychology, 72*, 810–820.

Hox, J. J. (1995). *Applied multilevel analysis*. Amsterdam: TT-Publikaties. Available on-line (March 6, 2000): <http://www.ioe.ac.uk/multilevel/workpap.html>.

Longford, N. T. (1993). *Random coefficient models*. Oxford, UK: Oxford University Press.

Marsh, H. W. (1980). The influence of student, course, and instructor characteristics on evaluations of university teaching. *American Educational Research Journal, 17*, 219–237.

Marsh, H. W. (1987). Students' evaluations of university teaching: Research findings, methodological issues, and directions for future research. *International Journal of Educational Research, 11*, 253–388.

Marsh, H. W., & Overall, J. U. (1979). *Validity of students' evaluations of teaching: A comparison with instructor self-evaluations by teaching assistants, undergraduate faculty, and graduate faculty*. Paper presented at the annual meeting of the American Educational Research Association, San Francisco, CA. (ERIC Document No. ED 177 205).

Marsh, H. W., & Roche, L. A. (1997). Making students' evaluations of teaching effectiveness effective: The critical issues of validity, bias, and utility. *American Psychologist, 52*, 1187–1197.

Marsh, H. W., & Roche, L. A. (2000). Effects of grading leniency and low workload on students' evaluations of teaching: Popular myth, bias, validity, or innocent bystanders? *Journal of Educational Psychology, 92*, 202–228.

McKeachie, W. J. (1997a). Good teaching makes a difference—and we know what it is. In R. P. Perry & J. C. Smart (Eds.), *Effective teaching in higher education: Research and practice* (pp. 396–408). New York: Agathon.

McKeachie, W. J. (1997b). Student ratings: The validity of use. *American Psychologist, 52*, 1219–1225.

Miller, D. T., & Ross, M. (1975). Self-serving biases in the attribution of causality: Fact or fiction? *Psychological Bulletin, 82*, 213–225.

Moses, I. (1986). Student evaluation of teaching in an Australian university staff perceptions and reactions. *Assessment & Evaluation in Higher Education, 11*, 117–129.

Murray, H. G. (1997). Effective teaching behaviors in the college classroom. In R. P. Perry & J. C. Smart (Eds.), *Effective teaching in higher education: Research and practice* (pp. 171–204). New York: Agathon.

Olivares, O. J. (2001). Student interest, grading leniency, and teacher ratings: A conceptual analysis. *Contemporary Educational Psychology, 26*, 382–399.

480  Ory, J. C. (1980). The influence of students' affective entry on instructor and course evaluations. *The*
481      *Review of Higher Education, 4*, 13–24.
482  Palmer, J., Carliner, G., & Romer, T. (1978). Leniency, learning, and evaluations. *Journal of Educational*
483      *Psychology, 70*, 855–863.
484  Pedhazur, E. J. (1997). *Multiple regression in behavioral research: Explanation and prediction* (3rd ed.).
485      Brace, New York: Harcourt.
486  Prave, R. S., & Baril, G. L. (1993). Instructor ratings: Controlling for bias from initial student interest.
487      *Journal of Education for Business, 68*, 362–366.
488  Schmelkin, L. P., Spencer, K. J., & Gellman, E. S. (1997). Faculty perspectives on course and teacher
489      evaluations. *Research in Higher Education, 38*, 575–592.
490  Snijders, T. A. B., & Bosker, R. J. (1999). *Multilevel analysis: An introduction to basic and advanced*
491      *multilevel modeling*. London: Sage.
492  Tata, J. (1999). Grade distributions, grading procedures, and students' evaluations of instructors: A justice
493      perspective. *Journal of Psychology Interdisciplinary & Applied, 133*, 263–271.
494  Wachtel, H. K. (1998). Student evaluation of college teaching effectiveness: A brief review. *Assessment &*
495      *Evaluation in Higher Education, 23*, 191–212.
496  Wilson, R. (1998). New research casts doubt on value of student evaluations of professors. *The Chronicle*
497      *of Higher Education*, A12–A14.
498  Worthington, A. G., & Wong, P. T. P. (1979). Effects of earned and assigned grades on student evaluations
499      of an instructor. *Journal of Educational Psychology, 71*, 764–775.