Artificial Intelligence formulated this projection for compatibility purposes from the original article published at Global Journals. However, this technology is currently in beta. *Therefore, kindly ignore odd layouts, missed formulae, text, tables, or figures.*

Student Age as an Impact Factor for Student Evaluations of Instruction Katja Specht Received: 6 June 2015 Accepted: 4 July 2015 Published: 15 July 2015

6 Abstract

⁷ Student Evaluations of Instruction (SEI) are an important issue in countries like the USA,

* where the evaluation results can impact professional promotion chances and salary of faculty.

9 According to Seldin [11], the percentage of American colleges using SEI grew from 29

10

11 Index terms— evaluation, extrinsic impacts, generalized linear models, regression, student age

12 1 Introduction

tudent Evaluations of Instruction (SEI) are very widespread and common practice in countries, where the evaluation results are applied for professional promotion chances and salary of faculty. In countries like Germany, SEI are used as an instrument for the internal quality management and teaching improvement process. This instrument is also of growing interest in the accreditation process of study programs and universities.

17 The intrinsic impact factors of the evaluation ratings are the single items of the evaluation questionnaire, which are answered by the students. But many statistical investigations have shown that there are undesirable 18 extrinsic factors, like class size or the quantitative exposition of the course, which are noninstructional by nature 19 and, therefore, should be eliminated for a fair comparison of the evaluation ratings. Costin, Greenough and 20 Menges [2] presented a review of empirical studies regarding student ratings. They concluded that SEI can 21 provide reliable and valid information on the quality of courses and instruction but for further interpretation 22 23 extrinsic factors should be taken into account. Already Heilman and Armentrout [6], Lovell and Haner [8], McDaniel and Feldhusen [9] and Hamilton ??5] have shown that teachers of large classes may receive lower 24 ratings. Hoefer, Yurkiewicz and Byrne [7] assessed significant differences between undergraduate and SEI. For 25 that matter, Brightman [1] states that it is unfair to compare a faculty member teaching a required core class 26 with another faculty member teaching a senior-level elective course. Peterson, Berenson, Misa and Radosevich 27 [10] have recommended to establish appropriate sets of norming reports in which possible semester factor effects 28 are considered. 29

It is tempting to perform a linear regression of the evaluation ratings on the non-instructional factors by 30 the least-squares principle and to use the estimated model for the compensation procedure. But, this will be 31 admissible, if the latent variable is normally distributed. This can be tested by using the residuals from the 32 regression as a proxy for the latent variable. Frequently, a dependent variable, like evaluation ratings, is skewed 33 34 to the right. This, in turn, usually prevents the residuals from being normal. At least, this occurs with our data. 35 Therefore, our investigation focuses on a proper methodical approach of estimating a non-linear model. After 36 a description of the data we shall present the Maximum-Likelihood (ML) estimation of a so-called Generalized Linear Model (GLM) step-by-step. The presentation is sufficiently detailed, so that the reader can, for instance, 37 apply the procedure to own data with a matrix-based programming software like MATHLAB or GAUSS. We 38 restrict our presentation to one non-instructional factor, namely 'student age' or, more precisely, the semester 39 counter of the evaluated course. The proposed procedure can easily and obviously be extended to more non-40 instructive factors. Eventually, we shall show how to use the estimated model to correct actual and future 41 evaluation ratings properly. 42

43 **2** II.

44 **3** Data

We have collected n = 140 evaluation ratings z i from seven-semester Bachelor programs from the Business Unit 45 of a German University of Applied Sciences together with the semester counter (one to seven), to which the 46 47 evaluated course regularly belongs. The evaluation ratings are means from a five-point Likert S scale, where 48 the choice 'one' is best and 'five' is worst. Unfortunately, the evaluation ratings are not normally distributed. 49 More precisely, the standardized measure of skewness is 1.05 and the standardized measure of kurtosis is 4.51, 50 indicating that the dependent variable is skewed to the right with a kurtosis much larger than that of the normal distribution. This results in non-normal residuals from a linear least-squares regression. And this prevents 51 inferential conclusions of such a regression, like t-values and p-values. The usual methodology is no longer valid 52 in this case. 53

Luckily, a Box Cox transformation of the evaluation ratings z i can convert the ratings in (approximately) normally distributed values y i :

⁵⁶ 4 III. Methodology and Exemplary Results

57 5 a) GLM estimation

The most general form of a regression model explains a variable by the sum of its (conditional) expected value and of some noise: X denotes the design matrix. In our example, it consists of a first column of ones, representing the constant, and a second column with the semester counts. Further columns may be appended for additional non-instructional factors. The latent variables ? i are independent and identically (iid) distributed, representing the noise.

In a GLM, the dependent variable must belong to the exponential family and its expected value, given the design matrix X, may be a non-linear function h of the linear predictor X?:

In our example, the column vector ? consists of two unknown parameters, ? 0 and ? 1 , and h is the inverse of the link function g and is called ' response function'.

The following ML estimation procedure is explained in more detail in Fahrmeir, Kneib, and Lang [4]. Let x i be the i-th row of the design matrix X. Then we need the following symbols:

69 Consequently, the following diagonal matrices depend on ?:

The goal is to receive a solution of the nonlinear equation system s(?) = 0, where s(?) is the functional vector of partial derivatives of the loglikelihood function:

Now, the ML estimator may be iteratively approximated by the following equations: i := g (z i) := z ? i ?73 1? ? N(μ , ? 2)

The value of ?, which minimizes the absolute ske wness of the transformed variables can be calculated numerically and is about 0.45 for our data. If we apply the rounded value 0.5, then we receive a standardized measure of kurtosis of about 2.85. The hypothesis of normality for the transformed variables y i cannot be rejected by any test. D'Agostino, Belanger, and D'Agostino Jr. [3] yields a pvalue above 90%. The transforming function g is called 'link function'.

The normal distribution belongs to the so-called 'exponent ial family'. This admits the estimation of a GLM, which will be specified in the next section. measure of skew-ness of about 0.03 and a standardized

85 The skew-ness-kurtosis test of

86 We have started the iterations with the leastsquares estimator

In order to estimate ? 2, which depends on , we first have to eliminate duplicate rows in X. We denote the reduced design matrix by . Note, that in our example it has just seven rows due to the seven semester counts. The y i have to be averaged to within the seven groups of identical rows of X. Let n j denote the number of observations in group j. Then, the variance can be estimated in each step of the iteration:

Here, p is the number of columns of X, in our example: p = 2.

- ⁹² Table ?? shows the five iterations, which are needed for convergence in our example.
- ⁹³ Table ?? : Iterations of the Fisher-Scoring algorithm Therefore, we receive the following estimated model:

The residuals from this model are clearly normal. Thus, they can be 'studentized' in order to eliminate outliers. In a first step, the residuals ? ?i have to be 'standardized':

- 96 In a second step, the standardized residuals will be transformed into a Student distribution:
- We choose to define an outlier as an observation with an absolute studentized residual above the percentage point of order 0.975. This yields a 5% probability of an error of first kind. In our example we have excluded ten
- $_{99}$ observations leading to n = 130 observations, to which the whole procedure is applied again. This final estimation

100 yields:

¹⁰¹ 6 b) Model diagnostics

For model diagnostics, we can test the hypothesis H 0 : C? = c by the asymptotically distributed Wald statistic: where r is the rank of C and X ? WX is the Fisher information matrix. In our example, the Wald statistic for H 0 : = 0 amounts to 34.32 with a p-value of almost ?2 (?(k+1)) = 1 7 ? p ? 7 j=1 n j ? ?j ? h x ? j ?(k+1) k ?(k) 0 ?(k) 1 s(?(k) 0) s(?(k) 1) 1r * i := r i ? n ? p ? 1 n ? p ? r 2 i ? t n?p?1 ? = h(X ?) with ? = (?0.5078, ?0.0555) ? w = (C ? ? c) ? (CF (?) ?1 C ?) ?1 (C ? ? c) ? ? 2 r

¹⁰⁷ zero. And the Wald statistic for H 0 := 0 amounts to 6.79 with a p-value of 0.0091. Thus, both coefficients ¹⁰⁸ are highly significant.r i := ?i ? ? ? 1 ? h ii with h ii = x ? i (X ? X) ?1 x i ? = h(X ?) with ? = (?0.5012, ¹⁰⁹ ?0.0524) ? ?(0) = (X ? X) ?1 X ? y . X ? y j , j = 1, ... 7, ? 2 r - F (?) = ? 1 ? 0 Year 2015

¹¹⁰ The ML estimator is (approximately) normally distibuted with covariance matrix . Then , the transformed

111 variable y may be estimated or predicted like this: c) Back transformation Eventually, we have to come back to

the original evaluation ratings z i. To this end, we apply a Taylor series approximation of the response function h, centered at

114 The Taylor series approximation of the response function enables the conclusion for the evaluation ratings:

Because the expectation values of odd powers in the Taylor series are zero, the approximation error (with some ?? [0, 1]) is limited to: This may be imagined to be negligible.

117 Table ?? shows the estimated evaluation ratings in the last column for each group of identical covariables:

¹¹⁸ 7 Table 2 : GLM-estimated evaluation ratings

119 It is clearly seen that the expected ratings in the last column are falling, and therefore getting better, with raising 120 semester count in the third column. Thus, advanced students tend to be more patient with instructors.

¹²¹ 8 d) Compensation

In the simple linear model, the elimination of the impact of the 'extrinsic' factors in X is realized by the correction of the mean value of the dependent variable by the individual residual ? ? of an actual or future observation h(y) $= (1 ? y / 2) ? 2 ? h(y) := (1 ? \mu / 2) ? 2 + (1 ? \mu / 2) ? 3 (y ? \mu) + 3 4 ? (1 ? \mu / 2) ? 4 (y ? \mu) 2 ? = E (z|X) =$ $E (h(y)|X) ? E (h(y)|X) = (1 ? \mu / 2) ? 2 + 3 4 ? (1 ? \mu / 2) ? 4 ? V (y | X) = (1 ? h(X ?) / 2) ? 2 + 3 4 ? (1 ?$ $h(X ?) / 2) ? 4 ? ? 2 (?) h ???? (?y + (1 ? ?) \mu) 4! ? (y ? \mu) 4 < 0.0004 i X i, 1 X i, 2 ? i E (z i | X) 1 1 1 0.6134$ 2.y ? = x ? ? ? + ?? ? y * ? = ? + ?? = ? + y ? ? x ? ? ? y ? , x ? ?) : (y = X ? + ? ? F (?) ? 1 ? = E (y) $|X) = h(X ?) \mu = h(X ?)$:

129 Table ?? illustrates the consequences of these compensations for some randomly chosen ratings.

130 Table ?? : Some examples of proper corrections of evaluation ratings.

The arbitrary ratings z ? are corrected into the expected direction and yield the values in the last column. The ratings of early semesters are lowered, thus improved, and ratings of late semesters are raised, thus penalized.

¹³³ 9 e) Semester dummies

Now, we are going to model the impact of the categorical variable ' semester count' by semester dummies. This will drop the assumption of a monotonous influence in favour of more flexibility. We choose the first semester as the reference category. The dummy variables S i , i = 2, ..., 7, are defined to be 'one', if the course is affiliated to semester i, and 'zero' otherwise. The related GLM reads: with the (n \times 7)-dimensional design matrix

The estimation procedure is the same as before. Six outliers can be identified in this model, leaving behind a sample number of n = 134 and the following vector of estimated coefficients: Again, the conclusion for the original ratings is performed by a Taylor series approximation of the response function. This yields the following expected evaluation rating values, dependent on the semester count:

Evidently, with our data the evaluation ratings are 'raising' in the beginning and in the last three semesters. And they are 'falling back' in the middle part of the study program. But, remember that evaluation ratings are like 'grades' in our example, meaning that a 'high rating' is equivalent to a 'low grade'.

The residuals of this regression are clearly normal. The p-value of the skewness-kurtosis test is about 45%. The simultaneous significance of the dummy variables may be tested by the hypothesis H 0 : C? = c with:

Table ?? demonstrates the way of compensation for the non-instructional factor ' semester count' for seven exemplary evaluation ratings. Observed ratings z ? have to be reduced (i.e. improved) in the first four semesters and else raised (i.e. deteriorated). The corrected rates are listed in the last column.

Table ?? : Some examples of proper corrections of evaluation ratings with semester dummies. z ? g (z ?) x156 ? [2] h(x ? ? ?) ? y * ? z * ? 2.7 0

157 **10** Conclusions

Evaluation ratings are an important instrument in quality management of teaching. Several noninstructional factors may bias the intended evaluation of the instructor. It is essential to assess the quantitative influence of those non-instructional factors in order to compensate the evaluation ratings for these extrinsic factors and achieve a fair comparison.

It is tempting to perform a linear least-squares regression of the evaluation ratings on the noninstructive factors. The estimated model could easily be used to eliminate the extrinsic impact. But, if the residuals from this regression are not normally distributed, the results will not be reliable. Another method of estimation has to be applied.

At least with our SEI data, the residuals from a linear least-squares regression on student's age are skewed and far from beeing normal. But a proper Box-Cox transformation of the evaluation ratings yields a normally distributed dependent variable. This, in turn, enables the maximum likelihood estimation of a GLM. This procedure is not quite common. Therefore, it is explained in detail in this paper.

170 Once we have estimated a valid model, we can use it to eliminate the impact of the considered covariable.

Due to the non-linear GLM approach, this task requires a Taylor series approximation of the response function, which can be fairly easily performed. In our example, the expected evaluation ratings are getting better with rising semester count. Students seem to get more indulgent with growing age.

Finally, we have conducted the GLM regression of the transformed evaluation ratings on semester dummy variables. Now we receive more flexible, nonmonotonic impacts of the semester count on evaluation ratings. Especially with small data sets, this might be the better approach.

177 **11** ""

- 178 An important message of this paper should be to carefully inspect the assumptions of an applied method. In
- 179 many cases, these assumptions may not be met by the data. In these cases a less familiar procedure may serve as an alternative.



Figure 1:

180

6	.9602	1	0.6088	0.5395	
			$0.8910 \ 3.2521$		
Volume XV	$1.5\ 0.3670\ 2.4\ 0.7090\ 1.4\ 0.3097$	1	0.6088	0.5395	
Issue IV Ver-		2	0.2977	1.3804	
sion I		7	0.5833	0.5395	
			0.6652	2.2452	
			0.4768	0.5395	
			$0.3724 \ 1.5099$		
() G					
Global Jour-	y = h(X ?) + ? = h(? 0 + ? 1 ? S 2 + ? ?	? ? + ? 6 ?	? S 7) +	? ? = (?0.6484, 0.1065, 0.054)	
nal of Man-					
agement and					
Business Re-					
search					
	@ 2015 Global Journals Inc. (US) 1				

Figure 2:

- [Agostino et al. ()] 'A suggestion for using powerful and informative tests of normality'. R B Agostino , A
 Belanger , R B D'agostinoJr . *The American Statistician* 1990. 44 (4) p. .
- [Peterson et al. ()] 'An Evaluation of Factors Regarding School'. R L Peterson , M Berenson , R B Misra , D J
 Radosevich . Decisio n Sciences Journal of Innovative Education 2008. 6 (2) p. .
- [Costin et al. ()] F Costin , W T Greenough , R J Menges . Student Ratings of College Teaching: Reliability,
 Validity, and Usefulness, 1972. 41 p. .
- [Lovell ()] 'Haner: Forced-choice applied to college faculty rating'. G D Lovell, CF. Educational and Psychological
 Measurement 1955. 15 p. .
- [Mcdaniel and Feldhusen ()] E D Mcdaniel , J F Feldhusen . Relationship between faculty ratings and indexes of service and scolarship, 1970. 5 p. .
- [Brightman ()] 'Mentoring faculty to improve teaching and student learning'. H J Brightman . Decision Sciences
 Journal of Innovative Education 2005. 3 p. .
- 193 [Fahrmeir et al. ()] Regression, L Fahrmeir, T Kneib, S Lang. 2009. Heidelberg: Springer. (2nd ed.)
- [Hoefer et al. ()] 'The Association between Students' Evaluation of Teaching a nd Grades'. P Hoefer , J
 Yurkiewicz , J C Byrne . Decision Sciences Journal of Innovative Education 2012. 10 (3) p. .
- [Seldin ()] 'The use and abuse of student ratings of professors'. P Seldin . The Chronicel of Higher Education,
 197 1993. 21.