Contents lists available at ScienceDirect

# International Review of Economics Education

# Teaching econometrics with data on coworker salaries and job satisfaction☆

Todd Easton

*Pamplin School of Business, University of Portland, 5000 N. Willamette Blvd, Portland, OR, 97203, United States*

ARTICLE INFO

ABSTRACT

Recent research studied how relative salary affects job satisfaction. It gave a random sample of University of California employees information about their coworkers' salaries and estimated the effect of this information on job satisfaction. This article suggests ways the dataset created by this research can be used in econometrics and statistics classes. It provides examples using these data to calculate frequency distributions, contingency tables, Chi-square tests, and linear probability models. It also explains how these examples can be used productively in class.

## 1. Introduction

This article introduces a valuable dataset for teaching introductory-level econometrics and statistics classes. In addition, it provides examples of how the dataset can be used to teach effectively. The dataset is valuable because it allows an instructor to use real data, take advantage of students' interest in income inequality, and connect his/her class to topics students study elsewhere.

The dataset was developed as part of seminal research by Card et al. (2012). They asked how earnings inequality within units of an organization affects employees' job satisfaction. I show how instructors can use the dataset to provide students opportunities to practice using frequency distributions, contingency tables, Chi-square tests, linear probability models, and logit models. I also show how an instructor might use the dataset to introduce students to experimental treatments to identify causal effects. The examples I give are taken from my teaching of a business statistics class for MBA students, one which emphasizes data literacy. However, most would be relevant in econometrics classes and other statistics classes.

The learning outcomes emphasized in my exposition are students' abilities to apply each statistical method and to interpret the output each method generates. For example, I show how one can use the dataset to illustrate the use of contingency tables to explore the association between a worker's knowledge of peer salaries and their level of job satisfaction. I also suggest ways one can help students interpret particular relative frequencies from a contingency table.

One reason this dataset is valuable is that it allows an instructor to connect their class to the real world and real research. Singer and Willett argue for using real data to teach statistics (1992). They maintain that students are interested in learning about the world they live in, so that data that are from that world motivate them. Becker and Greene suggest that instructors avoid "contrived situations with made-up data" and place more emphasis on the application of econometrics to answer real questions (2001). Neumann, Hood, and Neumann conducted interviews of undergraduate students in a first-year statistics class, asking them about real data sets that had been used teaching the class; 58% thought real data made learning statistics more interesting and enjoyable (2013). Using real data in class can increase student interest. It can also help prepare students to apply what they learn, by showing how methods have been used to answer research questions.

---

A second reason this dataset is good for teaching is that students care about income inequality. Real data increases student interest most when it connects to their lives in a substantial way (Willett and Singer, 1992). Income inequality has such a connection; income inequality affects the lives of students, their friends, and their families. In most developed economies, income inequality has been growing since the 1970s. Increases in the US have been especially large. American awareness of this increase is broad; for example, a 2012 Pew survey found that 68 % of college-age Americans knew income inequality had risen in the previous ten years.[1] The dataset presented here allows an instructor to take advantage of interest in inequality by helping students explore its consequences.

A third reason this dataset is useful is that it provides information on employees of a particular organization, including data on job satisfaction. It allows an instructor to draw valuable connections between an econometrics class and other classes students take, for example in labor economics, or business administration, or psychology. There are many good datasets available on the internet to help instructors teach econometrics and statistics, but very few provide information on an organization's employees. In addition, job satisfaction is an especially important employee characteristic. Increases in job satisfaction increase workers' well-being (Bowling et al., 2010), reduce job turnover (Wright and Bonett, 2007), and may increase organizational productivity (Böckerman and Ilmakunnas, 2012).

One alternative to the Card, et al. dataset is one containing responses to a survey evaluating management quality in an Australian educational institution, available on the website "SPSS Data Files and Exercises, 2019."[2] It provides responses from 536 employees to 10 survey items evaluating how well they think they are managed, their job satisfaction, and questions about age, length of service, and city. The Australian dataset would be a good choice for instructors wishing to expose students to a survey measuring different dimensions of good management, but it is of limited use for exploring job satisfaction: the satisfaction question is yes-no and only 10 % of respondents were not satisfied. The Card, et al. dataset also has the advantage of being well documented; students can learn the particulars of who created it, for what purpose, and when. The Australian dataset lacks this information.

The remainder of the paper is organized as follows. I introduce the dataset and the research that generated it in Section 2. Section 3 illustrates how the dataset might be used to teach six statistical methods frequently included in introductory classes, while also providing teaching tips and identifying potential pitfalls. In Section 4, I summarize, discuss assessment, and suggest ways the dataset might be used in advanced classes.

## 2. The research and the dataset

One important dimension of income inequality is earnings inequality within occupations. For example, Mouw and Kalleberg use 496 3-digit occupations to estimate changes in wage inequality from 1983 to 2008; by the end of their period, 57% of the variance of the log wage was within occupations, while 43% was between occupations (2010). How does this within-occupation inequality affect job satisfaction? For social scientists trying to understand human motivation and for managers crafting workplace policy, this is an important question. The dataset introduced here was used in seminal research by David Card, et al.; it provided evidence to help answer this question (2012).

In 2008, the *Sacramento Bee* established a website providing access to a database containing the pay for all state of California employees.[3] A few months later, the researchers sent an email to a random sample of employees at three University of California campuses; the email described the website and provided a link to it (Card et al., 2012).[4] Several days after this experimental treatment, the researchers surveyed all the employees at the campuses, asking about their job satisfaction and whether they had visited the site. Subsequently, the researchers calculated the position of each respondent in their unit's salary distribution. The resulting dataset allowed the researchers to estimate the effect on an employee's job satisfaction of knowing his/her relative salary.

The dataset provides a variety of relevant variables. The ones I discuss in this article are a variable indicating if the employee received the treatment email, a variable indicating if the employee visited the *Sacramento Bee* website, measures of employee job satisfaction, and a variable indicating whether an employee is in the bottom quartile of their unit's salary distribution. Table 1 lists these variables and provides descriptions. The dataset provides twenty-one additional variables, including indicator variables measuring campus of employment, additional measures of the employee's position in their unit's salary distribution, and, for employees who visited the *Bee* site, indicators of whose salaries they checked on their visit.

The text of the employee survey is available at the beginning of the Online Appendix published with the article. The survey response rate was 20.4 %, with 6411 total responses. The dataset itself is available on the article's supplementary materials page as a compressed Stata file. An abbreviated version of the dataset, including only the variables discussed in this article, is available as an Excel Workbook in Mendeley Data.

---

[1] The survey was conducted by the Pew Research Center (2013) in July 2012. The reported percentage is the share of people 19-22 in the sample choosing the "gotten larger" option in Question 13: "In the past ten years, do you happen to know if the income gap between the rich and the poor has [RANDOMIZE(gotten larger), (gotten smaller)] or has it stayed about the same?"

[2] The dataset's authors are not identified. It is available on a page of a website providing SPSS data files for *The SPSS Survival Manual*, by Julie Pallant.

[3] The link is to the current version of the Sacramento Bee site (2017).

[4] All information presented here about this research and the data it was based on are taken from Card et al. (2012). Additional detail is available in Card et al. (2010).

**Table 1**

Variables Discussed.

| Variable | Description |
| --- | --- |
| TREATED_IND | Equals 1 if the employee received the treatment email, 0 otherwise |
| SAC_BEE | Equals 1 if the employee visited the *Sacramento Bee* website, 0 otherwise. |
| WAGESAT2 | 4-level Likert item, the answer to: "How satisfied are you with your wage/salary on this job?" |
| JOBSAT2 | 4-level Likert item, the answer to: "All in all, how satisfied are you with your job?" |
| MOVE2 | 3-level Likert item, the answer to: ""Taking everything into consideration, how likely is it you will make a genuine effort to find a new job within the next year?" |
| LOW25 | Equals 1 if the employee is in the bottom quartile of his/her unit's salary distribution (calculated separately for faculty and staff), 0 otherwise |
| MOVE_VLIK | Equals 1 if MOVE2 equals 3 (Very Likely), 0 otherwise |

## 3. Teaching with the dataset

I use the dataset for in-class discussion, in-class exercises, and homework assignments. In what follows, I present questions an instructor can pose, techniques that can be used to answer those questions, and possible student misunderstandings.

### 3.1. Relative frequency distributions

One can calculate a great variety of interesting descriptive statistics with the dataset. I focus on relative frequency distributions, since understanding them well helps students grasp probability distributions when I introduce them.

The data were collected to answer questions about job satisfaction. Four variables in the dataset measure dimensions of satisfaction; I discuss the three that seem most strongly linked to how content an employee feels at work. The variable WAGESAT2 is a 4-level Likert item created from answers to this question: "How satisfied are you with your wage/salary on this job?" The variable JOBSAT2 is another 4-level Likert item; it was created with responses to this question: "All in all, how satisfied are you with your job?" A third variable, MOVE2, is based on 3-level Likert responses to this question: "Taking everything into consideration, how likely is it you will make a genuine effort to find a new job within the next year?"

*Teaching Tip:* Since the first two measures share the same response options, it is easy to ask students to compare them. For example, after I have introduced relative frequency distributions and the features of informative, easy-to-grasp graphs in class, on a homework assignment I ask students to create a good graph comparing the distributions of WAGESAT2 and JOBSAT2. As a follow-up, I ask, "Please write a description of an important relationship between the two variables, one that is apparent in your graph. Use less than 60 words and include two relative frequencies taken from the table of data you used to create your graph." I include short writing assignments like this one on homework to help students think about concepts and learn to communicate the results of statistical analysis. Making the question I ask specific and including a word limit eases my evaluation of student answers.

A good resulting graph might be like Fig. 1.

A good answer might be something like the following: "Based on the survey responses, employees seem more satisfied with their jobs than with their wages. For example, 37 % of employees say they are very satisfied with their jobs, but only 12 % say they are very satisfied with their wages."

After we discuss answers to the follow-up question in class, I try to get students to think more about it. After presenting the answer above, I might ask, "The proportion of employees who say they are satisfied with their job is quite similar to the proportion who say they are satisfied with their wages. Should we conclude from this that employees who are satisfied with their job are also likely to be
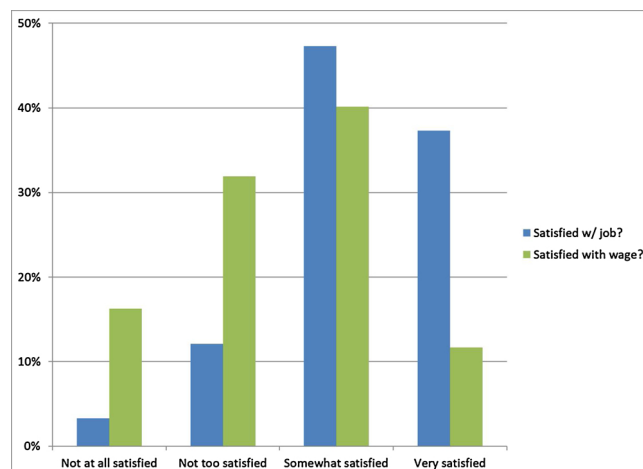


**Fig. 1.** Comparing Two Job Satisfaction Measures.

**Table 2**
Job Satisfaction and Wage Satisfaction.

| Count Satisfied with salary? | Satisfied with job? Not at all satisfied | Not too satisfied | Satisfied | Very satisfied | Grand Total |
|---|---|---|---|---|---|
| Not at all satisfied | 135 | 274 | 458 | 177 | 1044 |
| Not too satisfied | 41 | 324 | 1165 | 515 | 2045 |
| Satisfied with salary | 26 | 165 | 1260 | 1121 | 2572 |
| Very satisfied | 9 | 13 | 140 | 578 | 750 |
| **Grand Total** | **211** | **776** | **3033** | **2391** | **6411** |

satisfied with their wages? Why or why not?" A question like this gets students thinking about the joint distribution of the two variables, something which can help them understand contingency tables more quickly later.

Especially in a big class, it is worth including a question like this on a PowerPoint slide. Hearing the question and then being able to refer back to it saves class time, because students ask fewer questions of clarification. A visual version of the question also increases the chance that shy students, ones who will not ask a question, understand it. Finally, a visual version creates a stronger expectation that students actually come up with an answer; it makes it less likely they consider the question rhetorical.

### 3.2. Contingency tables

Once we have described variables individually, it is time to describe them jointly. Since we have already compared JOBSAT2 and WAGESAT2, a contingency table like Table 2 is a good place to begin. After reminding them how a contingency table is created, I ask students to interpret two numbers from the "Satisfied" column of the table: 1260 and 3033. After we discuss the difference between interior frequencies and marginal frequencies, I point out that those two numbers allow one to see that only 42 % of the employees who are satisfied with their job are also satisfied with their wages.

Having described the joint distribution of two job satisfaction measures, it is natural to ask what factors might influence job satisfaction. Card et al. (2012) investigated the link between an employee knowing coworkers' earnings and that employee's job satisfaction. Another contingency table can help students see this link. Having had students interpret counts from a contingency table, I begin their exploration by asking them to create one.

The dataset's only comprehensive measure of what a respondent knows about peer salaries is the SAC_BEE variable; it measures whether an employee has visited the *Sacramento Bee* website.[5] Students could create a contingency table classifying various of the dataset's measures of job satisfaction by SAC_BEE. However, the variable called MOVE2, which measures the intent-to-seek-a-job, produces the strongest associations, so I focus students' attention on it. I ask them to create a contingency table displaying counts, one that describes the relationship between the values of the SAC_BEE variable and the MOVE2 variable. To make it easier to compare the tables students produce, I ask them to put MOVE2 in the columns and SAC_BEE in the rows. Students produce a table that looks something like Table 3.

At this point, I mention that Table 3 poses a problem: far fewer employees visited the website than did not. This makes it hard to identify any relationship that might exist between visits to the website and intentions to seek a new job. To illustrate a solution to this problem, I present two new descriptions of the same data, the ones in Table 4. The first presents column percentages and the second row percentages.

*Potential Pitfall:* Students frequently confuse column percentages with row percentages, and *vice versa*. To help them understand the distinction, I have them explain a percentage from Table 4a and the corresponding percentage from Table 4b. For example, I might first ask them to interpret the third percentage in the second column of Table 4a: 26.1 %. A good answer might read something like this, "The 26.1 % in Table 4a is the percentage of employees visiting the website who said they were very likely to seek a new job." Then, I might ask the students to interpret the third percentage in the second column of Table 4b. Many students who struggle with the first question will answer the second well, since they now have a model to help structure their thinking. A possible good answer might be: "The 32.4 % in Table 4b is the percentage of employees saying they were very likely to seek a new job who visited the website."

Next, I get students thinking about the positive association between website visits and the expressed likelihood of seeking a new job. Why might the proportion of those very likely to seek a new job be lower for those who did not visit the website (20.8 %) than for those who did (26.1 %)? Why might the proportion of those visiting the website rise as one goes from those who answered, "not at all likely" to seek a new job (25.4 %) to those who answered, "very likely (32.4 %)?" Asking questions like these helps students see the connection between the contingency table and the sorts of questions researchers pose. Such questions engage some students' imagination and help students connect statistical methods with social science research.

Students come up with various explanations for these patterns; those explanations usually include some version of the following two. One, employees who are likely to seek a new job may be more likely to visit the website because they want more information about what they might earn in other jobs. Two, learning coworkers' salaries may affect an employee's job satisfaction. An employee

---

[5] There are 4 variables giving more specific information about what salaries the respondent checked when visiting the website, but the survey questions obtaining the data for those variables were only asked of respondents at UCLA.

**Table 3**

Visits to Sacramento Bee Website and Job Serach Intentions, Counts.

| Count<br>Seek new job? | Visit SB website<br>No | Yes | Grand Total |
|---|---|---|---|
| Not at all likely | 2246 | 765 | 3011 |
| Somewhat likely | 1433 | 541 | 1974 |
| Very likely | 964 | 462 | 1426 |
| **Grand Total** | **4643** | **1768** | **6411** |

**Table 4**

Visits to Sacramento Bee Website and Job Search Intentions, Relative frequencies.

a Column percentages

| Column %<br>Seek new job? | Visit SB website<br>No | Yes | Grand Total |
|---|---|---|---|
| Not at all likely | 48.4 % | 43.3 % | 47.0 % |
| Somewhat likely | 30.9 % | 30.6 % | 30.8 % |
| Very likely | 20.8 % | 26.1 % | 22.2 % |
| **Grand Total** | 100.0 % | 100.0 % | 100.0 % |

b Row percentages

| Row %<br>Seek new job? | Visit SB website<br>No | Yes | Grand Total |
|---|---|---|---|
| Not at all likely | 74.6 % | 25.4 % | 100.0 % |
| Somewhat likely | 72.6 % | 27.4 % | 100.0 % |
| Very likely | 67.6 % | 32.4 % | 100.0 % |
| **Grand Total** | 72.4 % | 27.6 % | 100.0 % |

who learns his/her pay is relatively low may become less satisfied and therefore more likely to leave. An employee whose pay is relatively high may become more satisfied and therefore less likely to leave.[6]

These two explanations are very different. The first says the intention to leave causes the website visit. The second says the website visit causes a change in job satisfaction, which affects the intention to leave. How could one distinguish between them?

*3.3. Experimental treatments*

Experimental treatment is a strategy to identify a causal effect. This dataset allows an instructor to introduce students to the use of experimental methods to identify causal effects, methods that have become more common in social sciences during the last twenty years. Card and his coauthors used a novel treatment to produce evidence about which of the two explanations mentioned above is correct (2012).

The treatment was the email researchers sent to a random sample of University of California employees. The email greatly increased the percentage of employees visiting the website, from 19 % for untreated employees to 49 % for treated ones. By providing easy access to the website and the information it contained, the email also reduced the influence of employee characteristics on website visits.

*Teaching Tip:* One way to illustrate the effect of the experimental treatment is to use a contingency table. One possible table is like Table 4a, but classifies respondents by whether they were treated, rather than by whether they visited the website.[7] Table 5b does this, while Table 5a repeats Table 4a. I show students these two tables together and suggest comparing them might lead to the conclusion that employee characteristics affected website visits. I ask them what in the tables might lead to this conclusion.

To answer the question: a comparison of the tables shows the differences between untreated and treated individuals are smaller than the differences between individuals who did not visit the website and individuals who did. As an example, take the percentages

---

[6] Of course, an instructor would be thrilled to hear a third explanation for these differences between proportions: sampling error. However, in this particular comparison the differences are unlikely to be due to sampling error. For example, if we evaluate the pairwise comparisons going down the columns in Table 4a, using a z-test with a 5% significance level and a Bonferroni correction, the difference between 20.8% and 26.1% is statistically significant.

[7] Here and throughout my discussion of experimental treatment, I focus attention on the treatment assigned, rather than the treatment received. Individuals receiving the email may choose to get the treatment (visit the website) or not. As Freedman (2008) points out, focusing analysis on the treatment-assigned group avoids potential bias. Employees who get the email and visit the website may differ in important ways from employees who get the email choose not to visit. Card and his coauthors argue that measuring the change in intent to leave from the untreated to the treated estimates the minimum average treatment effect.

**Table 5**
Website Visits, Email Treatment, and Job Search Intentions.

| Table 5a<br>Column %<br>Seek new job? | Visit SB website<br>No | Yes | Grand Total |
|---|---|---|---|
| Not at all likely | 48.4 % | 43.3 % | 47.0 % |
| Somewhat likely | 30.9 % | 30.6 % | 30.8 % |
| Very likely | 20.8 % | 26.1 % | 22.2 % |
| **Grand Total** | 100.0 % | 100.0 % | 100.0 % |
| Table 5b<br>Column %<br>Seek new job? | Treated individual<br>No | Yes | Grand Total |
| Not at all likely | 47.4 % | 45.8 % | 47.0 % |
| Somewhat likely | 30.7 % | 31.1 % | 30.8 % |
| Very likely | 21.9 % | 23.1 % | 22.2 % |
| **Grand Total** | 100.0 % | 100.0 % | 100.0 % |

who say they are very likely to seek a new job. In Table 5b, that percentage goes from 21.9 % for the untreated to 23.1 % for the treated. In Table 5a, it goes from 20.8 % for those who did not visit the website to 26.1 % for those who did. The fact that the first differences are smaller than the second suggests that differences in employee characteristics influenced visits to the *Sacramento Bee* website and the experimental treatment reduced the influence of this heterogeneity[8]

However, Table 5 groups all employees together. By doing that, it may cloak an effect of website visits. That could happen, for example, if an employee's position in his/her unit's salary distribution influences the effect of the visit on job satisfaction. For example, an economics professor with relatively low pay might become more likely to seek a new job, while the reverse might happen with an economics professor with relatively high pay. Table 6 allows us to investigate this possibility, by comparing employees in the bottom 25 % of their unit's salary distribution to employees in the top 75 % of their unit's salary distribution.

*Teaching Tip:* I show Table 6 to students and ask them to fix their attention on the last row. I ask them what evidence it provides about the effect of website visits on job satisfaction. The row shows that the treated are more likely to seek a new job, since the proportion saying they are very likely to seek a new job increases by 7.6 percentage points among the bottom quartile employees but decreases by 0.7 percentage points among employees in the top three quartiles. Hopefully, a student suggests that website visits seem to have decreased satisfaction among relatively low earning employees.

Once we have discussed the difference between Table 6a and Table 6b, I ask students about the evidence in Table 6a itself. What could have caused this decrease in satisfaction?

Typically, they suggest it might have resulted from employees becoming more aware that their pay is relatively low.

Finally, I ask students why there might be such a difference between the lowest-quartile group and the rest of employees. If low relative pay decreases satisfaction, why does not high relative pay increase satisfaction? This is a challenging question, one that research has not answered conclusively. After a discussion of possible reasons, I report that Card et al. argue that the result is evidence for inequality aversion; job satisfaction is influenced by relative pay, but satisfaction "is a concave function of relative pay" (2012, p. 2982).

### 3.4. Chi-square test

Once the course advances to inferential statistics, Table 6 allows an instructor to raise another issue with the dataset: sampling error. If the experiment were repeated, different employees would receive the email, very likely changing the sample proportions observed. Could the 7.6 percentage point difference discussed above result merely from sampling error?

Of course, one way to answer that question is to think about the difference between the proportions displayed in Table 6a and the population proportions they estimate. For example, we could focus our attention on the proportion of employees who said they were very likely to seek a new job. We could test the hypothesis that there is no difference between that proportion for employees who did not receive the email (the controls) and employees who did receive the email (the treated). More formally, our null and alternative hypothesis would be:

$H_0$: $\pi_C = \pi_T$

$H_1$: $\pi_C \neq \pi_T$

A simple way to implement this test is to create a table combining workers who are "somewhat likely" and "not at all likely" into one group, to create a binary measure of job satisfaction. Table 7 does this, by using a variable MOVE_VLIK which is equal to one for employees who report being very likely to move and zero otherwise. Table 7 mirrors the structure of Table 6, but presents counts, rather than relative frequencies, and adds statistics from the hypothesis test in the bottom half.

---

[8] To evaluate this hypothesis directly, we would like to know which employees visited the website before the treatment and which visited it after. The dataset does not provide this information.

**Table 6**
Email treament and likelihood of seeking new job.

a Employees in bottom 25 %

| Column % | Treated individual | | |
|---|---|---|---|
| Seek new job? | No | Yes | Yes-No |
| Not at all likely | 41.0 % | 35.1 % | −5.9 pp |
| Somewhat likely | 33.9 % | 32.2 % | −1.7 pp |
| Very likely | 25.1 % | 32.7 % | 7.6 pp |
| **Grand Total** | 100.0 % | 100.0 % | |

b Employees in upper 75 %

| Column % | Treated individual | | |
|---|---|---|---|
| Seek new job? | No | Yes | Yes-No |
| Not at all likely | 49.4 % | 49.0 % | −.5 pp |
| Somewhat likely | 29.6 % | 30.8 % | 1.2 pp |
| Very likely | 20.9 % | 20.2 % | −.7 pp |
| **Grand Total** | 100.0 % | 100.0 % | |

**Table 7**
Email treatment and likelihood of seeking new job, Chi-square.

| | Table 7a, Employees in bottom 25 % | | |
|---|---|---|---|
| *Seek new job?* | *Treated individual* | | *Total* |
| | No | Yes | |
| Not very likely | 829 | 276 | 1105 |
| Very likely | 278 | 134 | 412 |
| Total | 1107 | 410 | 1517 |
| ` | Value | df | Asymp. Sig. (2-sided) |
| Pearson Chi-square | 8.667[b] | 1 | .003 |
| Valid Cases | 1517 | | |

| | Table 7b, Employees in upper 75 % | | |
|---|---|---|---|
| *Seek new job?* | *Treated individual* | | *Total* |
| | No | Yes | |
| Not very likely | 2790 | 1090 | 3880 |
| Very likely | 738 | 276 | 1014 |
| Total | 3528 | 1366 | 4894 |
| ` | Value | df | Asymp. Sig. (2-sided) |
| Pearson Chi-square | .305[b] | 1 | .581 |
| Valid Cases | 4894 | | |

Once students have learned the Chi-square test for two proportions, I can show them Table 7 and ask them whether they would reject the null hypothesis for the bottom-quartile group and for the top-three-quartiles group. At the 1% level of significance, the data show the first answer is yes, since the *p*-value is .003, and the second answer is no, since the *p*-value is .581. These results support our earlier conclusion (using Table 6) that website visits reduced job satisfaction among employees in the low-earning group.

### 3.5. Linear probability models

Though contingency tables measure the consequences of experimental treatment, statisticians often estimate treatment effects with techniques that make it easier to control for non-treatment influences on the dependent variable. A linear probability model is one possibility, though such models do have defects. Horrace and Oaxaca show they frequently produce biased estimates and nearly always are inconsistent (2006). Woolridge points out such models always have heteroskedastic errors (2013).

Even so, linear probability models are a good way to introduce students to models that estimate treatment effects. They often make accurate predictions for central values of independent variables and their estimated coefficients have straightforward

**Table 8**

Predicting MOVE_VLIK with treatment.

| Variables | B | Std. Error | t | Sig. |
|---|---|---|---|---|
| (Constant) | .219 | .006 | 35.88 | .000 |
| TREATED_IND | .012 | .012 | 1.00 | .315 |

interpretations (Woolridge, 2013). Moreover, once a student understands linear probability models, they will find it easier to learn other models for predicting binary dependent variables. For example, think about the logit model. It will be easier for a student to understand that a coefficient from such a model estimates the effect of the independent variable on the log of the odds of an event occurring if the student first understands the coefficient of a linear probability model estimates its effect on the probability of an event.

Once students have studied multiple linear regression, dummy variables, and interaction terms, an instructor can illustrate their use with the dataset. One possibility is to focus the class's attention, as Table 7 did, on the MOVE_VLIK variable, the one that equals one if a respondent reports being very likely to search for a new job and zero otherwise. Using it as the dependent variable, the instructor can predict the probability that an employee reports being very likely to search for a new job. That model would look like this:

$$MOVE\hat{}\_VLIK = b_0 + b_1 TREATED\_IND + b_2 LOW25 + b_3 TREATED\_IND*LOW25,$$

where:

$MOVE\hat{}\_VLIK$ is the predicted value of the variable measuring whether the employee is very likely to search for a new job,
$TREATED\_IND$ is the variable measuring whether the employee received the treatment email,
$LOW25$ is the variable measuring whether the employee is in the bottom 25 % of their unit's salary distribution, and
$TREATED\_IND * LOW25$ is an interaction term.

*Potential pitfall:* Many students find interaction terms difficult to interpret. As a result, it might be best for an instructor to teach two preliminary models before interpreting the complete model represented by Equation 1. I illustrate that approach here.

The first preliminary model predicts MOVE_VLIK with only the first independent variable, the one measuring whether the employee received the treatment email. That estimation is presented in Table 8.

The instructor could present this estimation and ask about the interpretation of the constant and the slope coefficient, along with their significance levels. The value of the slope coefficient for the variable measuring whether the employee received the treatment nearly matches the change in relative frequencies in the "Very Likely" row of Table 5b.

The second preliminary model predicts MOVE_VLIK with only the second independent variable, the one measuring whether the employee was in the bottom 25 % of their unit's salary distribution. The estimation of that model is reported in Table 9.

In this case, having already discussed the constant term, an instructor might focus students' attention on the interpretation of the slope coefficient and its significance. In the subsequent discussion, it would be good to distinguish the description of the sample relationship provided here from the one provided in Table 6 and Table 7. Thus far, the distribution of the salary level variable has only been examined in tandem with two other variables, the ones measuring treatment and job satisfaction. Here, in contrast, we ignore treatment; the slope coefficient measures solely the average "effect" on job satisfaction of being in the bottom 25 %.[9]

After discussing this coefficient, an instructor might also ask students to provide possible explanations for it being positive and statistically significant. Why might employees in the bottom 25 % more often say they are very likely to search for a new job? Students should recognize various explanations are possible and that the coefficient, even though it is statistically significant, does not provide strong evidence for relatively low pay increasing interest in alternative employment. For example, the causality might run the other way. Low job satisfaction may lower productivity; lower productivity could cause lower pay.

Having helped students think about the independent variables measuring treatment and salary distribution position, the instructor has prepared them to interpret a variable interacting the two, so he/she can present the complete model. The results presented in Table 10 are from the estimation of a model predicting the probability an employee is very likely to search with the variable measuring treatment, the variable measuring position in the salary distribution, and a variable interacting treatment and position.

This specification is very close to the one Card et al. prefer (2012).[10] It provides evidence that, if you are in the bottom 25 %, discovering your position raises the chance you will say you are likely to search for a new job. The point estimate for the effect is 7.6 percentage points; that is, the coefficients in Table 10 imply that receiving the email increased the proportion of employees saying they are very likely to search for a new job by 30 %[11]

---

[9] This coefficient also presents an opportunity to strengthen students' intuition about the relationship between the linear probability models and contingency tables. An instructor could ask: "The coefficient on the LOW25 variable corresponds to the difference between two column proportions in the cross-tabulations presented in Table 6. Which two? (It corresponds to the difference between the lower-left cell of Table 6a [Untreated & Very Likely, 25.1%] and the lower-left cell of Table 6b [Untreated & Very Likely, 20.9%]: .251 - .209 = .042.)

[10] It is Specification (6) in Table 4 on p. 2994.

[11] In Fig. 2, this 7.6 percentage-point change corresponds to the difference between the predicted value for Individual 3 and the predicted value for Individual 2 (.327 - .251 = .076). The proportional rise in probability reported is calculated as follows: .076/.251 = .303.

**Table 9**

Predicting MOVE_VLIK with position in salary distribution.

| Variables | B | Std. Error | t | Sig. |
|---|---|---|---|---|
| (Constant) | .207 | .006 | 34.92 | .000 |
| LOW25 | .064 | .012 | 5.28 | .000 |

**Table 10**

Predicting MOVE_VLIK with treatment and position in salary distribution.

| R Square | Adjusted R Square | | | SEE |
|---|---|---|---|---|
| .006 | .005 | | | .415 |

| Variables | B | Std. Error | t | Sig. |
|---|---|---|---|---|
| (Constant) | .209 | .007 | 29.95 | .000 |
| TREATED_IND | −.007 | .013 | −.54 | .589 |
| LOW25 | .042 | .014 | 2.94 | .003 |
| TREATED_IND*LOW25 | .083 | .027 | 3.03 | .002 |

After discussing the coefficients and their levels of statistical significance, the instructor might show students how to calculate a predicted value, for example the predicted probability for an employee in the top 75 % of their unit's salary distribution who has not received an email (Individual 1 in Fig. 2).

*Teaching tip:* Because interaction terms are difficult for many students to understand, I like to give students time to think about them with a partner. To do that, I use a think-pair-share exercise like the handout illustrated in Fig. 2. I distribute the handout (with the probabilities for Individual 1 blank), before showing students how to calculate the predicted probability for Individual 1. Having a handout to mark up as I lecture helps students grasp the calculation. Including four specific questions on the handout allows me to structure partners' work, to make it more productive. Having two paragraphs explaining the calculation of the predicted values gives students a resource, besides me, to help them answer the handout's questions. Especially in a large class, anticipating questions in this way can substantially reduce "dead time" students spend waiting for the instructor to answer a question.

### 3.6. Logit models

In case an instructor wishes to discuss a logit model, Table 11 estimates Equation 1) using logistic regression. The results are very similar to the linear probability model. The R-squared surrogates in Table 11 are very close to the R-squared in Table 10.[12] Measures of statistical significance are also quite close between the two tables. Predicted probabilities for the four groups (represented by the four individuals in Fig. 2) are nearly identical.

### 4. Closing discussion

This article introduces a valuable dataset, one that allows an instructor to link the study of statistics to students' interest in income inequality. The data were generated by Card et al. to study the effect of occupational earnings inequality on job satisfaction (2012). The dataset describes the responses to a survey of staff and faculty at three University of California campuses.

Besides addressing a topic of interest to students, the dataset allows an instructor to create relevant and clear examples of the use of statistical tools included in many introductory classes. This article shows how an instructor might illustrate the use of frequency distributions, contingency tables, Chi-square tests, linear probability models, and logit models. Table 12 provides additional reminders of the topics discussed. The dataset might be particularly valuable for classes including the study of experimental treatments, since the values of a crucial variable were generated using a creative experiment, one that might open students' eyes to the potential for low-cost trials to provide important information.

The illustrations included in the article highlight teaching students to a) apply statistical methods and b) interpret the output those methods generate. The article also illustrates how one might assess student learning of interpretation skills on the fly, by asking questions of the class and having students do in-class exercises. The questions I suggest can also be used in problem sets and exams to assess interpretation skills more formally. Application skills are more challenging to assess. I do that using homework assignments and an end-of-semester data analysis project.

The dataset could also be useful in more advanced classes. For the study of experimental treatments, the data could be used to illustrate the implications of heterogeneous treatment effects; they also could be used to illustrate the difference between average

---

[12] One can also evaluate goodness of fit using a classification rule. One plausible rule predicts a value of 1 (the employee will report they are very likely to search for a new job) if the predicted probability is .251 or larger, and otherwise predicts a value of 0. For that rule, the sensitivity is 27.2% and the specificity is 79.3%.

| | | TREATED_IND | LOW25 | TREATED*LOW25 | Constant | |
|---|---|---|---|---|---|---|
| | **Coefficients** | **-.007** | **.042** | **.083** | **.209** | **Predicted** |
| | **1** | 0 | 0 | 0 | | **Values** |
| | | 0.000 | 0.000 | 0.000 | 0.209 | **0.209** |
| | | | | | | |
| | **2** | 0 | 1 | 0 | | |
| | | 0.000 | 0.042 | 0.000 | 0.209 | **0.251** |
| **Individuals** | | | | | | |
| | **3** | 1 | 1 | 1 | | |
| | | -0.007 | 0.042 | 0.083 | 0.209 | **0.327** |
| | | | | | | |
| | **4** | 1 | 0 | 0 | | |
| | | -0.007 | 0.000 | 0.000 | 0.209 | **0.202** |

The table above describes four employees (Individual 1, 2, 3, and 4) using the values of two independent variables: TREATED_IND and LOW25. They tell us if a person received an email and if they have "low earnings." The third variable in the table, TREATED*LOW25, is the product of the first two. For example, Individual 1 is described by TREATED_IND=0 and LOW25=0; she is an employee who did not receive the email and who is not in the bottom 25% of her unit's salary distribution. In turn, TREATED*LOW25=0 because $0 \times 0 = 0$.

To predict the likelihood of each individual saying he/she wants to look for a new job, the table takes two steps:
    1) It multiplies each independent variable by its associated coefficient. (The resulting products are listed just below the values of the independent variables.)
    2) It adds the products to the constant. (These sums are listed in the final column, labeled "Predicted Values")

Please answer the following four questions about the table:
    a) What does the number .209 represent?
    b) How would you describe Individual 2?
    c) How are Individual 2 and Individual 3 the same? How are they different?
    d) Going from Individual 2 to Individual 3, the predicted value rises by 7.6 percentage points (.327-.251=.076). In the context of this model, what causes this rise?

**Fig. 2.** Think-pair-share exercise on how to interpret the linear probability model.

**Table 11**
Predicting MOVE_VLIK using logistic regression.

| −2 Log likelihood | | Cox & Snell R Square | | | Nagelkerke R Square | |
|---|---|---|---|---|---|---|
| 6759.6 | | .006 | | | .008 | |

| Coefficients | B | S.E. | Wald | df | Sig. | Exp(B) |
|---|---|---|---|---|---|---|
| (Constant) | − 1.330 | .041 | 1032.14 | 1 | .000 | .265 |
| TREATED_IND | − .044 | .079 | .31 | 1 | .581 | .957 |
| LOW25 | .237 | .081 | 8.64 | 1 | .003 | 1.268 |
| TREATED*LOW25 | .414 | .149 | 7.73 | 1 | .005 | 1.512 |

**Table 12**

Methods and topics the dataset is suited to illustrate.

| Methods | Topics |
| --- | --- |
| Frequency distributions | Relative frequency histograms |
| Contingency tables | Absolute frequency distributions, relative frequency distributions (including the interpretation of row and column percentages) |
| Chi-square tests | The connection between relative frequency distributions and counts in the Chi-square table, interpreting *p*-values |
| Linear probability models | Interpreting coefficients (including those on interaction terms), understanding the connections between coefficients and relative frequencies in contingency tables |
| Logit models | Shortcomings of linear probability models, interpretation of logit coefficients |

treatment effects and local average treatment effects.[13] Classes studying latent variables could utilize the four measures of job satisfaction the dataset provides. Research methods classes could assess the authors' use of a placebo treatment to evaluate the effect of the treatment email.

**Sole author**

I am the sole author and did all the necessary work for this paper.

**References**

Becker, W.E., Greene, W.H., 2001. Teaching statistics and econometrics to undergraduates. J. Econ. Perspect. 15 (4), 169–182.

Böckerman, P., Ilmakunnas, P., 2012. The job satisfaction-productivity nexus: a study using matched survey and register data. Ind. Labor Relat. Rev. 65 (2), 244–262.

Bowling, N.A., Eschleman, K.J., Wang, Q., 2010. A meta-analytic examination of the relationship between job satisfaction and subjective well-being. J. Occup. Organ. Psychol. 83 (4), 915–934.

Card, D., Mas, A., Moretti, E., Saez, E., 2010. Inequality at Work: the Effect of Peer Salaries on Job Satisfaction. National Bureau of Economic Research Working Paper 16396 Available at https://escholarship.org/uc/item/48z7z9dn. Accessed on August 1, 2014.

Card, D., Mas, A., Moretti, E., Saez, E., 2012. Inequality at work: the effect of peer salaries on job satisfaction. Am. Econ. Rev. 102, 2981–3003.

Freedman, D., 2008. Randomization does not justify logistic regression. Stat. Sci. 23, 237–249.

Horrace, W.C., Oaxaca, R.L., 2006. Results on the bias and inconsistency of ordinary least squares for the linear probability model. Econ. Lett. 90 (3), 321–327.

Mouw, T., Kalleberg, A.L., 2010. Occupations and the structure of wage inequality in the United States, 1980s to 2000s. Am. Sociol. Rev. 75 (3), 402–431.

Neumann, David L., Hood, Michelle, Neumann, Michelle M., 2013. Using real-life data when teaching statistics: student perceptions of this strategy in an introductory statistics course. Statistics Education Research Journal 12 (2).

Pew Research Center, 2013. Middle Class II.   Dataset available at http://www.pewsocialtrends.org/dataset/middle-class-ii/. Downloaded June 27, 2017. .

Sacramento Bee, 2017. State Worker Salary Database.   Available at http://www.sacbee.com/site-services/databases/state-pay/article2642161.html. Accessed on June 18, 2018. .

SPSS Data Files and Exercises, 2019. SPSS Data Files and Exercises.  Accessed on May 25, 2019.  http://spss.allenandunwin.com.s3-website-ap-southeast-2.amazonaws.com/data-files.html#.XOmi9YhKiUk.

Willett, J.B., Singer, J.D., 1992. Providing a statistical model': teaching applied statistics using real-world data. Statistics for the twenty-first century 83–98.

Woolridge, J., 2013. Introductory Econometrics: A Modern Approach. South-Western.

Wright, T.A., Bonett, D.G., 2007. Job satisfaction and psychological well-being as nonadditive predictors of workplace turnover. J. Manage. 33 (2), 141–160.

---

[13] The article discusses heterogeneous treatment effects on pages 2990-92 and the distinction between the two types of treatment effects on pages 2994-95.