

Data Exercise 3

Devin Judge-Lord

April 11, 2019

Women at the Deer Valley Utility Company claim that their job performances are not rewarded to the same degree as the job performances of men. Is there statistical evidence to support this complaint? This summary analysis for the Director of the Office of Equal Opportunity includes findings and a brief discussion of other factors we may want to investigate before issuing a final report.

We have the following data for 60 people:

Salary: thousands of dollars.

Gender: "1" for men and "0" for women.

Rating: The employee's average performance rating over the last two years. The scale has a top score of 100. The company claims that performance rating is the primary factor in the determination of salary.

Credits earned either in college courses or company programs.

Retrieve data from *Stata* and load them into **R**.

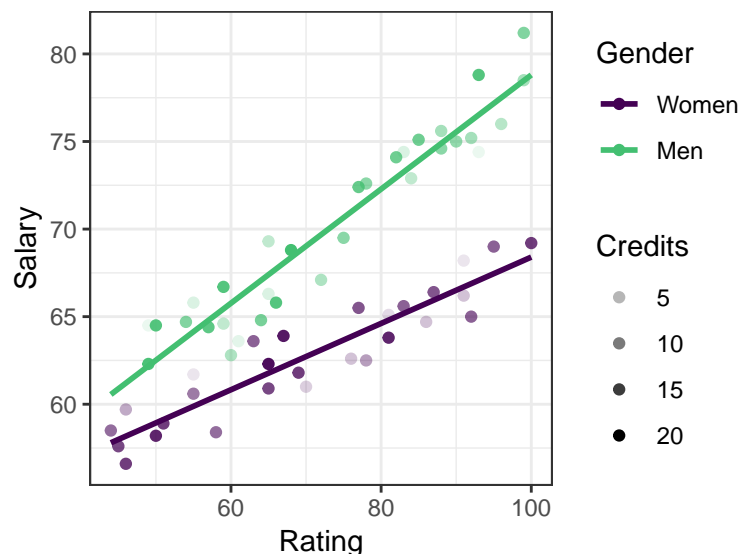
```
net install PS813_EX3, from(https://faculty.polisci.wisc.edu/weimer)
PS813_EX3 1234
save "data/EX3.dta"
```

```
d <- read_dta("data/EX3.dta") %>%
  mutate(Gender = factor(Sex, labels = c("Women", "Men")))
```

First, let us plot the raw data.

```
# scatterplot
p <- ggplot(d) +
  aes(x = Rating, y = Salary, color = Gender) +
  geom_point(aes(alpha = Credits)) + scale_color_discrete()

# quick y ~ mx linear regression per group
p + geom_smooth(method = lm, se = F, fullrange = T)
```



Hypotheses

H1: Job performances of women are rewarded differently than the job performances of men. That is, the relationship between salary and performance differs by gender.

H0: There is no difference in how men's performance and women's performance are rewarded. That is, the relationship between salary and performance does not differ by gender.

(There are least two other ways to write this hypothesis and at least one slightly different hypothesis that might better address the question.)

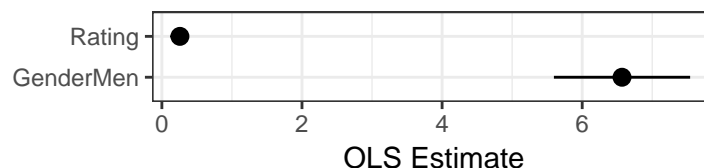
Model

The dependent variable is salary. For employee i , let their salary be y_i in the model $y_i = \beta_0 + \dots + \epsilon_i$. β_0 is the predicted salary, \hat{y} , when all other variables in the model are 0.

Note: The model, $y_i = \beta_0 + \beta_1 * Gender + \beta_2 * Rating + \epsilon_i$, does *not* test the relationship of interest.

```
m1 <- lm(Salary ~ Gender + Rating, data = d) %>%
  tidy(conf.int = TRUE)

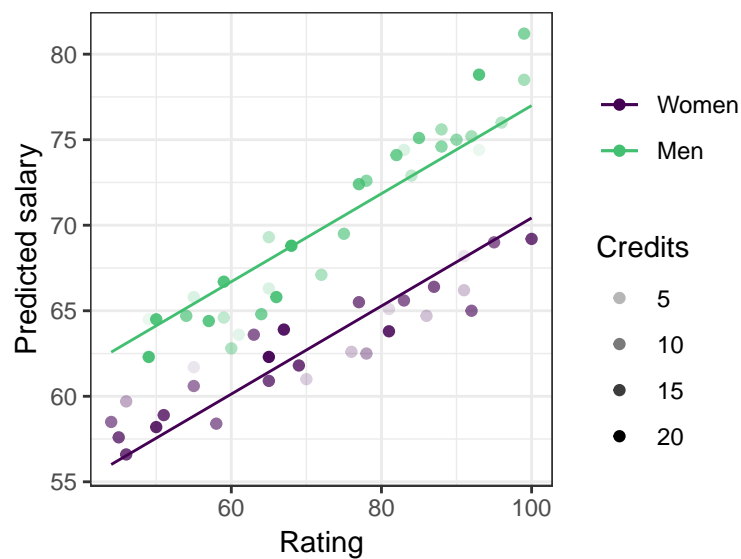
ggplot(m1 %>% filter(term != "(Intercept)")) +
  aes(x = term,
      y = estimate,
      ymin = conf.low,
      ymax = conf.high) +
  geom_pointrange() +
  coord_flip() +
  labs(x="", y="OLS Estimate")
```



Why is this wrong? Let's plot the results against our data!

```
b0 <- m1$estimate[1]
b1 <- m1$estimate[2]
b2 <- m1$estimate[3]

p +
  geom_line(aes(color = "Men",
                 y = b0 + b1*1 + b2*Rating) ) +
  geom_line(aes(color = "Women",
                 y = b0 + b1*0 + b2*Rating) ) +
  scale_color_discrete() +
  labs(color = "",
       y = "Predicted salary")
```



What is wrong with this picture? Why does this model fail to test the hypothesis? What did it test?

Hypothesis test

Lorem ipsum $\beta_? = 0$

Lorem ipsum $\beta_? \neq 0$

Findings

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

Discussion

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.