

Problem Set 4

Applied Stats/Quant Methods 1

Due: December 4, 2022

Question 1: Economics

Using the `prestige` dataset in the `car` library

```
1 data(Prestige)
```

We would like to study whether individuals with higher levels of income have more prestigious jobs. Moreover, we would like to study whether professionals have more prestigious jobs than blue and white collar workers.

- (a) Created a new variable `professional` by recoding the variable `type` so that professionals are coded as 1, and blue and white collar workers are coded as 0.¹

```
1 Prestige[ 'professional' ] <- ifelse(Prestige$type == 'prof', 1, 0)
```

- (b) A linear model with `prestige` as an outcome and `income`, `professional`, and the interaction of the two as predictors was run.

```
1 pres_inc_prof <- lm(prestige ~ income + professional +  
2 income:professional, data = Prestige)
```

The results are in Table 1. The *pvalue* for β_3 is 8.831093×10^{-05} , so we reject the hypothesis that the two variables do not interact, ie we conclude that a change in income does not have the same effect on prestige for professionals and non-professionals (see Figure 1).

¹There were 4 nas in type - athletes, newsboys, babysitters, farmers, these were excluded from the model. See models in Table 2

Table 1: Prestige as a function of professional job and income

| | <i>Dependent variable:</i> |
|-------------------------|-----------------------------|
| | prestige |
| income | 0.003171*** (0.000499) |
| professional | 37.781280*** (4.248274) |
| income:professional | −0.002326*** (0.000567) |
| Constant | 21.142260*** (2.804426) |
| Observations | 98 |
| R ² | 0.787154 |
| Adjusted R ² | 0.780361 |
| Residual Std. Error | 8.011644 (df = 94) |
| F Statistic | 115.877800*** (df = 3; 94) |
| <i>Note:</i> | *p<0.1; **p<0.05; ***p<0.01 |

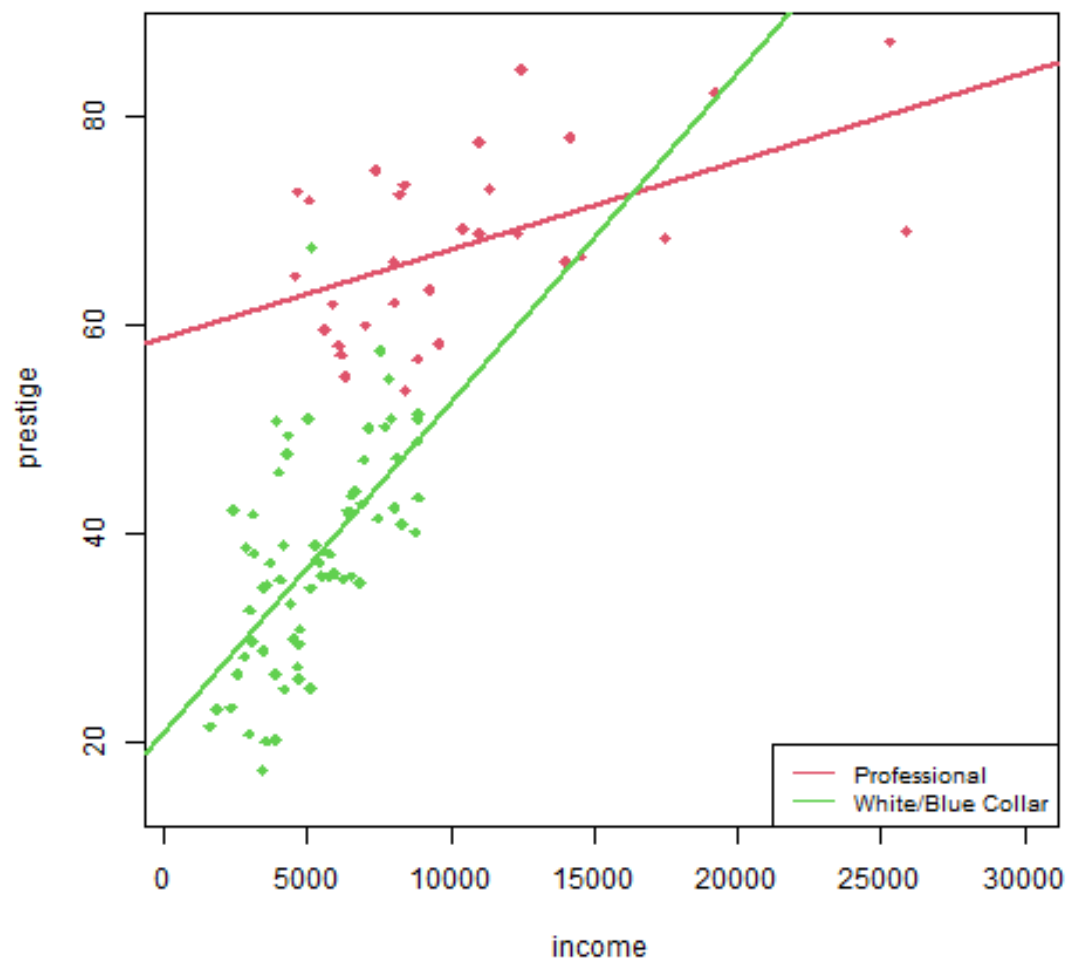


Figure 1: Prestige as a function of income, professional job type and their interaction

(c) **Prediction Equation:**

if $y = \text{prestige}$, $x = \text{income}$, $z = \text{professional}$

$$\hat{y} = \beta_0 + \delta_1 z + \beta_1 x + \delta_2 xz + \epsilon \quad (1)$$

$$\beta_0 = 21.14226, \beta_1 = 0.003171, \delta_1 = 37.78128, \delta_2 = -0.002326$$

ie

$$\text{prestige} = 21.14226 + 37.78128 \times \text{professional} + 0.003171 \times \text{income} + (-0.002326) \times \text{professional} \times \text{income}$$

(d) The coefficient for **income** is $\beta_1 + \delta_2 z = 0.003171 - 0.002326 \times \text{professional}$

Assuming professional status is constant, a \$1 rise in **income** results in a predicted rise in prestige of 0.003171 if in a blue collar or white collar job; prestige rises by 0.000845 per \$ (ie 0.003171-0.002326) if in a professional job.

(e) The coefficient for **professional** is $\delta_1 + \delta_2 x = 37.78128 - 0.002326 \times \text{income}$

If **professional** switches to 1 (ie **type** = **prof**) then **prestige** increases by 37.78128 - 0.002326 × **income**, if **income** is held constant. If **professional** changes to 1 and **income** changes, the change in prestige is 37.78128 - 0.002326 × old **income** + 0.000845 × change in **income**.

(f) If **professional** = 1, and $\Delta \text{income} = 1000$ then the equation becomes

$$\Delta \hat{y} = ((\beta_1 + \delta_2 \times 1) \times 1000$$

The marginal effect on **prestige** of an increase of \$1,000 in income, for professional occupations: = $(0.003171 - 0.002326) \times 1000 = 0.845$

(g) From equation 1, if **income** x is 6000, then when **professional** z_0 was 0, **prestige** was:

$$\hat{y}_0 = (\beta_0) + (\beta_1)x = 21.14226 + 0.003171 \times 6000 = 40.16826$$

when **professional** z_1 is 1, **prestige** is:

$$\hat{y}_1 = (\beta_0 + \delta_1) + (\beta_1 + \delta_2)xz_1$$

$$= (21.14226 + 37.78128) + (0.003171 - 0.002326) * 6000 = 63.99354$$

so,

$$\Delta \hat{y} = \hat{y}_1 - \hat{y}_0 = 63.99354 - 40.16826 = 23.82528$$

(or, $\Delta \hat{y} = \delta_1 + \delta_2 x = 37.78128 - 0.002326 \times 6000$)

The marginal effect of professional jobs when the variable **income** takes the value of \$6,000 is an increase in **prestige** of 23.82528.

Question 2: Political Science

Researchers are interested in learning the effect of all of those yard signs on voting preferences.² Working with a campaign in Fairfax County, Virginia, 131 precincts were randomly divided into a treatment and control group. In 30 precincts, signs were posted around the precinct that read, “For Sale: Terry McAuliffe. Don’t Sellout Virginia on November 5.”

Below is the result of a regression with two variables and a constant. The dependent variable is the proportion of the vote that went to McAuliffe’s opponent Ken Cuccinelli. The first variable indicates whether a precinct was randomly assigned to have the sign against McAuliffe posted. The second variable indicates a precinct that was adjacent to a precinct in the treatment group (since people in those precincts might be exposed to the signs).

Impact of lawn signs on vote share

| | |
|--|------------------|
| Precinct assigned lawn signs (n=30) | 0.042 (0.016) |
| Precinct adjacent to lawn signs (n=76) | 0.042 (0.013) |
| Constant | 0.302 (0.011) |

Notes: $R^2=0.094$, $N=131$

```
1  n <- 131
2  est_coeffs <- 3
3  df <- n - est_coeffs
4  R2 <- 0.094
5
```

- (a) To determine whether having these yard signs in a precinct affects vote share we conduct a hypothesis test with $\alpha = .05$. The null hypothesis is that having yard signs in a precinct have no effect on the voting in that precinct, ie $H_0 : \beta_1 = 0$, $H_{alt} : \beta_1 \neq 0$.

```
1  beta1 <- 0.042
2  n1 <- 30
3  se1 <- 0.016
4
5  t_precinct <- beta1 / se1
6  pval_precinct <- 2*pt(abs(t_precinct), df, lower.tail=FALSE)
7
```

²Donald P. Green, Jonathan S. Krasno, Alexander Coppock, Benjamin D. Farrer, Brandon Lenoir, Joshua N. Zingher. 2016. “The effects of lawn signs on vote outcomes: Results from four randomized field experiments.” Electoral Studies 41: 143-150.

$$tvalue = 0.042/0.016 = 2.625$$

$$pvalue = Pr(tvalue = 2.625, df = 128) = 0.00972002$$

$pvalue < \alpha$, so we reject the null hypothesis and conclude that the presence of signs in yards in a precinct is associated with an average increase in vote for Cuccinelli of 4.2% in that precinct.

- (b) To determine whether being next to precincts with these yard signs affects vote share we conduct a hypothesis test with $\alpha = .05$. The null hypothesis is that having yard signs in a precinct have no effect on the voting in the adjacent precincts, ie $H_0 : \beta_2 = 0, H_{alt} : \beta_2 \neq 0$.

```

1  n2 <- 76
2  beta2 <- 0.042
3  se2 <- 0.013
4
5  t_adj <- beta2 / se2
6  pval_adj <- 2*pt(abs(t_adj), df, lower.tail = FALSE)
7

```

$$tvalue = 0.042/0.013 = 3.230769$$

$$pvalue = 0.00156946$$

$pvalue < \alpha$, so we reject the null hypothesis and conclude that the presence of signs in yards in a precinct is associated with an average increase in vote for Cuccinelli of 4.2% in the adjacent precincts.

- (c) If not in a precinct with yard signs and not in an adjacent precinct, Cuccinelli averages 30.2% of the vote.

At the 95% confidence level, the vote for Cuccinelli in a precinct which is considered to be unaffected by signs is between 28.02% and 32.38%.

```

1  beta0 <- 0.302
2  se0 <- 0.011
3
4  tscore <- qt(0.975, df) # get tscore for df 128,
5
6  CI0.L <- beta0 - tscore*se0
7  CI0.U <- beta0 + tscore*se0
8

```

t-statistic for two-sided test with $\alpha = 0.05$, 128 degrees of freedom = 1.978671

- (d) If the assumptions are valid, the presence of yard signs (even if they don't mention the candidate or their party) can predict an increase in votes of 4.2%³. That could be the difference between winning an losing a close election, so the signs would be worth the

³CI β_1 : 1.03% to 7.37%; CI β_2 : 1.63% to 6.77%, at 95%

cost, particularly as they increased vote share in 106 precincts, even though they were only placed in 30.

The model's value is based on the test being set up with a effectively random/representative assignment of yard signs (it couldn't be actually random, or there would have been precincts with yard signs adjacent to each other). This assumes precincts are homogeneous with regard to party affiliation, voter turnout, through traffic, size (affecting proportion of voters who would have seen the signs), etc.

The model assumes predictors are independent, ie that the vote share in the adjacent precincts is not being affected by some other characteristic based on proximity to the areas with yard signs, and that the change in vote share is not due to a confounding variable which affects the sign and adjacent precincts, but which is not being modelled.

However, the R^2 value is 0.094, ie only 9.4% of the variation in the dependent variable can be explained by the yard sign model. Most of the variation in vote share results from factors which are not modelled. If those missing variables were accounted for they could change both the values and the significance of the coefficients. Alternatively, the relationship between the dependent variable (vote share) and the predictors (yard signs) may not be linear, and an alternative model, with the same predictors, might be appropriate.

1 Appendix

1.0.1 Code

PS04_ImeldaFinn.R

Q1 - model variations depending on treatment of missing values

Went with conservative option of ignoring NAs rather than recoding. The difference to the coefficients and p-values wasn't significant.

```
1 # class all nas as non-professional (based on education)
2 Prestige2 <- Prestige
3 Prestige2['professional'] <- ifelse(is.na(Prestige2$type) ,0,Prestige2$
  professional)
4
5 # class all nas as non-professional, apart from athletes (based on income)
6 Prestige3 <- Prestige2
7 Prestige3["athletes","professional"] <- 1
8
9 pres_nas <- lm(prestige ~ income + professional + income:professional ,
10               data = Prestige2)
11 pres_athletes <- lm(prestige ~ income + professional + income:professional ,
12                    data = Prestige3)
```


Table 2: Effect of na job types on model

| | <i>Dependent variable:</i> | | |
|-------------------------|----------------------------|--------------------------|--------------------------|
| | | prestige | |
| | (1) | (2) | (3) |
| income | 0.0032*** (0.0005) | 0.0033*** (0.0005) | 0.0032*** (0.0005) |
| professional | 37.7813*** (4.2483) | 38.1200*** (4.0798) | 37.1860*** (4.0920) |
| income:professional | -0.0023*** (0.0006) | -0.0024*** (0.0005) | -0.0023*** (0.0005) |
| Constant | 21.1423*** (2.8044) | 20.8035*** (2.5387) | 21.0533*** (2.5748) |
| Observations | 98 | 102 | 102 |
| R ² | 0.7872 | 0.7893 | 0.7863 |
| Adjusted R ² | 0.7804 | 0.7828 | 0.7797 |
| Residual Std. Error | 8.0116 (df = 94) | 8.0181 (df = 98) | 8.0747 (df = 98) |
| F Statistic | 115.8778*** (df = 3; 94) | 122.3380*** (df = 3; 98) | 120.1705*** (df = 3; 98) |

Note:

*p<0.1; **p<0.05; ***p<0.01

Q2

A full F-test was used to evaluate hypothesis that all the coefficients are 0; at 5% reject null hypothesis.

```
1 F.test <-(R2/(k-1))/((1-R2)/(n-k))
2 #F test statistic F = 6.640 with 2 and 129 degrees of freedom
3 df1 <- k - 1
4 df2 <- n-k
5
6 F.pvalue <-df(F.test , df1 , df2)
7 # 0.001634304
8
```