

Linear Probability Models

2018-04-10

Preparation

We will use the data from the *graduation.csv* file in these notes. These data include student-level attributes for $n = 2344$ randomly sampled students who were first-time, full-time freshman from the 2002 cohort at a large, midwestern research university classified by the Carnegie Foundation as having very high research activity. Any students who transferred to another institution were removed from the data. The source of these data is: Jones-White, Radcliffe, Lorenz, & Soria (2014). We will use these data to explore predictors of college graduation.

The attributes in the *graduation.csv* file include:

- degree: Did the student graduate from the institution? (0 = No; 1 = Yes)
- act: Student's ACT score (If the student reported a SAT score, a concordance table was used to transform the score to a comparable ACT score.)
- scholarship: Amount of scholarship offered to student (in thousands of dollars)
- ap: Number of Advanced Placement credits at time of enrollment
- firstgen: Is the student a first generation college student? (0 = No; 1 = Yes)
- nontrad: Is the student a non-traditional student (older than 19 years old at the time of freshman enrollment)? (0 = No; 1 = Yes)

```
# Load libraries
library(broom)
library(corr)
library(dplyr)
library(ggplot2)
library(readr)
library(sm)

# Read in data
grad = read_csv(file = "~/Dropbox/epsy-8252/data/graduation.csv")
head(grad)
```

degree	act	scholarship	ap	firstgen	nontrad
1	21	0.0	0	0	0
1	19	0.0	0	0	0
1	27	0.0	0	1	0
1	25	0.5	0	1	0
0	28	0.0	17	1	0
1	21	0.0	0	0	1

Data Exploration

To begin the analysis, we will explore the outcome variable `degree`. Since this is a dichotomous variable, we can look at counts/proportions. The analysis suggests that most freshmen who enroll at this institution tend to graduate (73%).

```
grad %>% group_by(degree) %>% summarize(Count = n(), Prop = n() / nrow(grad))
```

degree	Count	Prop
0	627	0.267
1	1717	0.733

We will also explore the `act` variable, which we will use as a predictor in the analysis.

```
sm.density(grad$act, xlab = "ACT score")
```

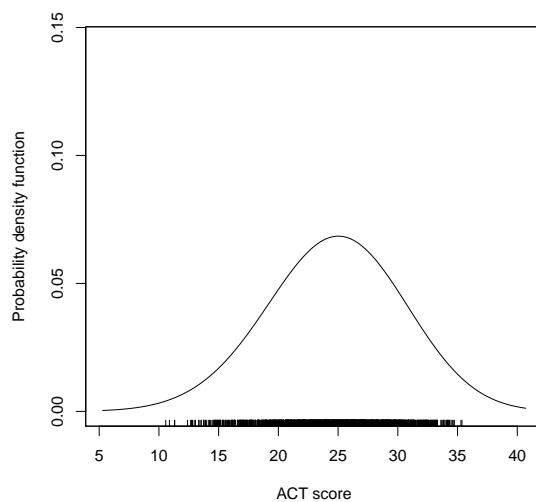


Figure 1: Density plot of the ACT scores.

```
grad %>% summarize( M = mean(act), SD = sd(act) )
```

M	SD
24.8	4.15

The plot is approximately normally distributed and indicates that freshmen at this institution have a mean ACT score near 25. While there is a great deal of variation in ACT scores (scores range from 10 to 36), most students have a score between 21 and 29.

Relationship between ACT Score and Graduation

When the outcome variable was continuous, we examined relationships graphically via a scatterplot. If we try that when the outcome is dichotomous, we run into problems.

```
ggplot(data = grad, aes(x = act, y = degree)) +  
  geom_point() +  
  theme_bw() +  
  xlab("ACT score") +  
  ylab("Graduated")
```

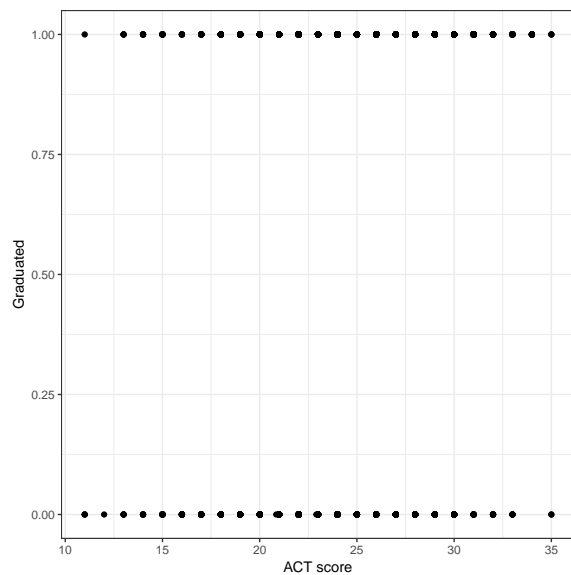


Figure 2: Scatterplot of whether a student graduated versus ACT score.

Since there are only two values for the outcome, many of the observations are over-plotted, and it is impossible to tell anything about the relationship between the variables. One solution to this problem is to add (or subtract) a tiny amount to each student's degree value. This is called “jittering” the observations“. To do this, use the `jitter()` function. (Note that the amount of jittering can be adjusted by including an additional argument to the `jitter()` function.)

```
ggplot(data = grad, aes(x = act, y = jitter(degree))) +
  geom_point() +
  geom_smooth(method = "lm", se = FALSE) +
  theme_bw() +
  xlab("ACT score") +
  ylab("Graduated")
```

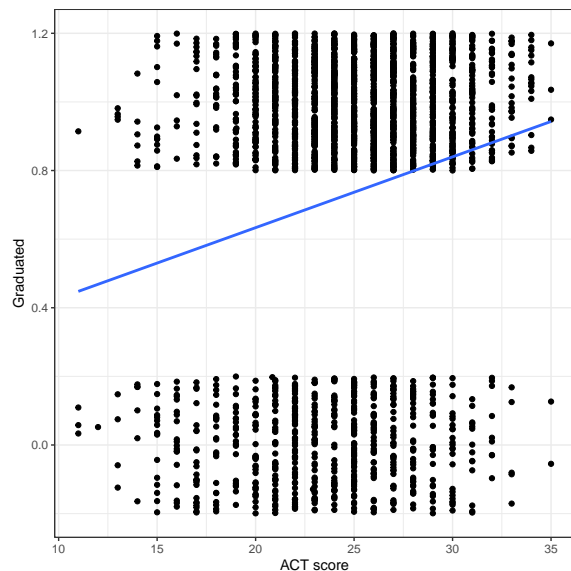


Figure 3: Scatterplot of whether a student graduated versus ACT score. The points have been jittered to alleviate over-plotting, and the regression smoother has also been added to the plot.

This spreads out the observations vertically, so we no longer have the problem of overplotting. But, the relationship is still difficult to discern. Adding the regression smoother helps us see that there is a positive relationship between ACT scores and graduation.

What does this mean? To understand this, we have to think about what the linear regression is modeling. Remember that the regression model is predicting the average Y at each X . Since our Y is dichotomous, the average represents the proportion of students with a 1. In other words, the regression predicts the proportion of students who graduate for a particular ACT score. The positive relationship between ACT score and graduation suggests that higher ACT scores are associated with higher proportions of students who graduate. We can also see this relationship by examining the correlation matrix between the two variables.

```
grad %>%
  select(degree, act) %>%
  correlate()
```

rowname	degree	act
degree	NA	0.195
act	0.195	NA

Fitting the Linear Probability (Proportion) Model

We can fit the linear model shown in the plot using the `lm()` function as we always have.

```
lm.1 = lm(degree ~ 1 + act, data = grad)
summary(lm.1)

##
## Call:
## lm(formula = degree ~ 1 + act, data = grad)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.946 -0.550  0.221  0.283  0.554
##
## Coefficients:
##              Estimate Std. Error t value      Pr(>|t|)
## (Intercept)  0.21665    0.05431   3.99      0.000068 ***
## act          0.02084    0.00216   9.63 < 0.0000000000000002 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.434 on 2342 degrees of freedom
## Multiple R-squared:  0.0381, Adjusted R-squared:  0.0377
## F-statistic: 92.8 on 1 and 2342 DF, p-value: <0.0000000000000002
```

Differences in ACT score account for 3.8% of the variation in graduation. This is statistically significant, $F(1, 2342) = 92.8, p < .001$. The fitted model is

$$\hat{\pi}_i = 0.22 + 0.02(\text{ACT Score}_i)$$

where $\hat{\pi}_i$ is the predicted proportion of students who graduate. Interpreting the coefficients,

- On average, 0.22 of students having an ACT score of 0 are predicted to graduate. (Extrapolation)
- Each one-point difference in ACT score is associated with an additional 0.02 predicted improvement in the proportion of students graduating, on average.

Let's examine the model assumptions.

```
out = augment(lm.1)
#head(out)

# Examine normality assumption
sm.density(out$.std.resid, xlab = "Standardized residuals")
```

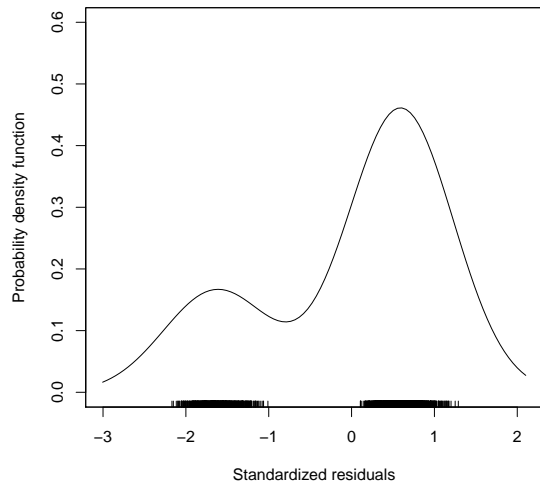


Figure 4: Density plot of the standardized residuals from the linear probability model.

```
# Examine linearity and homoskedasticity
ggplot(data = out, aes(x = .fitted, y = .std.resid)) +
  geom_point() +
  geom_hline(yintercept = 0) +
  theme_bw() +
  xlab("Fitted values") +
  ylab("Standardized residuals")
```

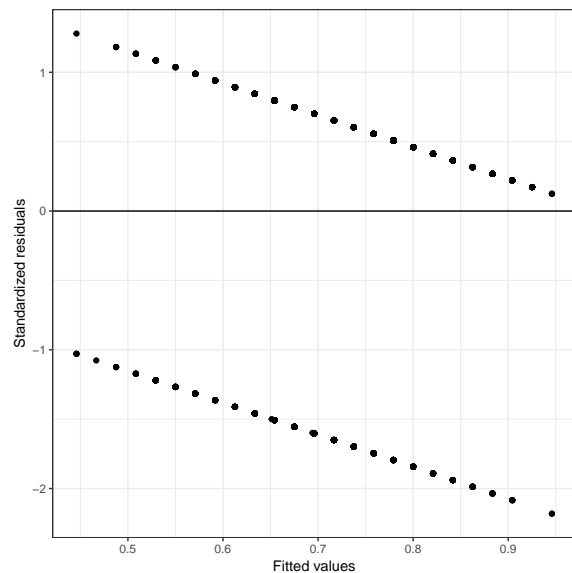


Figure 5: Scatterplot of the standardized residuals versus the fitted values from the linear probability model.

It is clear that the assumptions associated with linear regression are violated. First off, the residuals are not normally distributed. They are in fact, bimodal. The scatterplot of the residuals versus the fitted values also indicates violation of the linearity assumption.

Understanding the Residuals from the Linear Probability Model

Look closely at the scatterplot of the residuals vs. the fitted values. At each fitted value, there are only two residual values. Why is this? Recall that residuals are computed as $\epsilon_i = Y_i - \hat{Y}_i$. Now, remember that Y_i can only be one of two values, 0 or 1. Also remember that in the linear probability model $\hat{Y}_i = \hat{\pi}_i$. Thus, for $Y = 0$,

$$\begin{aligned}\epsilon_i &= 0 - \hat{Y}_i \\ &= -\hat{\pi}_i\end{aligned}$$

And, if $Y = 1$,

$$\begin{aligned}\epsilon_i &= 1 - \hat{Y}_i \\ &= 1 - \hat{\pi}_i\end{aligned}$$

This means that the residual computed using a particular fitted value (or from a particular ACT value given that the fitted values is a function of X) can only take on one of two values: $-\hat{\pi}_i$ or $1 - \hat{\pi}_i$. (Plotting the residuals for all fitted values, the marginal distribution of the residuals, results in a bimodal distribution for this same reason.)

References

Jones-White, D. R., Radcliffe, P. M., Lorenz, L. M., & Soria, K. M. (2014). Priced out?: The influence of financial aid on the educational trajectories of first-year students starting college at a large research university. *Research in Higher Education*, 55(4), 329–350.