

Nonlinearity: Polynomial Effects

2017-04-09

Preparation

In this set of notes, you will learn one method of dealing with nonlinearity. Specifically, we will look at the including polynomial effects into a model. The data we will use in this set of notes, *mnSchools.csv*, contains institutional data for several Minnesota colleges and universities collected in 2011. The variables are:

- **id**: Institution ID number
- **name**: Institution name
- **gradRate**: Six-year graduation rate. This measure represents the proportion of first-time, full-time, bachelor's or equivalent degree-seeking students who started in Fall 2005 and graduated within 6 years.
- **public**: Dummy variable indicating educational sector (0 = private institution; 1 = public institution)
- **sat**: Estimated median SAT score for incoming freshmen at the institution
- **tuition**: Cost of attendance for full-time, first-time degree/certificate-seeking in-state undergraduate students living on campus for academic year 2013-14.

These source of these data is: <http://www.collegeresults.org>. We will examine use these data to examine if (and how) academic “quality” of the student-body (measured by SAT score) is related to institutional graduation rate.

```
mn = read.csv(file = "~/Google Drive/Documents/epsy-8251/data/mnSchools.csv")
head(mn)
```

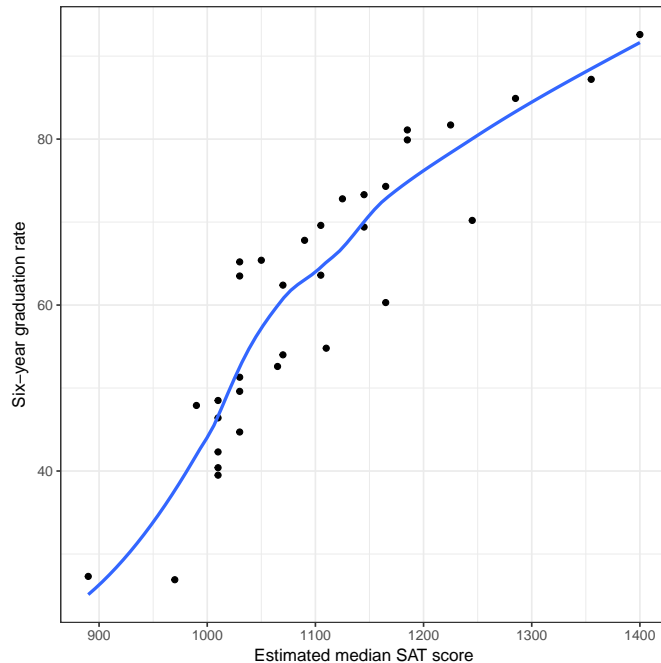
	id	name	gradRate	public	sat	tuition
1	1	Augsburg College	65.2	0	1030	39294
2	3	Bethany Lutheran College	52.6	0	1065	30480
3	4	Bethel University, Saint Paul, MN	73.3	0	1145	39400
4	5	Carleton College	92.6	0	1400	54265
5	6	College of Saint Benedict	81.1	0	1185	43198
6	7	Concordia College at Moorhead	69.4	0	1145	36590

```
# Load libraries
library(ggplot2)
library(sm)
```

Examine Relationship between Graduation Rate and SAT Scores

As always, we begin the analysis by graphing the data.

```
ggplot(data = mn, aes(x = sat, y = gradRate)) +
  geom_point() +
  geom_smooth(se = FALSE) +
  theme_bw() +
  xlab("Estimated median SAT score") +
  ylab("Six-year graduation rate")
```



The loess smoother suggests that the relationship between SAT scores and graduation rate is non-linear. Nonlinearity implies that the effect of SAT on graduation rates is not constant across the range of X ; for colleges with lower values of SAT (say $\text{SAT} < 1100$) the effect of SAT has a rather high, positive effect (steep slope), while for colleges with higher values of SAT (≥ 1100) the effect of SAT is positive and moderate (the slope is less steep). Another way of saying this is that for schools with lower SAT scores, a one-unit difference in SAT is associated with a larger change in graduation rates than the same one-unit change for schools with higher SAT values.

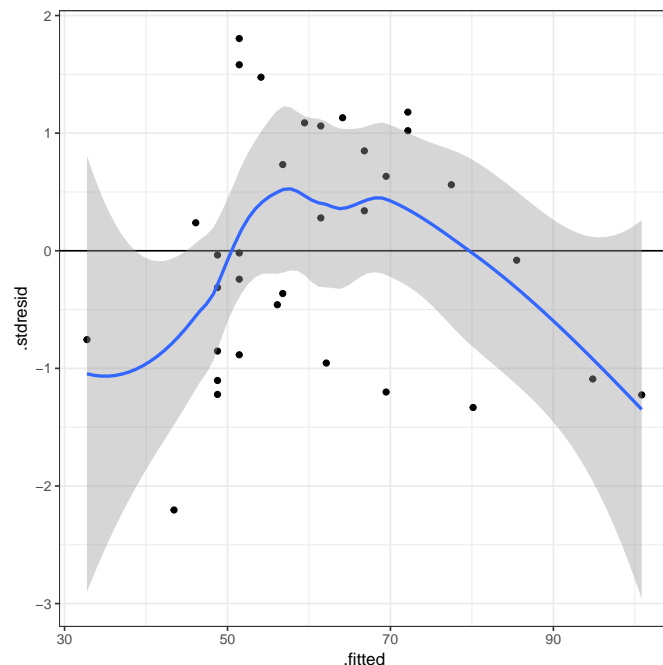
Residual Plot: Another Way to Spot Nonlinearity

Sometimes, the nonlinear relationship is difficult to detect from the scatterplot of Y versus X . Often it helps to fit the linear model and then examine the assumption of linearity in the residuals. It is sometimes easier to detect nonlinearity in the plot of the residuals versus the fitted values.

```
# Fit linear model
lm.1 = lm(gradRate ~ 1 + sat, data = mn)

# Obtain residuals
out = fortify(lm.1)

# Examine residuals for linearity
ggplot(data = out, aes(x = .fitted, y = .stdresid)) +
  geom_point() +
  geom_hline(yintercept = 0) +
  geom_smooth() +
  theme_bw()
```



This plot suggests that the assumption of linearity may be violated. There is systematic over-estimation for low fitted values, systematic under-estimation for moderate fitted values, and systematic over-estimation for high fitted values.

Polynomial Models

One way of modeling non-linearity is by including polynomial effects. In regression, a polynomial effects are predictors that have a power greater than one. For example, x^2 (quadratic term), or x^3 (cubic term). Note that

$$x^2 = x \times x.$$

So the quadratic term, x^2 is a product of x times itself. Recall that products are how we express interactions. Thus the quadratic term of x^2 is really the interaction of x with itself. To model this, we simply (1) create the product term, and (2) include the product term and all constituent main-effects in the regression model.

```
mn$sat_quadratic = mn$sat * mn$sat
head(mn)
```

	id	name	gradRate	public	sat	tuition
1	1	Augsburg College	65.2	0	1030	39294
2	3	Bethany Lutheran College	52.6	0	1065	30480
3	4	Bethel University, Saint Paul, MN	73.3	0	1145	39400
4	5	Carleton College	92.6	0	1400	54265
5	6	College of Saint Benedict	81.1	0	1185	43198
6	7	Concordia College at Moorhead	69.4	0	1145	36590

	sat_quadratic
1	1060900
2	1134225
3	1311025
4	1960000

```
5      1404225
6      1311025
```

```
# Fit model
lm.2 = lm(gradRate ~ 1 + sat + sat_quadratic, data = mn)
summary(lm.2)
```

Call:

```
lm(formula = gradRate ~ 1 + sat + sat_quadratic, data = mn)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-12.7867	-5.0969	0.3968	5.0011	13.6869

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-366.34287223	98.62031987	-3.715	0.000831 ***
sat	0.62716548	0.17270382	3.631	0.001040 **
sat_quadratic	-0.00021503	0.00007507	-2.864	0.007559 **

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 7.019 on 30 degrees of freedom

Multiple R-squared: 0.8351, Adjusted R-squared: 0.8241

F-statistic: 75.97 on 2 and 30 DF, p-value: 1.81e-12

Since this is an interaction model, we start by examining the interaction term; the quadratic coefficient. This term is statistically reliable ($p = .008$), suggesting that the quadratic term explains variation above and beyond the linear term. This suggests that we should keep the quadratic term in the model.

Interpretation of a Significant Polynomial Term

How do we interpret the quadratic term? First, we will write out the fitted model.

$$\widehat{\text{Graduation Rate}} = -366.3 + 0.63(\text{SAT}) - 0.0002(\text{SAT}^2)$$

From algebra, you may remember that the coefficient in front of the quadratic term (-0.0002) informs us of whether the quadratic is an upward-facing U-shape, or a downward-facing U-shape. Since our term is negative, the U-shape is downward-facing. It also indicates whether the U-shape is skinny or wide. The intercept and linear terms help us locate the U-shape in the coordinate plane (moving it right, left, up, or down from the origin). You could work these out algebraically, but typically, we will just plot the predicted values and interpret from the plot.

Refit the model using the I() function

Before we create the plot, we use a different method of fitting polynomial terms in a regression. Rather than create a new variable in the data set, we insert the polynomial directly into the model using the I() function.

```
lm.2 = lm(gradRate ~ 1 + sat + I(sat ^ 2), data = mn)
summary(lm.2)
```

Call:

```
lm(formula = gradRate ~ 1 + sat + I(sat^2), data = mn)
```

Residuals:

Min	1Q	Median	3Q	Max
-12.7867	-5.0969	0.3968	5.0011	13.6869

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	-366.34287223	98.62031987	-3.715	0.000831	***
sat	0.62716548	0.17270382	3.631	0.001040	**
I(sat^2)	-0.00021503	0.00007507	-2.864	0.007559	**

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 7.019 on 30 degrees of freedom

Multiple R-squared: 0.8351, Adjusted R-squared: 0.8241

F-statistic: 75.97 on 2 and 30 DF, p-value: 1.81e-12

The I() function forces R to create a separate term for the quadratic (which is essentially what we do by adding the product term to the model). Without it, R will do the computation `sat + sat^2` and use those values as a single variable...not what we want. The other advantage for plotting is that we have only used a single predictor, `sat` in specifying the model. Thus we only need to include `sat` in our plot data rather than both `sat` and `sat_quadratic`. (Note: You cannot use the colon notation (`:`) to fit a polynomial term in the model.)

```
# Set up data
plotData = expand.grid(
  sat = seq(from = 890, to = 1400, by = 10)
)

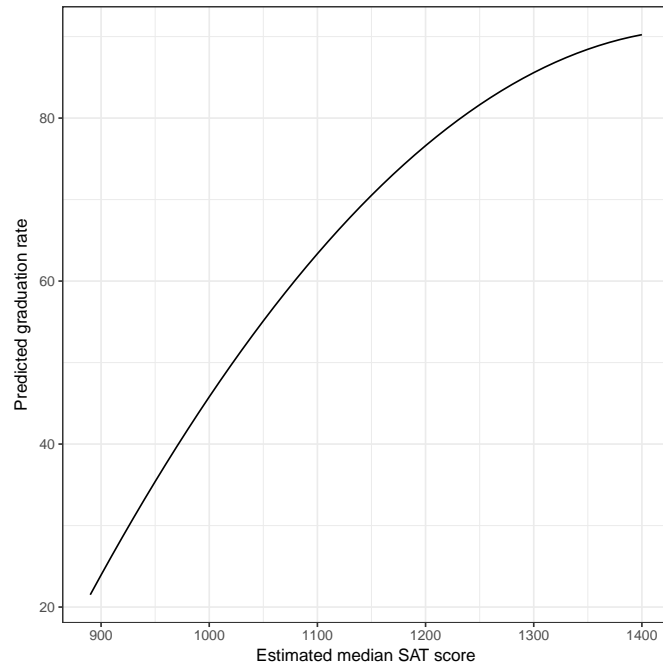
# Predict
plotData$yhat = predict(lm.2, newdata = plotData)

# Examine data
head(plotData)
```

	sat	yhat
1	890	21.50981
2	900	23.93244
3	910	26.31207
4	920	28.64869
5	930	30.94230
6	940	33.19291

```
# Plot
ggplot(data = plotData, aes(x = sat, y = yhat)) +
```

```
geom_line() +
theme_bw() +
xlab("Estimated median SAT score") +
ylab("Predicted graduation rate")
```



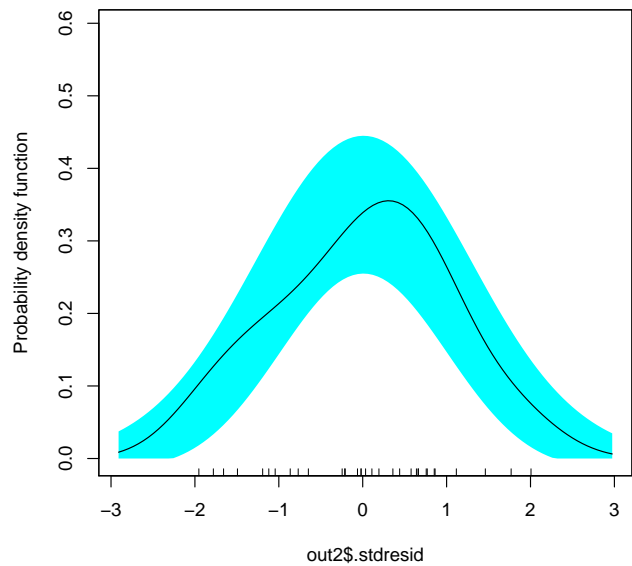
The quadratic relationship is expressed in the predicted values. Aside from plotting them versus SAT scores, there is nothing further we need to do to get the quadratic relationship to appear. The plot, more importantly, helps us interpret the relationship between SAT scores and graduation rates. The effect of SAT on graduation rate depends on SAT score (definition of an interaction). For schools with low SAT scores, the effect of SAT score on graduation rate is positive and fairly high. For schools with high SAT scores, the effect of SAT score on graduation rate remains positive, but it has a smaller effect on graduation rates.

Re-Examining the Residuals for the Quadratic Model

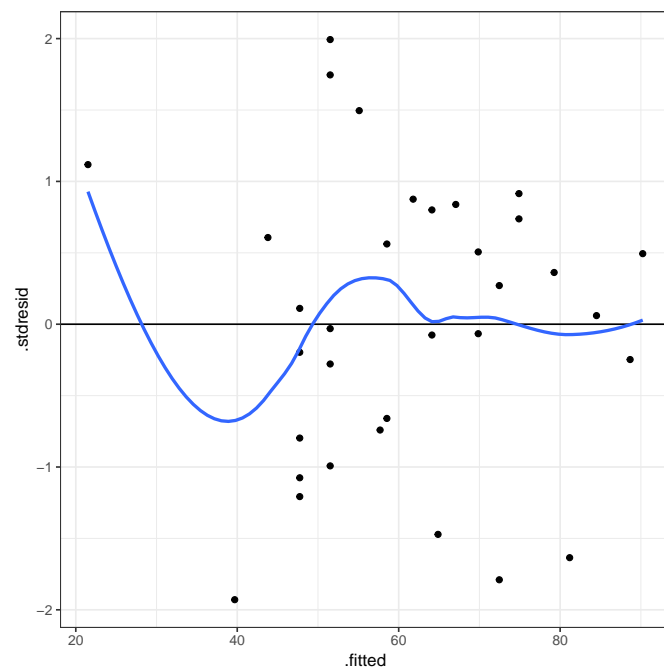
Since we fitted a different model, we should examine the residuals to see whether the assumptions for the model seem satisfied.

```
out2 = fortify(lm.2)

# Check normality
sm.density(out2$.stdresid, model = "normal")
```



```
# Check other assumptions
ggplot(data = out2, aes(x = .fitted, y = .stdresid)) +
  geom_point() +
  geom_hline(yintercept = 0) +
  geom_smooth(se = FALSE) +
  theme_bw()
```



Based on the plots, the quadratic model seems to meet the assumptions for regression. At the very least, it clearly does so better than the linear model.

Adding Covariates

We can also include covariates in a polynomial model (to control for other predictors), the same way we do in a linear model, by including them as additive terms in the `lm()` model. Below we include the `public` dummy-coded predictor to control for the effects of sector.

```
lm.3 = lm(gradRate ~ 1 + sat + I(sat ^ 2) + public, data = mn)
summary(lm.3)
```

Call:

```
lm(formula = gradRate ~ 1 + sat + I(sat^2) + public, data = mn)
```

Residuals:

Min	1Q	Median	3Q	Max
-16.139	-3.207	0.687	2.902	10.388

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-384.16062215	79.41299015	-4.838	0.0000398 ***
sat	0.67037086	0.13925237	4.814	0.0000425 ***
I(sat^2)	-0.00023707	0.00006059	-3.912	0.000507 ***
public	-9.12481278	2.18743115	-4.171	0.000251 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 5.644 on 29 degrees of freedom

Multiple R-squared: 0.897, Adjusted R-squared: 0.8863

F-statistic: 84.14 on 3 and 29 DF, p-value: 2.048e-14

Here there is still a quadratic effect of SAT on graduation rates, even after controlling for differences in sector. Similarly, there are differences in sector (public schools have a graduation rate 9.1% lower than private schools, on average), even after controlling for the linear and quadratic effects of SAT. Plot the fitted model to aid interpretation.

```
# Set up data
plotData = expand.grid(
  sat = seq(from = 890, to = 1400, by = 10),
  public = c(0, 1)
)

# Predict
plotData$yhat = predict(lm.3, newdata = plotData)

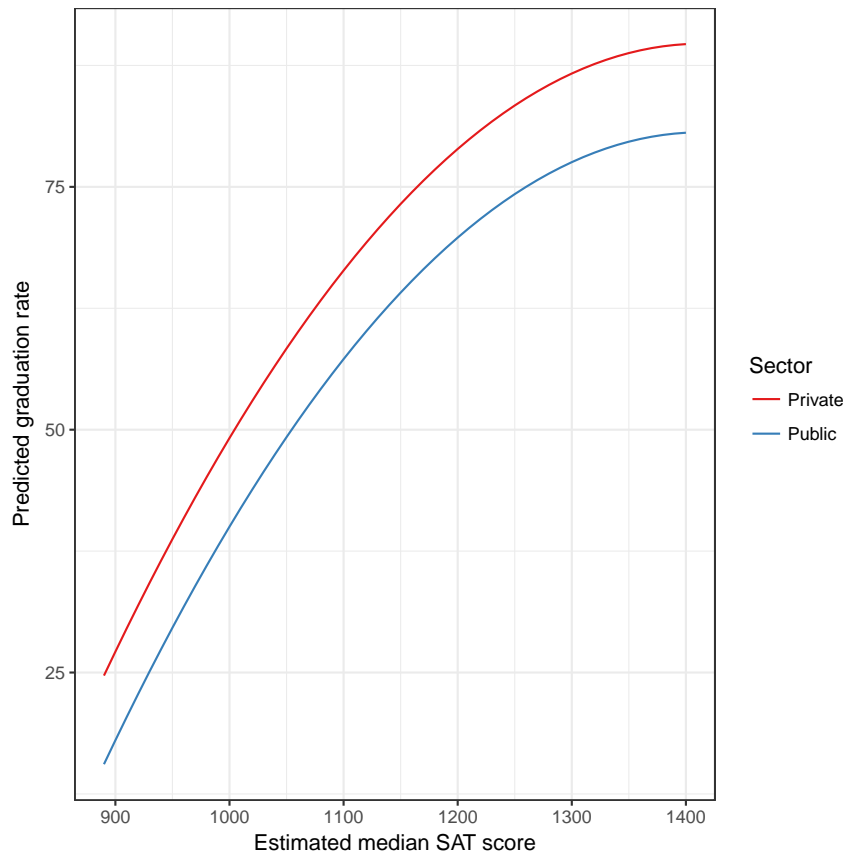
# Coerce public into a factor for better plotting
plotData$public = factor(plotData$public, levels = c(0, 1), labels = c("Private", "Public"))

# Examine data
head(plotData)
```

	sat	public	yhat
1	890	Private	24.68509
2	900	Private	27.14522
3	910	Private	29.55794
4	920	Private	31.92324
5	930	Private	34.24112

6 940 Private 36.51159

```
# Plot
ggplot(data = plotData, aes(x = sat, y = yhat, group = public, color = public)) +
  geom_line() +
  theme_bw() +
  xlab("Estimated median SAT score") +
  ylab("Predicted graduation rate") +
  scale_color_brewer(name = "Sector", palette = "Set1")
```



The plot shows the quadratic effect of SAT scores on graduation rate; the effect of SAT on graduation rates is positive, but this effect declines for increasingly higher SAT scores, after controlling for sector differences. Private schools have higher graduation rates, on average, than public schools for all levels of SAT score.