

Notebook 4 - Machine Learning

May 24, 2024

Reference Guide for R (student resource) - Check out our reference guide for a full listing of useful R commands for this project.

0.1 Data Science Project: Use data to determine the best and worst colleges for conquering student debt.

0.1.1 Notebook 4: Machine Learning

Does college pay off? We'll use some of the latest data from the US Department of Education's College Scorecard Database to answer that question.

In this notebook (the 4th of 4 total notebooks), you'll use R to add polynomial terms to your multiple regression models (i.e. polynomial regression). Then, you'll use the principles of machine learning to tune models for a prediction task on *unseen* data.

```
[1]: ## Run this code but do not edit it. Hit Ctrl+Enter to run the code.  
# This command downloads a useful package of R commands  
library(coursekata)
```

```
CourseKata packages  
coursekata 0.15.0  
  dslabs           0.8.0           Metrics  
  0.1.4  
  Lock5withR       1.2.2           lsr  
  0.5.2  
  fivethirtyeightdata 0.1.0       mosaic  
  1.9.1  
  fivethirtyeight   0.6.2           supernova  
  3.0.0
```

0.1.2 The Dataset (four_year_colleges.csv)

General description - In this notebook, we'll be using the `four_year_colleges.csv` file, which only includes schools that offer four-year bachelors degrees and/or higher graduate degrees. Community colleges and trade schools often have different goals (e.g. facilitating transfers, direct career education) than institutions that offer four-year bachelors degrees. By comparing four-year colleges only to other four-year colleges, we'll have clearer analyses and conclusions.

This data is a subset of the US Department of Education's College Scorecard Database. The data is current as of the 2020-2021 school year.

Description of all variables: See [here](#)

Detailed data file description: See [here](#)

```
[2]: ## Run this code but do not edit it. Hit Ctrl+Enter to run the code.
# This command downloads data from the file 'colleges.csv' and stores it in an
  ↪ object called `dat`
dat <- read.csv('https://skewthescript.org/s/four_year_colleges.csv')
```

```
[3]: str(dat)
```

```
'data.frame':  1053 obs. of  26 variables:
 $ OPEID          : int  100200 105200 105500 100500 105100 831000
100900 101200 100300 101900 ...
 $ name           : chr   "Alabama A & M University" "University of
Alabama at Birmingham" "University of Alabama in Huntsville" "Alabama State
University" ...
 $ city           : chr   "Normal" "Birmingham" "Huntsville" "Montgomery"
...
 $ state          : chr   "AL" "AL" "AL" "AL" ...
 $ region         : chr   "South" "South" "South" "South" ...
 $ median_debt    : num   15.2 15.1 14 17.5 17.7 ...
 $ default_rate   : num   12.1 4.8 4.7 12.8 4 8.2 2.6 4.4 9.9 10 ...
 $ highest_degree : chr   "Graduate" "Graduate" "Graduate" "Graduate" ...
 $ ownership      : chr   "Public" "Public" "Public" "Public" ...
 $ locale         : chr   "Small City" "Small City" "Small City" "Small
City" ...
 $ hbcu           : chr   "Yes" "No" "No" "Yes" ...
 $ admit_rate     : num   89.7 80.6 77.1 98.9 80.4 ...
 $ SAT_avg        : int   959 1245 1300 938 1262 1061 1302 1202 1068 1101
...
 $ online_only    : chr   "No" "No" "No" "No" ...
 $ enrollment     : int   5090 13549 7825 3603 30610 4301 24368 1129 1834
917 ...
 $ net_price      : num   15.5 16.5 17.2 19.5 20.9 ...
 $ avg_cost       : num   23.4 25.5 24.9 21.9 30 ...
 $ net_tuition    : num   8.1 11.99 8.28 9.3 14.71 ...
 $ ed_spending_per_student: num   4.84 14.69 8.32 9.58 9.65 ...
 $ avg_faculty_salary : num   7.6 11.38 9.7 7.19 10.35 ...
 $ pctPELL        : num   71 34 24 73.7 17.2 ...
 $ pct_fed_loan   : num   75 46.9 38.5 78 36.4 ...
 $ grad_rate      : num   28.7 61.2 57.1 31.8 72.1 ...
 $ pct_firstgen   : num   36.6 34.1 31 34.3 22.6 ...
 $ med_fam_income : num   23.6 34.5 44.8 22.1 66.7 ...
 $ med_alum_earnings : num   36.3 47 54.4 32.1 52.8 ...
```

0.1.3 1.0 - Motivating non-linear regression

So far, we've focused entirely on **linear regression** and **multiple linear regression** models, which use linear functions to relate predictors (e.g. `net_tuition`, `grad_rate`, `pctPELL`) to the outcome (`default_rate`).

In this notebook, we're going to investigate ways to model **non-linear** relationships. To make this task a bit more manageable at the start, let's reduce the size of our dataset by taking a random sample of 20 colleges from the `dat` dataframe. We will store our sample in a new R dataframe called `sample_dat`.

```
[ ]: ## Run this code but do not edit it
# create a dataset to train the model with 20 randomly selected observations
set.seed(2)
sample_dat <- sample(dat, size = 20)
```

```
[ ]: #library("dplyr")
#sample_dat = select_if(sample_dat, is.numeric)
#for (col in colnames(sample_dat)){
#  val = sapply(sample_dat[col], typeof)
#}
#sapply(sample_dat, class)
nums <- unlist(lapply(sample_dat, is.numeric), use.names = FALSE)
sample_dat = sample_dat[, nums]
print(nums)
#str(sample_dat)
```

```
[4]: set.seed(2)
sample_dat <- sample(dat, size = 300)
```

Note: When getting a random sample, we'll get different results each time we run our code because it's ... well ... random. This can be quite annoying. So, in the code above, we used the command `set.seed(2)`. This ensures that each time the code is executed, we get the same results for our random sample - the results stored in seed 2. We could have also set the seed to 1 or 3 or 845 or 12345. The seed numbers serve merely as a unique ID that corresponds to a certain result from a random draw. By setting a certain seed, we'll always get a certain random draw.

1.1 Let's take a look at our sample data set. Print out the `head` and `dim` of `sample_dat`.

```
[100]: # Your code goes here
head(sample_dat)
```

		OPEID	name	city	state	region
		<int>	<chr>	<chr>	<chr>	<chr>
A data.frame: 6 × 27	975	379400	Saint Martin's University	Lacey	WA	Far West
	710	316100	Northeastern State University	Tahlequah	OK	Rockies & Southw
	774	330400	Muhlenberg College	Allentown	PA	Northeast
	416	231800	Spring Arbor University	Spring Arbor	MI	Midwest
	392	223400	Adrian College	Adrian	MI	Midwest
	273	189200	University of Iowa	Iowa City	IA	Midwest

```
[8]: # Your code goes here
dim(sample_dat)
```

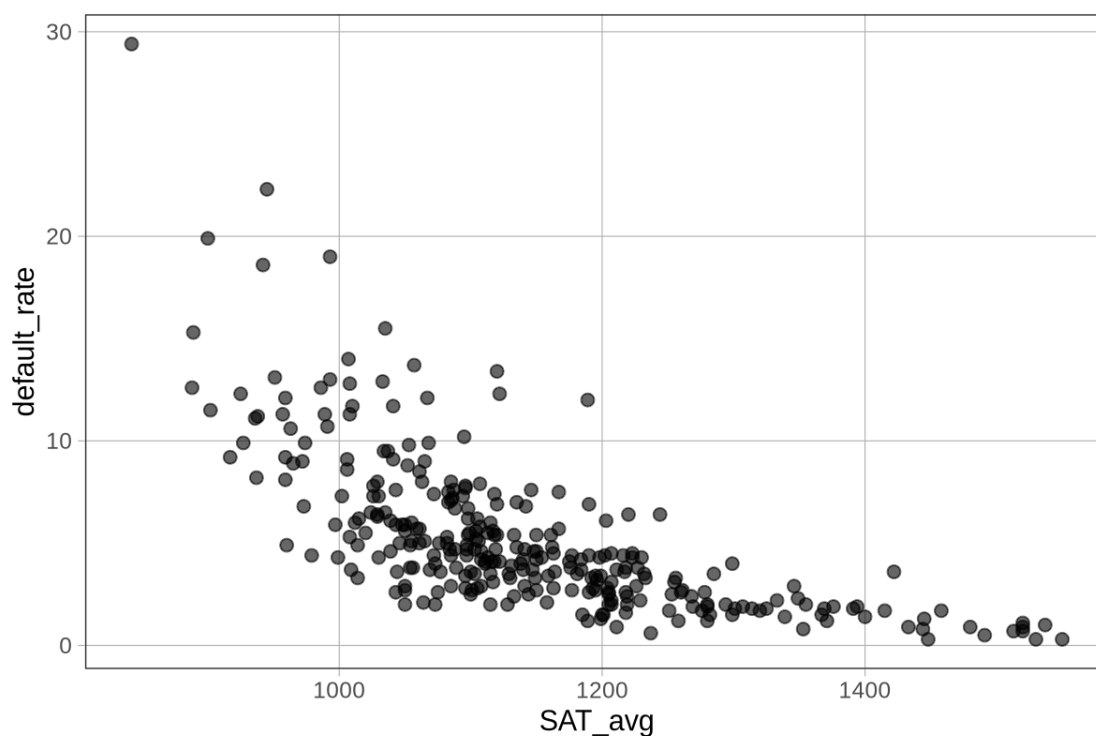
1. 300 2. 27

Check yourself: The dimensions of `sample_dat` should be 20 rows and 27 columns.

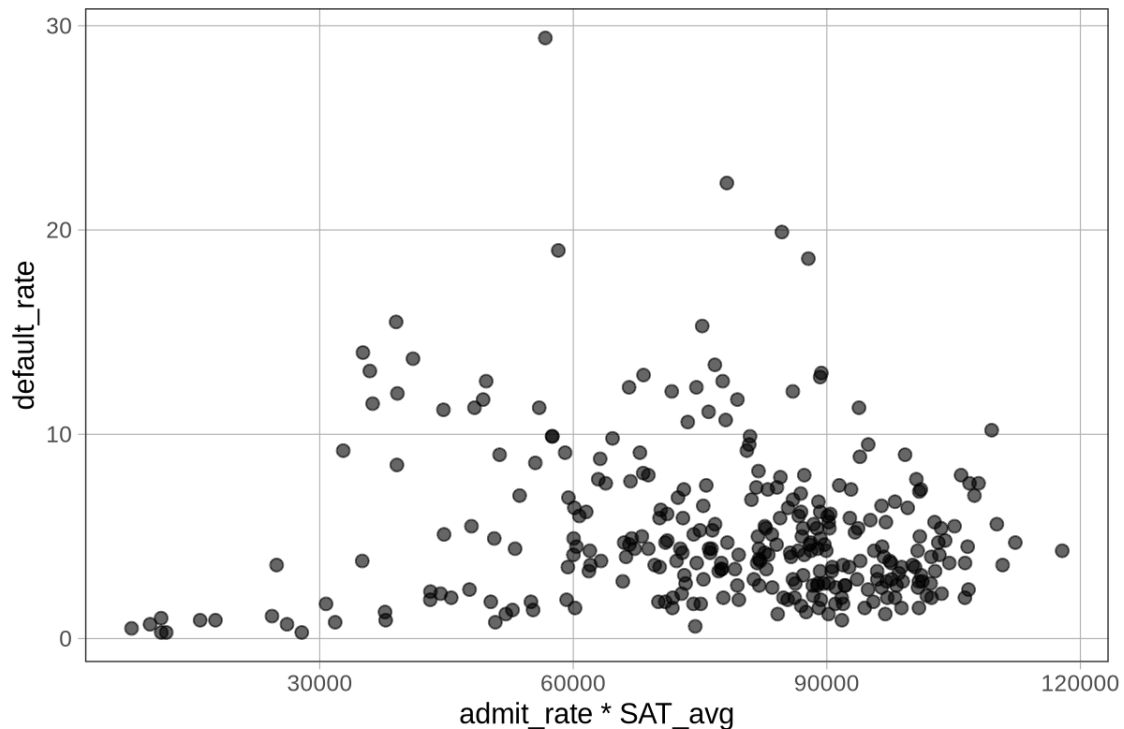
In prior notebooks, we focused on institutional and economic predictors of student loan default rates. In this notebook, we'll begin by analyzing an *academic* variable: `SAT_avg`. This variable shows the average SAT score of students who matriculate to a college.

The following code creates a scatterplot of the relationship between `SAT_avg` (predictor) and `default_rate` (outcome) from the dataset `sample_dat`:

```
[5]: ## Run this code but do not edit it
# create scatterplot: default_rate ~ SAT_avg
gf_point(default_rate ~ SAT_avg, data = sample_dat)
```



```
[13]: gf_point(default_rate ~ admit_rate, data = sample_dat)
```



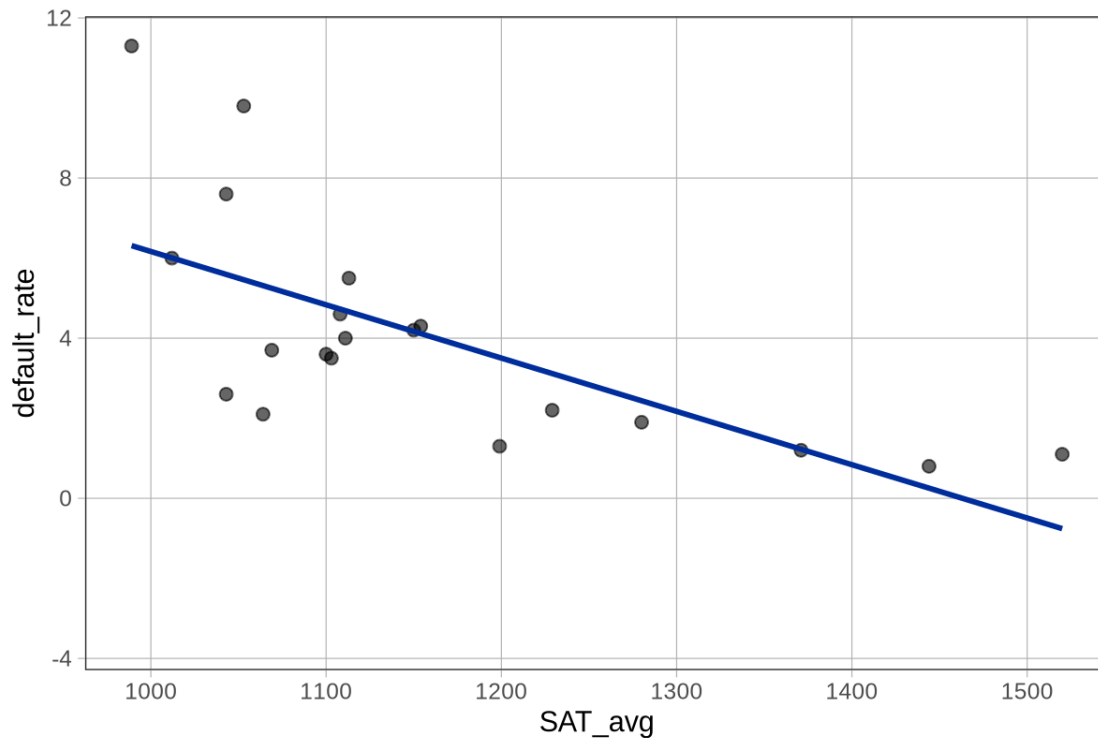
1.2 - Describe the direction of the relationship between `SAT_avg` and `default_rate`. Is it positive or negative? Why do you think this is?

Double-click this cell to type your answer here: There is a weak, negative, nonlinear relationship between `SAT_avg` and `default_rate`. This is evident from the fact that as `SAT_avg` increases, `default_rate` decreases.

1.3 - Create the same scatterplot as above, but with the simple linear model between `default_rate` (outcome) and `SAT_avg` (predictor) overlayed on top.

Hint: Recall the `gf_lm` command from notebook 2.

```
[26]: # Your code goes here
gf_point(default_rate ~ SAT_avg, data = sample_dat) %>% gf_lm(default_rate ~ SAT_avg, data = sample_dat)
```



```
[262]: for (col in
  list("admit_rate",
    "SAT_avg",
    "enrollment",
    "net_price",
    "avg_cost",
    "net_tuition",
    "ed_spending_per_studen",
    "avg_faculty_salary",
    "pctPELL",
    "pct_fed_loan",
    "grad_rate",
    "pct_firstgen",
    "med_fam_income",
    "med_alum_earnings")) {
  print(col)
  #gf_point(default_rate ~ sample_data$col, data = sample_dat) %>%
  ↪gf_lm(default_rate ~ sample_data$col, data = sample_dat)
}
gf_point(default_rate ~ admit_rate, data = sample_dat) %>% gf_lm(default_rate ~
  ↪admit_rate, data = sample_dat)
```

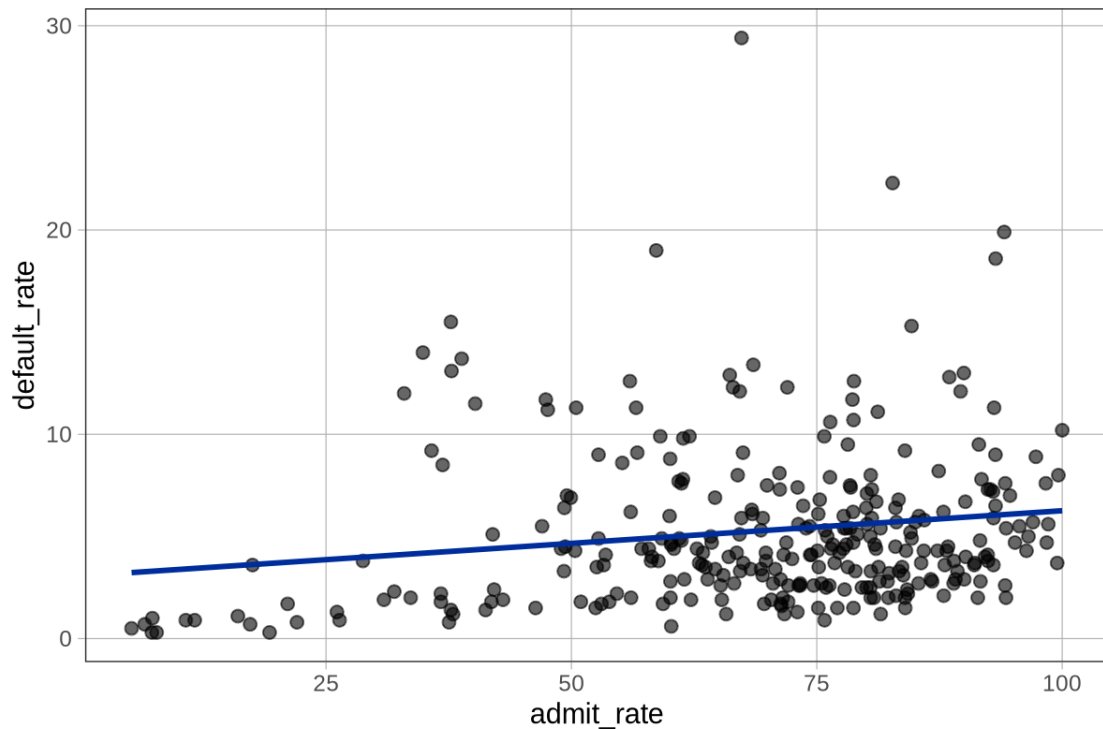
```
[1] "admit_rate"
```

```
[1] "SAT_avg"
```

```

[1] "enrollment"
[1] "net_price"
[1] "avg_cost"
[1] "net_tuition"
[1] "ed_spending_per_studen,\n avg_faculty_salary"
[1] "pctPELL"
[1] "pct_fed_loan"
[1] "grad_rate"
[1] "pct_firstgen"
[1] "med_fam_income"
[1] "med_alum_earnings"

```



1.4 - Would you say that this model provides a “good” fit for this dataset? Explain.

Double-click this cell to type your answer here: This model appears to provide a good fit for this dataset, as most of the points on the scatterplot are relatively close to the least squares regression line.

1.5 - Use the `lm` command to fit the linear regression model, where we use `SAT_avg` (predictor) to predict `default_rate` (outcome) in the dataset `sample_dat`. Store the model in a variable named `sat_model_1` and use the `summary` command to print out information about the model fit.

```

[330]: # Your code goes here
sat_model_1 <- lm(default_rate ~ 1/SAT_avg, data = sample_dat)
summary(sat_model_1)

```

```

Call:
lm(formula = default_rate ~ 1/SAT_avg, data = sample_dat)

Residuals:
    Min       1Q   Median       3Q      Max
-4.9743 -2.5993 -0.9743  1.5257 24.1257

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)    5.274      0.224   23.55  <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.88 on 299 degrees of freedom

```

Check yourself: The R^2 value shown in the model summary should be 0.4571

1.6 - Does the model's R^2 value indicate that this model provides a strong fit for this dataset? Explain.

Double-click this cell to type your answer here: The model's R^2 value indicates that this model does not provide a strong fit for this dataset because it does not explain much of the variation in the dependent variable

1.7 - If this model were curved, rather than linear, do you believe the R^2 could be higher? Explain.

Double-click this cell to type your answer here: If the data has a curved relationship, then the value of R^2 could be higher.

0.1.4 2.0 - Polynomial regression

Recall that simple linear regression follows this formula:

$$\hat{y} = \beta_0 + \beta_1 x$$

Where:

- β_0 is the intercept
- β_1 is the slope (coefficient of x)
- \hat{y} is the predicted `default_rate`
- x is the value of `SAT_avg`

If we want to capture the curvature in a scatter plot by creating a non-linear model, we can use a technique called **polynomial regression**. For example, we could use a degree 2 polynomial (quadratic), which looks like this:

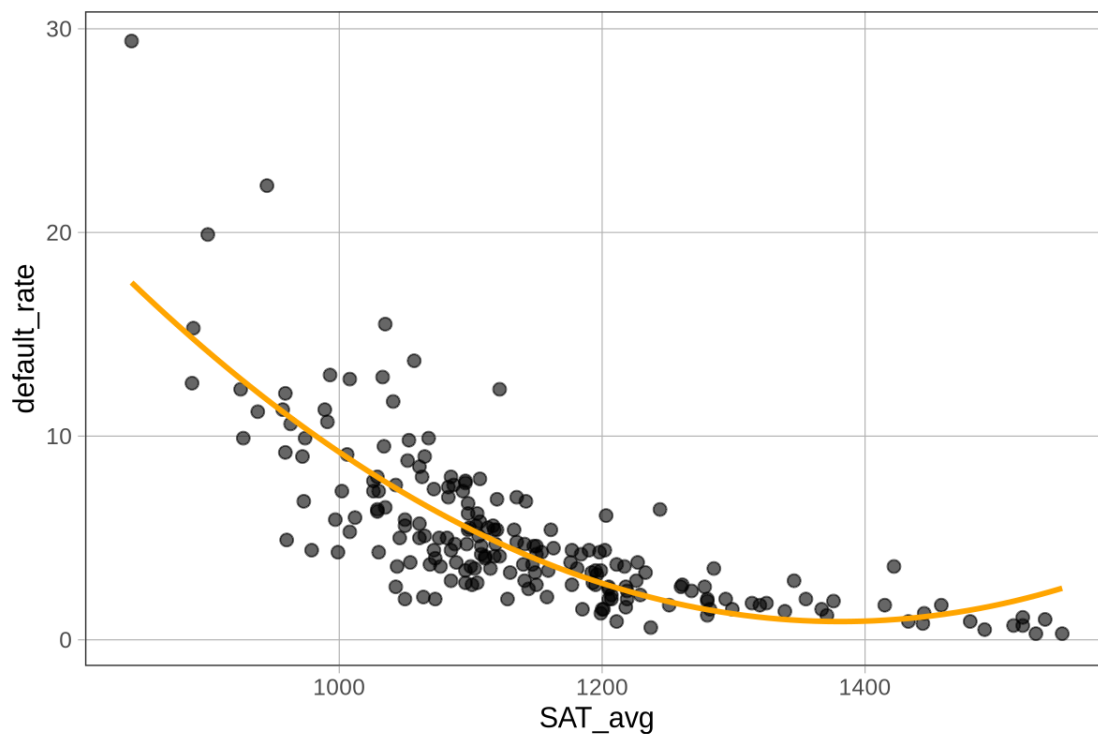
$$\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2$$

Where:

- β_0 is the intercept
- β_1 is the coefficient of x (linear term)
- β_2 is the coefficient of x^2 (squared term)
- \hat{y} is the predicted `default_rate`
- x is the `SAT_avg`

Below, we visualize the fit of this degree-2 polynomial (quadratic) model between `SAT_avg` and `default_rate`:

```
[106]: ## Run this code but do not edit it
# create scatterplot: default_rate ~ SAT_avg, with degree 2 polynomial model
# overlaid
gf_point(default_rate ~ SAT_avg, data = sample_dat) %>% gf_lm(formula = y ~
# poly(x, 2), color = "orange")
```



```
[274]: for (col in
list("admit_rate",
"SAT_avg",
"enrollment",
"net_price",
"avg_cost",
```

```

"net_tuition",
"ed_spending_per_studen,
avg_faculty_salary",
"pctPELL",
"pct_fed_loan",
"grad_rate",
"pct_firstgen",
"med_fam_income",
"med_alum_earnings")) {
  print(col)
  #gf_point(default_rate ~ sample_data$col, data = sample_dat) %>%
  ↪gf_lm(default_rate ~ sample_data$col, data = sample_dat)
}

```

```

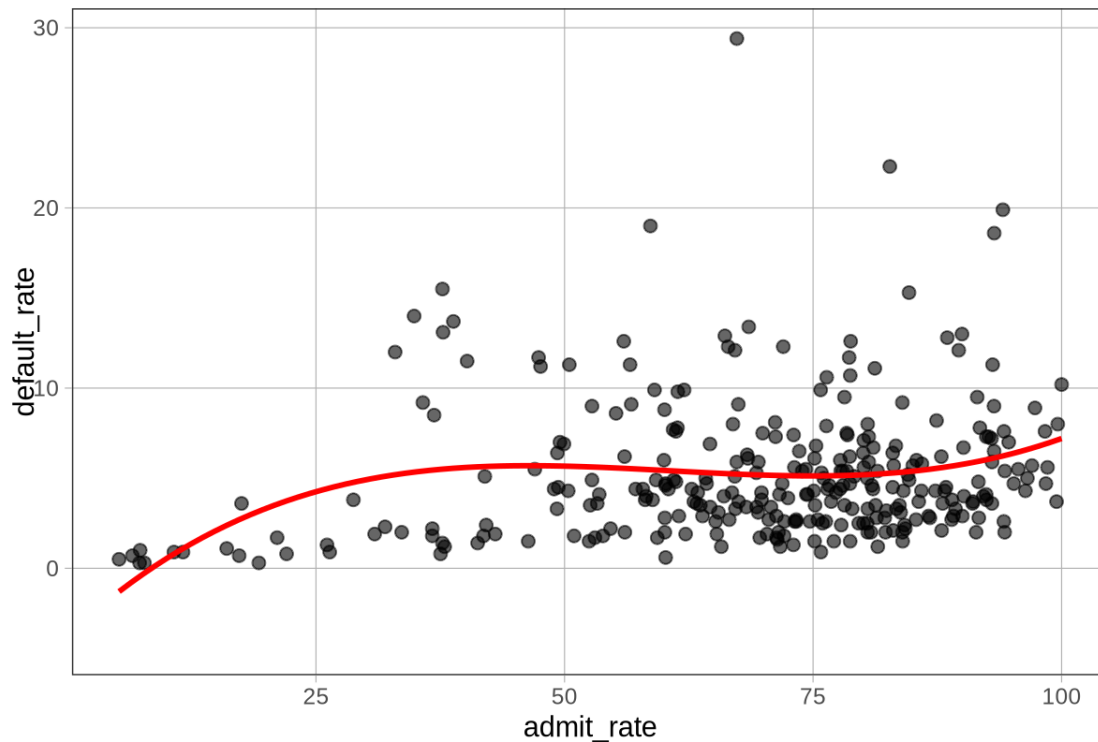
[1] "admit_rate"
[1] "SAT_avg"
[1] "enrollment"
[1] "net_price"
[1] "avg_cost"
[1] "net_tuition"
[1] "ed_spending_per_studen,\n avg_faculty_salary"
[1] "pctPELL"
[1] "pct_fed_loan"
[1] "grad_rate"
[1] "pct_firstgen"
[1] "med_fam_income"
[1] "med_alum_earnings"

```

```

[343]: gf_point(default_rate ~ admit_rate, data = sample_dat) %>% gf_lm(formula = y ~
  ↪poly(x,3), color="red")

```



2.1 - Make a prediction: Will this polynomial regression model have a higher or lower R^2 value than the linear regression model? Justify your reasoning.

Double-click this cell to type your answer here: Since the scatterplot shows that there is a weak, negative, nonlinear relationship between SAT_avg and default_rate, this polynomial regression model will have a higher R^2 value than the linear regression model. If the data does have a curved relationship, then the value of R^2 could be higher.

Let's test your prediction. To do so, we'll first need to fit the polynomial model. We can fit a degree 2 polynomial to the data using the `poly()` function inside of the `lm()` function. Run the cell below to see how it's done.

```
[329]: ## Run this code but do not edit it
# degree 2 polynomial model for default_rate ~ SAT_avg
sat_model_2 <- lm(default_rate ~ poly(SAT_avg, 2), data = sample_dat)
sat_model_2
```

Call:

```
lm(formula = default_rate ~ poly(SAT_avg, 2), data = sample_dat)
```

Coefficients:

```
(Intercept)  poly(SAT_avg, 2)1  poly(SAT_avg, 2)2
      5.274             -45.859             22.134
```

```
[331]: sat_model_2 <- lm(default_rate ~ poly(1/SAT_avg, 2), data = sample_dat)
sat_model_2
```

Call:

```
lm(formula = default_rate ~ poly(1/SAT_avg, 2), data = sample_dat)
```

Coefficients:

(Intercept)	poly(1/SAT_avg, 2)1	poly(1/SAT_avg, 2)2
5.274	49.067	17.374

The equation for this model would be

$$\hat{y} = 4.065 - 8.391x + 4.355x^2$$

Where:

- $\beta_0 = 4.065$ is the intercept
- $\beta_1 = -8.391$ is the coefficient of x , the linear term
- $\beta_2 = 4.355$ is the coefficient of x^2 , the squared term
- \hat{y} is the predicted default_rate
- x is the SAT_avg

2.2 - Use the `summary` command on `sat_model_2` to see summary information about the quadratic model.

```
[109]: # Your code goes here
summary(sat_model_2)
```

Call:

```
lm(formula = default_rate ~ poly(1/SAT_avg, 2), data = sample_dat)
```

Residuals:

Min	1Q	Median	3Q	Max
-6.3686	-1.2063	-0.2513	0.9104	10.0588

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	5.1550	0.1618	31.861	< 2e-16 ***
poly(1/SAT_avg, 2)1	42.3880	2.2881	18.525	< 2e-16 ***
poly(1/SAT_avg, 2)2	17.5926	2.2881	7.689	6.91e-13 ***

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

Residual standard error: 2.288 on 197 degrees of freedom

Multiple R-squared: 0.6713, Adjusted R-squared: 0.6679
F-statistic: 201.1 on 2 and 197 DF, p-value: < 2.2e-16

Check yourself: The R^2 value shown in the model summary should be 0.5802

2.3 - How does this model's R^2 value compare to that of the linear model? Was your prediction right? Explain.

Double-click this cell to type your answer here: My prediction was right, as this model's R^2 value is greater than that of the linear model.

This analysis raises a natural question: Why stop at degree 2? By raising the degree, we can add more curves to our model, potentially better fitting the data! Let's visualize what happens when we increase the degree in our polynomial regression models.

Degree 3 Polynomial Model

$$\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$$

```
[110]: ## Run this code but do not edit it
# create scatterplot: default_rate ~ SAT_avg, with degree 3 polynomial model,
# overlaid
gf_point(default_rate ~ SAT_avg, data = sample_dat) %>% gf_lm(formula = y ~
# poly(x, 3), color = "orange") + ylim(-4,12)
```

Warning message:

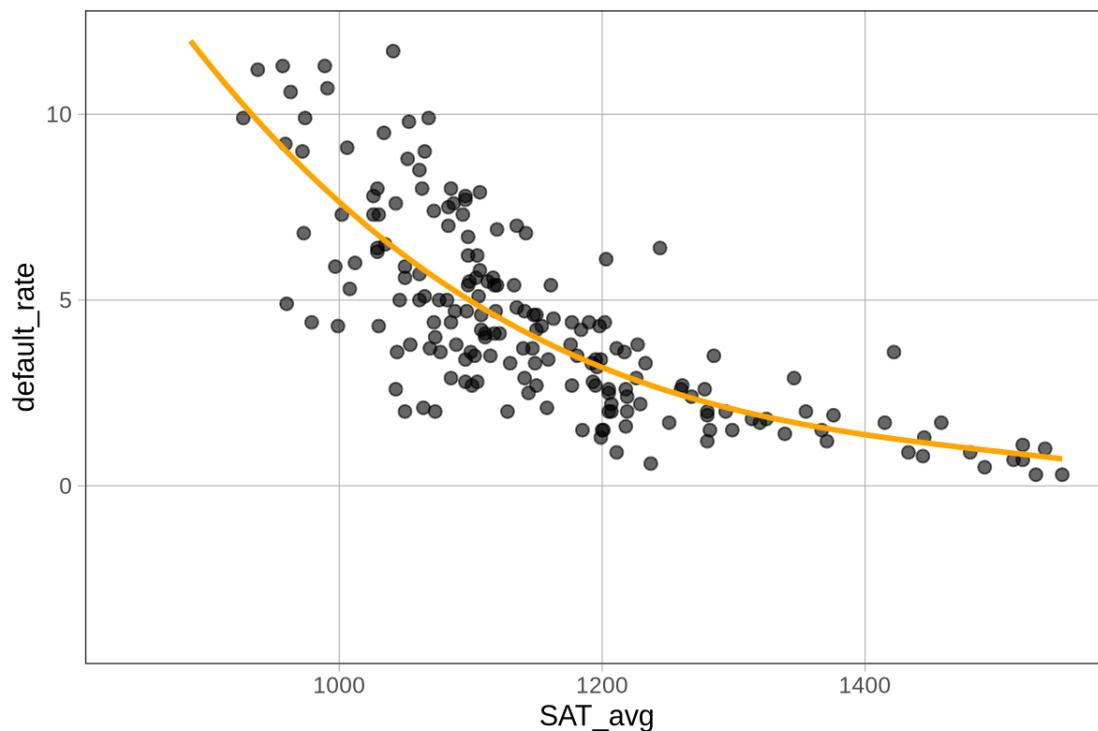
"Removed 13 rows containing non-finite outside the scale range
(`stat_lm()`)."

Warning message:

"Removed 13 rows containing missing values or values outside the scale
range
(`geom_point()`)."

Warning message:

"Removed 5 rows containing missing values or values outside the scale
range
(`geom_lm()`)."



Degree 5 Polynomial Model

$$\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 x^4 + \beta_5 x^5$$

```
[111]: ## Run this code but do not edit it
# create scatterplot: default_rate ~ SAT_avg, with degree 5 polynomial model
# overlaid
gf_point(default_rate ~ SAT_avg, data = sample_dat) %>% gf_lm(formula = y ~
  poly(x, 5), color = "orange") + ylim(-4,12)
```

Warning message:

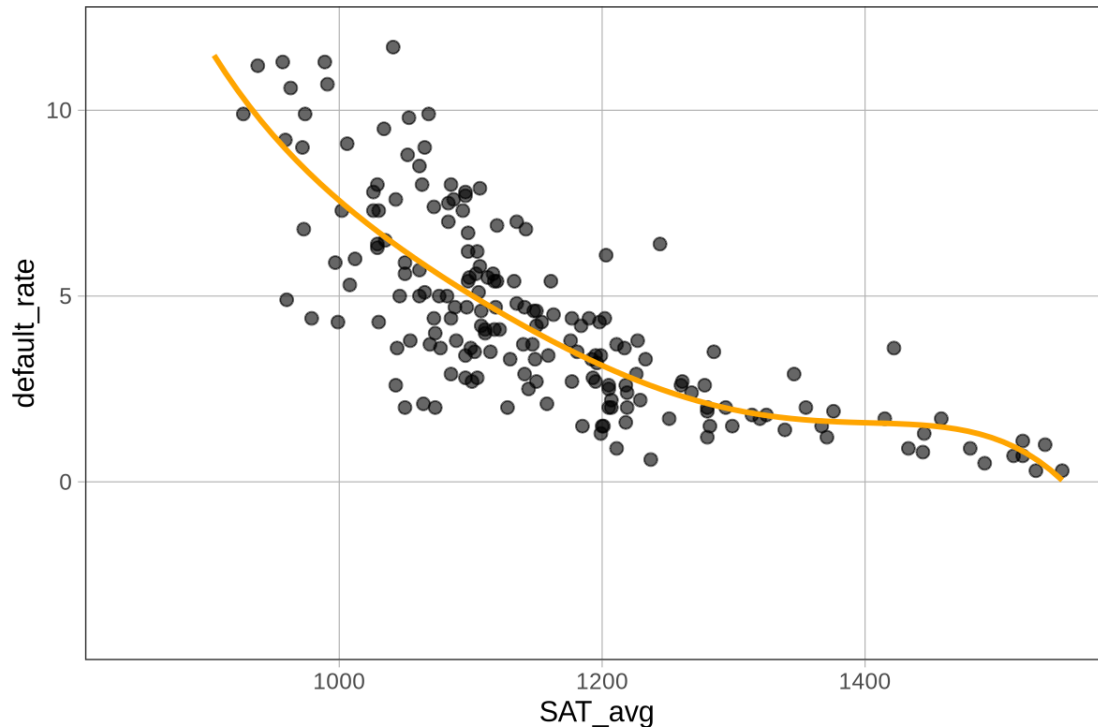
"Removed 13 rows containing non-finite outside the scale range
(`stat_lm()`)."

Warning message:

"Removed 13 rows containing missing values or values outside the scale
range
(`geom_point()`)."

Warning message:

"Removed 7 rows containing missing values or values outside the scale
range
(`geom_lm()`)."



Degree 12 Polynomial Model

$$\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 x^4 + \beta_5 x^5 + \beta_6 x^6 + \dots + \beta_{12} x^{12}$$

Note: The following code is pre-run, to save computer space.

```
[ ]: ## Note: This code was pre-run, to save computer space
# create scatterplot: default_rate ~ SAT_avg, with degree 12 polynomial model
# gf_point(default_rate ~ SAT_avg, data = sample_dat) %>% gf_smooth(method =
# "lm", formula = y ~poly(x,12), color = "orange") + ylim(-4,14)
```

2.4 - Examine each plot for the polynomial models with degrees 3, 5, 12. Which model do you think would have the largest R^2 value? Why?

Double-click this cell to type your answer here: The polynomial model with degree 12, since the points on the scatterplot are closest to the LSRL for this model than with the previous 2 models.

To determine which polynomial model fits the data the best, we will fit models for each degree (3, 5, 12).

```
[328]: ## Run this code but do not edit it
# degree 3, 5, and 12 polynomial models for default_rate ~ SAT_avg
sat_model_3 <- lm(default_rate ~ poly(SAT_avg, 3), data = sample_dat)
```

```
sat_model_5 <- lm(default_rate ~ poly(SAT_avg, 5), data = sample_dat)
sat_model_12 <- lm(default_rate ~ poly(SAT_avg, 12), data = sample_dat)
```

Now we can compare each model's R^2 value. Normally, we use the `summary` command and read the R^2 value. However, since we've fit so many models, we don't want to print out the entire summary for each one.

Instead, we'll use commands like this: `summary(sat_model_1)$r.squared`. The `$` operator is used to extract just the `r.squared` element from the full `summary`. We execute this command for each model, then print the results for ease of comparison.

```
[40]: ## Run this code but do not edit it
# r-squared value for each model
r2_sat_model_1 <- summary(sat_model_1)$r.squared
r2_sat_model_2 <- summary(sat_model_2)$r.squared
r2_sat_model_3 <- summary(sat_model_3)$r.squared
r2_sat_model_5 <- summary(sat_model_5)$r.squared
r2_sat_model_12 <- summary(sat_model_12)$r.squared

# print each model's r-squared value
print(paste("The R squared value for the degree 1 model is", r2_sat_model_1))
print(paste("The R squared value for the degree 2 model is", r2_sat_model_2))
print(paste("The R squared value for the degree 3 model is", r2_sat_model_3))
print(paste("The R squared value for the degree 5 model is", r2_sat_model_5))
print(paste("The R squared value for the degree 12 model is", r2_sat_model_12))
```

```
[1] "The R squared value for the degree 1 model is 0.457055427196517"
[1] "The R squared value for the degree 2 model is 0.580193597490376"
[1] "The R squared value for the degree 3 model is 0.60314009577391"
[1] "The R squared value for the degree 5 model is 0.647445002110733"
[1] "The R squared value for the degree 12 model is 0.775449820710613"
```

Check yourself: The R^2 for the degree 5 model should be about 0.647

2.5 - The degree 12 model has the highest R^2 value. Does that mean it's the “best” model? Why or why not?

Hint: Think about which model would do the best for predicting *the rest of the data* from the original full dataset.

Double-click this cell to type your answer here: A high degree polynomial like this degree 12 one is prone to overfitting because it has memorized the training data and not the general relationship between the variables.

0.1.5 3.0 - Prediction, model tuning, & machine learning

In prior notebooks, we've used our models to make inferences about default rates. However, sometimes in data science, we care more about predictions than we do about inferences. In particular, many data science tasks ask for making accurate predictions on *new* data - data that hadn't yet been collected when we first fit the model. This process of building models to predict new data, especially when it's automated, is called **machine learning**.

The key to machine learning is building models that make accurate predictions on **test** data - unseen data that weren't used when fitting the model. Let's see how this works. First, let's create a test dataset of 10 randomly sampled colleges. Importantly, these are colleges that **our models didn't see while fitting**:

```
[41]: ## Run this code but do not edit it
# create a data set to test the model with 10 new, randomly selected
# observations
# not used to train the model
set.seed(23)
test_dat <- sample(dat, size = 10)
```

3.1 - Use the `head` command on the `test_dat` data set.

```
[42]: # Your code goes here
head(test_dat)
```

		OPEID <int>	name <chr>	city <chr>	state <chr>	reg <chr>
A data.frame: 6 × 27	925	364600	Texas Woman's University	Denton	TX	Ro
	284	1025600	Benedictine College	Atchison	KS	Mi
	456	244900	Avila University	Kansas City	MO	Mi
	615	291400	Catawba College	Salisbury	NC	Sou
	913	363200	Texas A & M University-College Station	College Station	TX	Ro
	1015	2136600	Wisconsin Lutheran College	Milwaukee	WI	Mi

The following code visualizes the new test data alongside the training data (the data we used to originally fit our models).

```
[ ]: ## Run this code but do not edit it
# label train and test sets
sample_dat$phase <- "train"
test_dat$phase <- "test"

# concatenate two datasets
full_dat <- rbind(sample_dat, test_dat)

# create scatterplot: default_rate ~ SAT_avg, with degree 5 polynomial model
# overlaid
gf_point(default_rate ~ SAT_avg, data = full_dat, color = ~phase, shape = ~phase)
```

3.2 - Of all the polynomial models we fit before, which do you think would do best in predicting the default rates in the test dataset?

Note: Use your gut and intuition here. No calculations required.

Double-click this cell to type your answer here: A polynomial model with a degree of 4.

Let's see how good one of our models is at predicting default rates. The R code in the next cell uses the `predict` function to make predictions on the test dataset. In this case, output shows the

predicted default rates for the 10 test set colleges, as predicted by our degree 5 model.

```
[45]: ## Run this code but do not edit it
# get predictions for degree 5 model
pred_deg5 <- predict(sat_model_5, newdata = data.frame(SAT_avg =
  ↪test_dat$SAT_avg))
pred_deg5
```

```
1    4.71195777211723 2    2.80117510935355 3    4.76610631136996 4    5.83135112143679 5
1.1525539865398 6    3.41638983019025 7    9.85190742866518 8    18.1060265766056 9
1.29744229099948 10   4.28267983941371
```

So, how can we interpret these values? Well, the last college in our test set is **University of New Orleans**, which has an `SAT_avg` value of 1088 and a `default_rate` of 6.6. It's shown here on the graph, alongside our degree 5 model.

Note: The following code is pre-run.

```
[46]: ### Run this code but do not edit it
## create scatterplot: default_rate ~ SAT_avg, with degree 5 polynomial model
  ↪overlaid
#gf_point(default_rate ~ SAT_avg, data = sample_dat, color = ~phase, shape =
  ↪~phase) %>% gf_lm(formula = y ~poly(x, 5), color = "orange") %>%
  ↪gf_point(default_rate ~ SAT_avg, data = full_dat, color = ~phase, shape =
  ↪~phase) + ylim(0,19)
```

Our degree 5 model's predicted default rate for this first data point was 4.28. That means that our degree 5 model under-estimates the actual value for default rate by...

$$6.6 - 4.28 = 2.32$$

The model's prediction and error is visualized in the plot below:

So, its predicted default rate is “off” by about 2 percentage points! This is pretty amazing, considering the model had only 20 training data values and the **University of New Orleans** wasn't included among them. This is the power of machine learning! Predicting previously *unseen* data!

This is just one prediction. We're really interested in how this model performed across all its predictions. For that, let's measure its R^2 (prediction strength) on the test set!

We can use the `cor` function to correlate the predictions with the actual default rates (r) and then square that value to get R^2 , which gets us the prediction strength!

```
[47]: ## Run this code but do not edit it
# Get correlation between predicted and actual default rates in test set
cor(test_dat$default_rate, pred_deg5) ^ 2
```

```
0.616546630372038
```

We can now repeat this same process for all polynomial degrees.

```
[48]: ## Run this code but do not edit it
# Storing test set predictions for all models
pred_deg1 <- predict(sat_model_1, newdata = data.frame(SAT_avg =
  ↪test_dat$SAT_avg))
pred_deg2 <- predict(sat_model_2, newdata = data.frame(SAT_avg =
  ↪test_dat$SAT_avg))
pred_deg3 <- predict(sat_model_3, newdata = data.frame(SAT_avg =
  ↪test_dat$SAT_avg))
pred_deg5 <- predict(sat_model_5, newdata = data.frame(SAT_avg =
  ↪test_dat$SAT_avg))
pred_deg12 <- predict(sat_model_12, newdata = data.frame(SAT_avg =
  ↪test_dat$SAT_avg))

# print each model's r-squared value
print(paste("The test R squared value for the degree 1 model is",
  ↪cor(test_dat$default_rate, pred_deg1) ^ 2))
print(paste("The test R squared value for the degree 2 model is",
  ↪cor(test_dat$default_rate, pred_deg2) ^ 2))
print(paste("The test R squared value for the degree 3 model is",
  ↪cor(test_dat$default_rate, pred_deg3) ^ 2))
print(paste("The test R squared value for the degree 5 model is",
  ↪cor(test_dat$default_rate, pred_deg5) ^ 2))
print(paste("The test R squared value for the degree 12 model is",
  ↪cor(test_dat$default_rate, pred_deg12) ^ 2))
```

```
[1] "The test R squared value for the degree 1 model is 0.55824446697698"
[1] "The test R squared value for the degree 2 model is 0.70025122602337"
[1] "The test R squared value for the degree 3 model is 0.733729012084851"
[1] "The test R squared value for the degree 5 model is 0.616546630372038"
[1] "The test R squared value for the degree 12 model is 0.176121561427524"
```

```
[62]: set.seed(23)
test_dat <- sample(dat, size = 100)
```

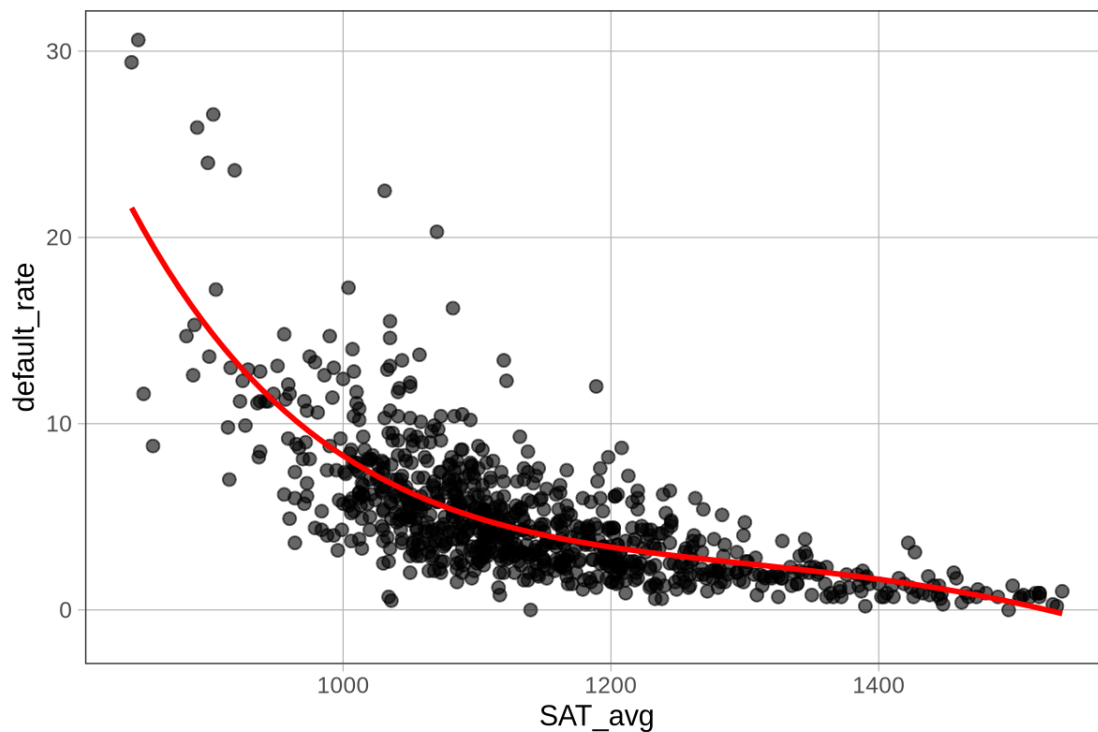
```
[ ]: ## Run but do not edit this code

# set training data to be 80% of all colleges
train_size <- floor(0.8 * nrow(dat))

## sample row indeces
set.seed(2024)
train_ind <- sample(seq_len(nrow(dat)), size = train_size)

sample_dat <- dat[train_ind, ]
test_dat <- dat[-train_ind, ]
```

```
[8]: gf_point(default_rate ~ SAT_avg, data = sample_dat) %>% gf_lm(formula = y ~  
      ↪poly(log(x),3), color="red")
```

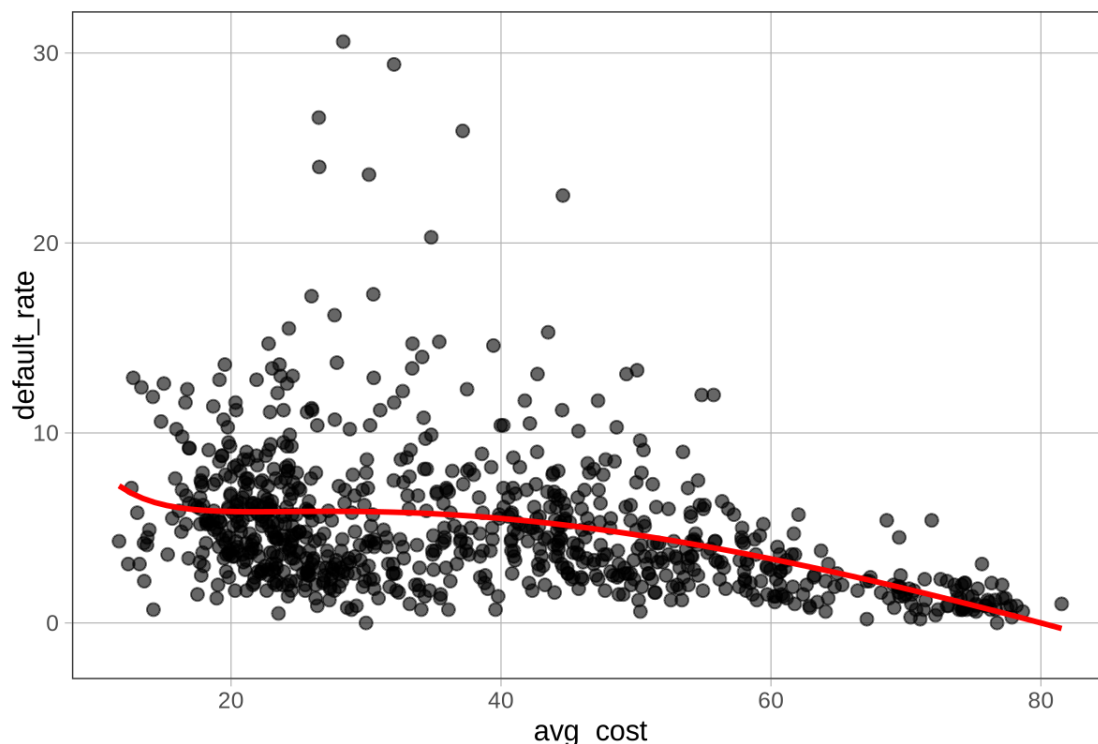


```
[9]: model_1 <- lm(default_rate ~ poly(SAT_avg, 1), data = sample_dat)
model_2 <- lm(default_rate ~ poly(SAT_avg, 2), data = sample_dat)
model_3 <- lm(default_rate ~ poly(SAT_avg, 3), data = sample_dat)
pred_deg1 <- predict(model_1, newdata = data.frame(SAT_avg = test_dat$SAT_avg))
pred_deg2 <- predict(model_2, newdata = data.frame(SAT_avg = test_dat$SAT_avg))
pred_deg3 <- predict(model_3, newdata = data.frame(SAT_avg = test_dat$SAT_avg))
print(paste("The test R squared value for the degree 1 model is",  
      ↪cor(test_dat$default_rate, pred_deg1) ^ 2))
print(paste("The test R squared value for the degree 2 model is",  
      ↪cor(test_dat$default_rate, pred_deg2) ^ 2))
print(paste("The test R squared value for the degree 3 model is",  
      ↪cor(test_dat$default_rate, pred_deg3) ^ 2))
print(paste("The test R value for the degree 1 model is",  
      ↪cor(test_dat$default_rate, pred_deg1)))
print(paste("The test R value for the degree 2 model is",  
      ↪cor(test_dat$default_rate, pred_deg2)))
print(paste("The test R value for the degree 3 model is",  
      ↪cor(test_dat$default_rate, pred_deg3)))
#summary(model_1)
#summary(model_2)
```

```
#summary(model_3)
```

```
[1] "The test R squared value for the degree 1 model is 0.43699889591149"  
[1] "The test R squared value for the degree 2 model is 0.531450285931865"  
[1] "The test R squared value for the degree 3 model is 0.53251552552625"  
[1] "The test R value for the degree 1 model is 0.661058920151214"  
[1] "The test R value for the degree 2 model is 0.729006368924076"  
[1] "The test R value for the degree 3 model is 0.729736613804083"
```

```
[158]: gf_point(default_rate ~ avg_cost, data = sample_dat) %>% gf_lm(formula = y ~  
  ↪poly(log(x),3), color="red")
```



```
[159]: model_1 <- lm(default_rate ~ poly(avg_cost, 1), data = sample_dat)  
model_2 <- lm(default_rate ~ poly(avg_cost, 2), data = sample_dat)  
model_3 <- lm(default_rate ~ poly(log(avg_cost), 3), data = sample_dat)  
pred_deg1 <- predict(model_1, newdata = data.frame(avg_cost =  
  ↪test_dat$avg_cost))  
pred_deg2 <- predict(model_2, newdata = data.frame(avg_cost =  
  ↪test_dat$avg_cost))  
pred_deg3 <- predict(model_3, newdata = data.frame(avg_cost =  
  ↪test_dat$avg_cost))  
print(paste("The test R squared value for the degree 1 model is",  
  ↪cor(test_dat$default_rate, pred_deg1) ^ 2))
```

```

print(paste("The test R squared value for the degree 2 model is",
  ↪cor(test_dat$default_rate, pred_deg2) ^ 2))
print(paste("The test R squared value for the degree 3 model is",
  ↪cor(test_dat$default_rate, pred_deg3) ^ 2))
print(paste("The test R value for the degree 1 model is",
  ↪cor(test_dat$default_rate, pred_deg1)))
print(paste("The test R value for the degree 2 model is",
  ↪cor(test_dat$default_rate, pred_deg2)))
print(paste("The test R value for the degree 3 model is",
  ↪cor(test_dat$default_rate, pred_deg3)))

```

```

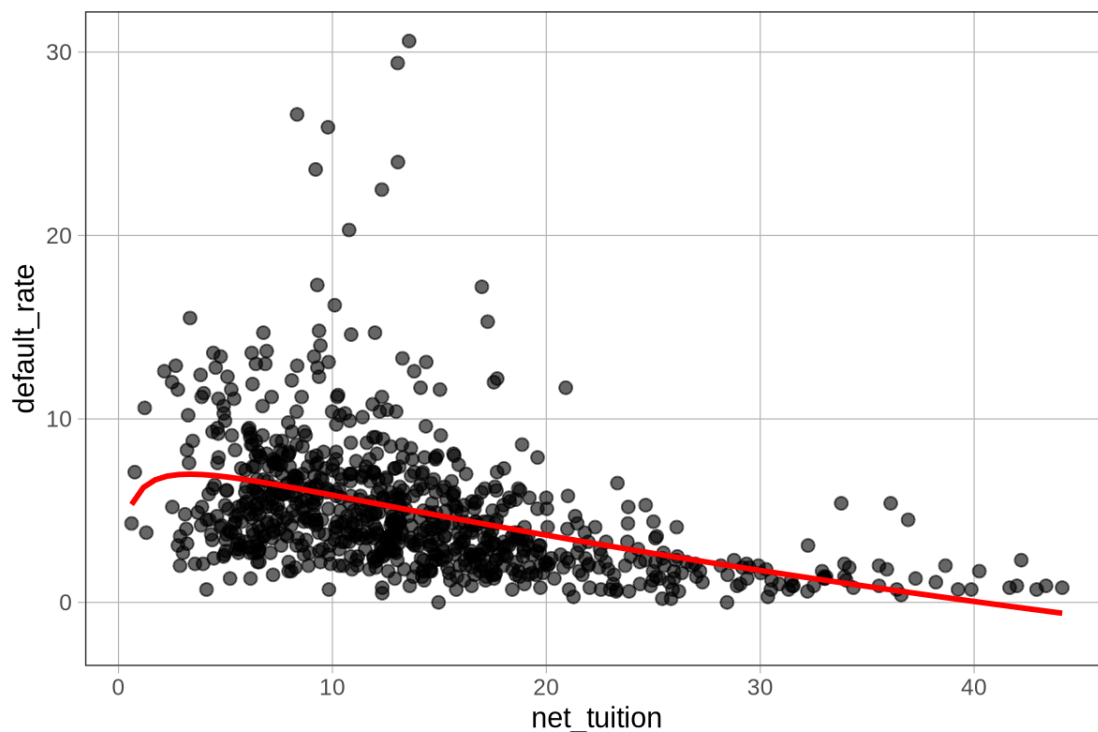
[1] "The test R squared value for the degree 1 model is 0.177023675341143"
[1] "The test R squared value for the degree 2 model is 0.196225709398749"
[1] "The test R squared value for the degree 3 model is 0.196094894541057"
[1] "The test R value for the degree 1 model is 0.420741815536729"
[1] "The test R value for the degree 2 model is 0.442973711859687"
[1] "The test R value for the degree 3 model is 0.442826031914404"

```

```

[205]: gf_point(default_rate ~ net_tuition, data = sample_dat) %>% gf_lm(formula = y ~
  ↪poly(log(x),3), color="red")

```



```

[204]: model_1 <- lm(default_rate ~ poly(log(net_tuition), 2), data = sample_dat)
model_2 <- lm(default_rate ~ poly(net_tuition, 2), data = sample_dat)

```

```

model_3 <- lm(default_rate ~ poly(net_tuition, 4), data = sample_dat)
pred_deg1 <- predict(model_1, newdata = data.frame(net_tuition =
  ↪test_dat$net_tuition))
pred_deg2 <- predict(model_2, newdata = data.frame(net_tuition =
  ↪test_dat$net_tuition))
pred_deg3 <- predict(model_3, newdata = data.frame(net_tuition =
  ↪test_dat$net_tuition))
print(paste("The test R squared value for the degree 1 model is",
  ↪cor(test_dat$default_rate, pred_deg1) ^ 2))
print(paste("The test R squared value for the degree 2 model is",
  ↪cor(test_dat$default_rate, pred_deg2) ^ 2))
print(paste("The test R squared value for the degree 3 model is",
  ↪cor(test_dat$default_rate, pred_deg3) ^ 2))
print(paste("The test R value for the degree 1 model is",
  ↪cor(test_dat$default_rate, pred_deg1)))
print(paste("The test R value for the degree 2 model is",
  ↪cor(test_dat$default_rate, pred_deg2)))
print(paste("The test R value for the degree 3 model is",
  ↪cor(test_dat$default_rate, pred_deg3)))

```

```

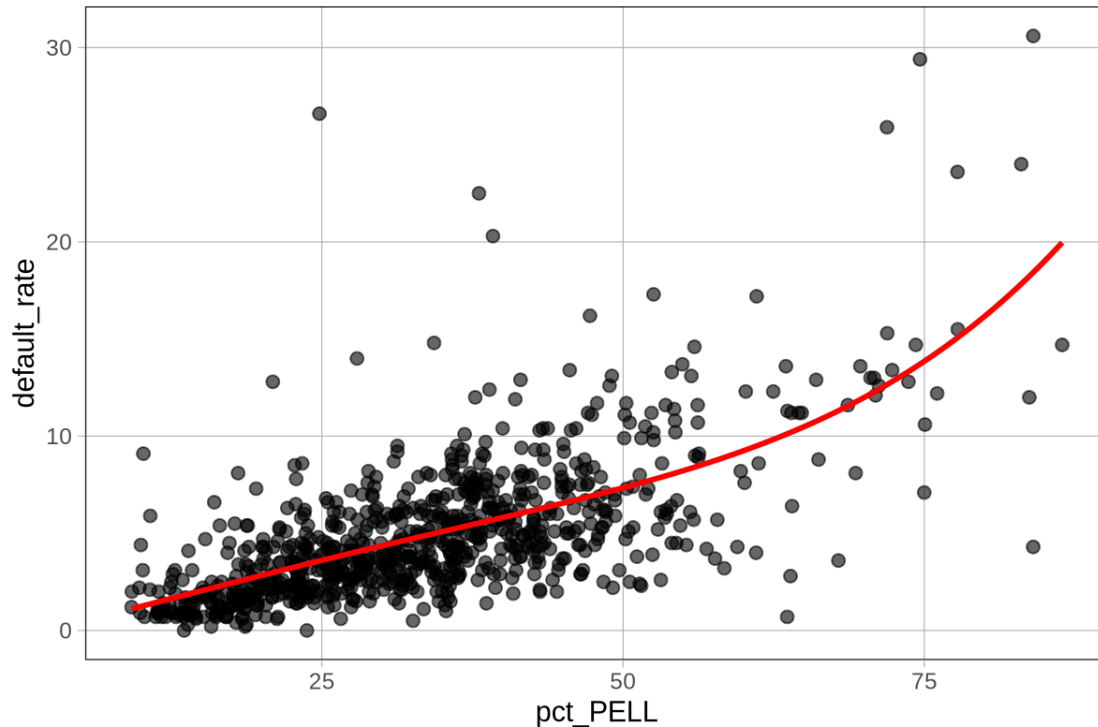
[1] "The test R squared value for the degree 1 model is 0.255538754081114"
[1] "The test R squared value for the degree 2 model is 0.243661705531134"
[1] "The test R squared value for the degree 3 model is 0.251670597654293"
[1] "The test R value for the degree 1 model is 0.505508411484037"
[1] "The test R value for the degree 2 model is 0.493621014069634"
[1] "The test R value for the degree 3 model is 0.501667816043936"

```

```

[114]: gf_point(default_rate ~ pctPELL, data = sample_dat) %>% gf_lm(formula = y ~
  ↪poly(x,4), color="red")

```



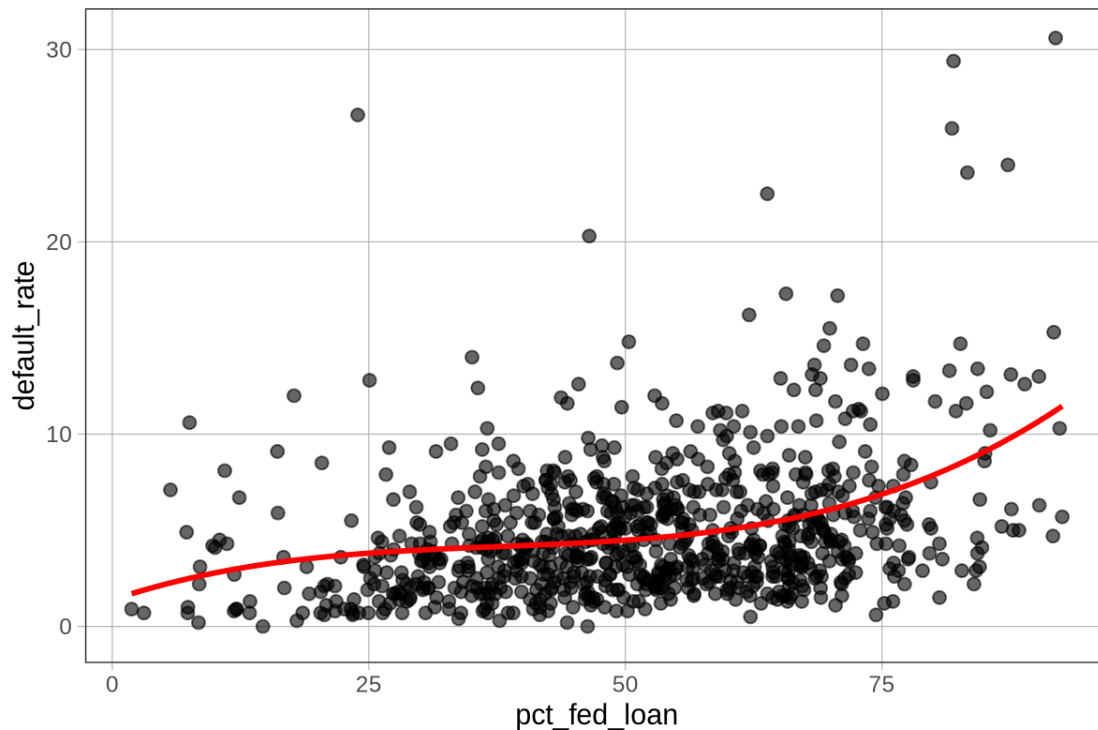
```
[7]: model_1 <- lm(default_rate ~ poly(pctPELL, 1), data = sample_dat)
model_2 <- lm(default_rate ~ poly(pctPELL, 2), data = sample_dat)
model_3 <- lm(default_rate ~ poly(pctPELL, 4), data = sample_dat)
pred_deg1 <- predict(model_1, newdata = data.frame(pctPELL =
  ↳ test_dat$pctPELL))
pred_deg2 <- predict(model_2, newdata = data.frame(pctPELL =
  ↳ test_dat$pctPELL))
pred_deg3 <- predict(model_3, newdata = data.frame(pctPELL =
  ↳ test_dat$pctPELL))
print(paste("The test R squared value for the degree 1 model is",
  ↳ cor(test_dat$default_rate, pred_deg1) ^ 2))
print(paste("The test R squared value for the degree 2 model is",
  ↳ cor(test_dat$default_rate, pred_deg2) ^ 2))
print(paste("The test R squared value for the degree 3 model is",
  ↳ cor(test_dat$default_rate, pred_deg3) ^ 2))
print(paste("The test R value for the degree 1 model is",
  ↳ cor(test_dat$default_rate, pred_deg1)))
print(paste("The test R value for the degree 2 model is",
  ↳ cor(test_dat$default_rate, pred_deg2)))
print(paste("The test R value for the degree 3 model is",
  ↳ cor(test_dat$default_rate, pred_deg3)))
```

```
[1] "The test R squared value for the degree 1 model is 0.45615433624854"
```



```
[1] "The test R squared value for the degree 2 model is 0.472725609178097"
[1] "The test R squared value for the degree 3 model is 0.486800870614671"
[1] "The test R value for the degree 1 model is 0.675391987107146"
[1] "The test R value for the degree 2 model is 0.687550441188206"
[1] "The test R value for the degree 3 model is 0.697711165608428"
```

```
[115]: gf_point(default_rate ~ pct_fed_loan, data = sample_dat) %>% gf_lm(formula = y ~
  ↪~ poly(x,3), color="red")
```



```
[123]: model_1 <- lm(default_rate ~ poly(pct_fed_loan, 1), data = sample_dat)
model_2 <- lm(default_rate ~ poly(pct_fed_loan, 2), data = sample_dat)
model_3 <- lm(default_rate ~ poly(log(pct_fed_loan), 3), data = sample_dat)
pred_deg1 <- predict(model_1, newdata = data.frame(pct_fed_loan =
  ↪test_dat$pct_fed_loan))
pred_deg2 <- predict(model_2, newdata = data.frame(pct_fed_loan =
  ↪test_dat$pct_fed_loan))
pred_deg3 <- predict(model_3, newdata = data.frame(pct_fed_loan =
  ↪test_dat$pct_fed_loan))
print(paste("The test R squared value for the degree 1 model is",
  ↪cor(test_dat$default_rate, pred_deg1) ^ 2))
print(paste("The test R squared value for the degree 2 model is",
  ↪cor(test_dat$default_rate, pred_deg2) ^ 2))
```

```

print(paste("The test R squared value for the degree 3 model is",
  ↪cor(test_dat$default_rate, pred_deg3) ^ 2))
print(paste("The test R value for the degree 1 model is",
  ↪cor(test_dat$default_rate, pred_deg1)))
print(paste("The test R value for the degree 2 model is",
  ↪cor(test_dat$default_rate, pred_deg2)))
print(paste("The test R value for the degree 3 model is",
  ↪cor(test_dat$default_rate, pred_deg3)))

```

```

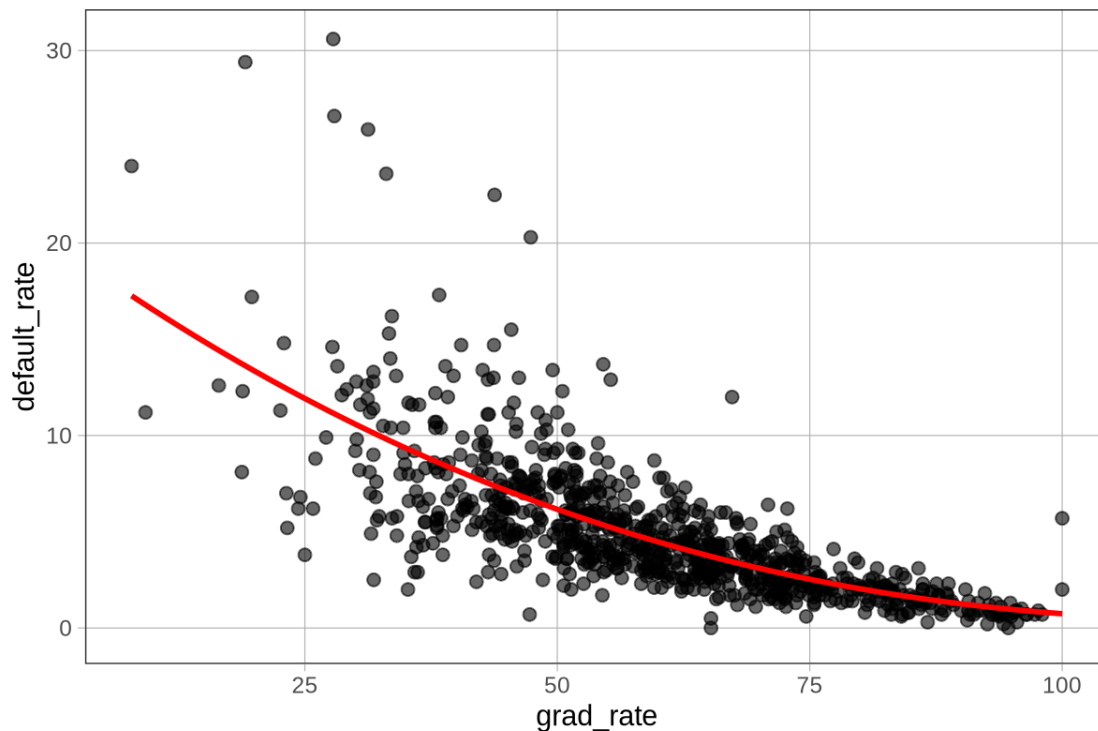
[1] "The test R squared value for the degree 1 model is 0.181931724114206"
[1] "The test R squared value for the degree 2 model is 0.194954156700041"
[1] "The test R squared value for the degree 3 model is 0.192449880191799"
[1] "The test R value for the degree 1 model is 0.426534552075452"
[1] "The test R value for the degree 2 model is 0.441536132949548"
[1] "The test R value for the degree 3 model is 0.438691098829004"

```

```

[156]: gf_point(default_rate ~ grad_rate, data = sample_dat) %>% gf_lm(formula = y ~
  ↪poly(x,3), color="red")

```



```

[164]: model_1 <- lm(default_rate ~ poly(1/grad_rate, 4), data = sample_dat)
model_2 <- lm(default_rate ~ poly(grad_rate, 2), data = sample_dat)
model_3 <- lm(default_rate ~ poly(grad_rate, 3), data = sample_dat)

```

```

pred_deg1 <- predict(model_1, newdata = data.frame(grad_rate =
  ↪test_dat$grad_rate))
pred_deg2 <- predict(model_2, newdata = data.frame(grad_rate =
  ↪test_dat$grad_rate))
pred_deg3 <- predict(model_3, newdata = data.frame(grad_rate =
  ↪test_dat$grad_rate))
print(paste("The test R squared value for the degree 1 model is",
  ↪cor(test_dat$default_rate, pred_deg1) ^ 2))
print(paste("The test R squared value for the degree 2 model is",
  ↪cor(test_dat$default_rate, pred_deg2) ^ 2))
print(paste("The test R squared value for the degree 3 model is",
  ↪cor(test_dat$default_rate, pred_deg3) ^ 2))
print(paste("The test R value for the degree 1 model is",
  ↪cor(test_dat$default_rate, pred_deg1)))
print(paste("The test R value for the degree 2 model is",
  ↪cor(test_dat$default_rate, pred_deg2)))
print(paste("The test R value for the degree 3 model is",
  ↪cor(test_dat$default_rate, pred_deg3)))

```

```

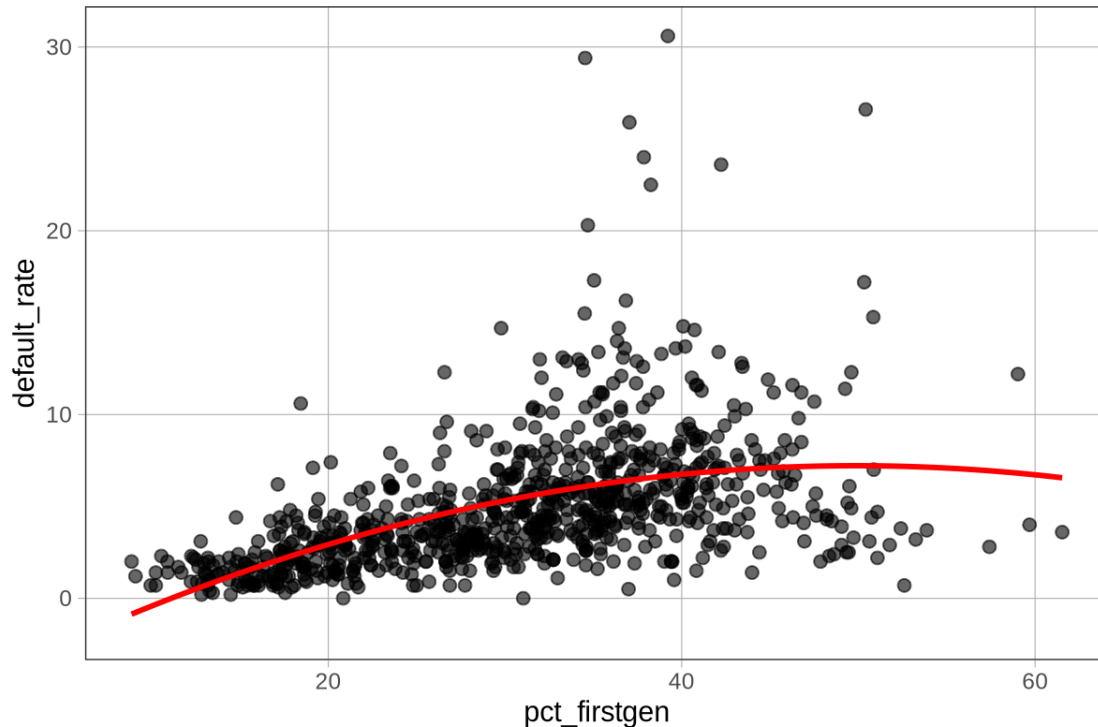
[1] "The test R squared value for the degree 1 model is 0.584460491123511"
[1] "The test R squared value for the degree 2 model is 0.57894791089609"
[1] "The test R squared value for the degree 3 model is 0.579011715674716"
[1] "The test R value for the degree 1 model is 0.764500157700122"
[1] "The test R value for the degree 2 model is 0.760886266728537"
[1] "The test R value for the degree 3 model is 0.760928193507584"

```

```

[169]: gf_point(default_rate ~ pct_firstgen, data = sample_dat) %>% gf_lm(formula = y
  ↪~ poly(x,2), color="red")

```

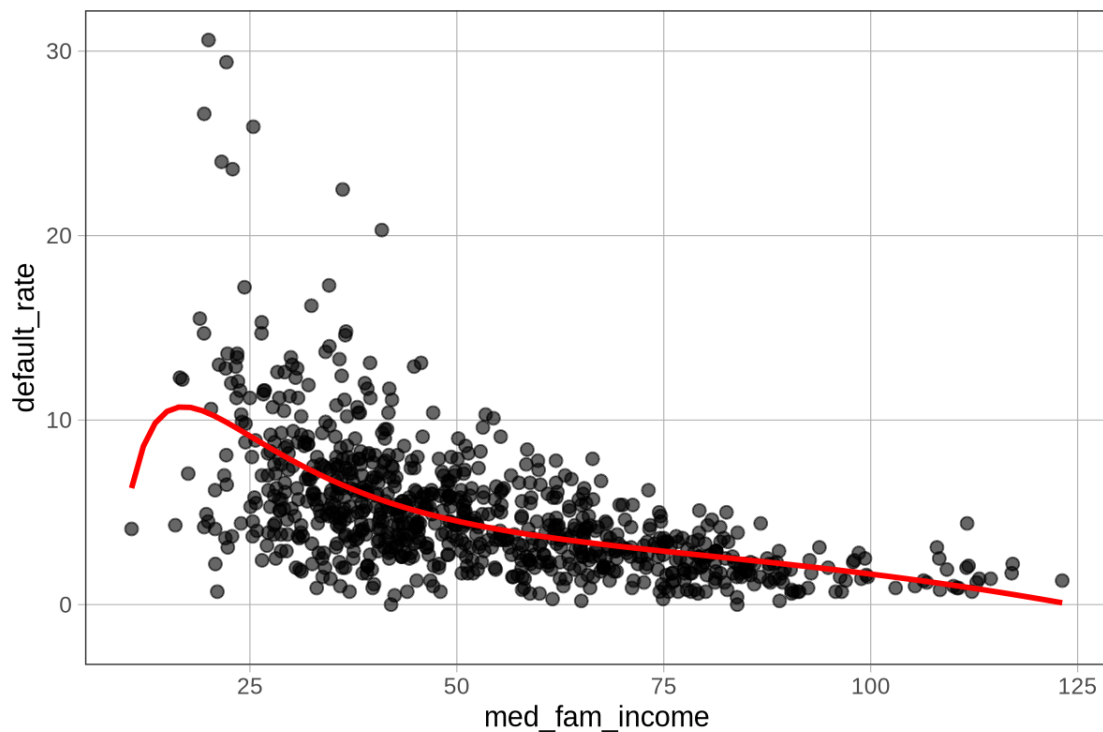


```
[165]: model_1 <- lm(default_rate ~ poly(pct_firstgen, 4), data = sample_dat)
model_2 <- lm(default_rate ~ poly(pct_firstgen, 2), data = sample_dat)
model_3 <- lm(default_rate ~ poly(pct_firstgen, 3), data = sample_dat)
pred_deg1 <- predict(model_1, newdata = data.frame(pct_firstgen =
  ↳ test_dat$pct_firstgen))
pred_deg2 <- predict(model_2, newdata = data.frame(pct_firstgen =
  ↳ test_dat$pct_firstgen))
pred_deg3 <- predict(model_3, newdata = data.frame(pct_firstgen =
  ↳ test_dat$pct_firstgen))
print(paste("The test R squared value for the degree 1 model is",
  ↳ cor(test_dat$default_rate, pred_deg1) ^ 2))
print(paste("The test R squared value for the degree 2 model is",
  ↳ cor(test_dat$default_rate, pred_deg2) ^ 2))
print(paste("The test R squared value for the degree 3 model is",
  ↳ cor(test_dat$default_rate, pred_deg3) ^ 2))
print(paste("The test R value for the degree 1 model is",
  ↳ cor(test_dat$default_rate, pred_deg1)))
print(paste("The test R value for the degree 2 model is",
  ↳ cor(test_dat$default_rate, pred_deg2)))
print(paste("The test R value for the degree 3 model is",
  ↳ cor(test_dat$default_rate, pred_deg3)))
```

```
[1] "The test R squared value for the degree 1 model is 0.22541383565914"
```

```
[1] "The test R squared value for the degree 2 model is 0.222707749063853"
[1] "The test R squared value for the degree 3 model is 0.225080490909004"
[1] "The test R value for the degree 1 model is 0.474777669714088"
[1] "The test R value for the degree 2 model is 0.471919218790518"
[1] "The test R value for the degree 3 model is 0.474426486306366"
```

```
[199]: gf_point(default_rate ~ med_fam_income, data = sample_dat) %>% gf_lm(formula =
  ↳ y ~ poly(log(x),4), color="red")
```



```
[202]: model_1 <- lm(default_rate ~ poly(med_fam_income, 1), data = sample_dat)
model_2 <- lm(default_rate ~ poly(med_fam_income, 2), data = sample_dat)
model_3 <- lm(default_rate ~ poly(med_fam_income, 4), data = sample_dat)
pred_deg1 <- predict(model_1, newdata = data.frame(med_fam_income =
  ↳ test_dat$med_fam_income))
pred_deg2 <- predict(model_2, newdata = data.frame(med_fam_income =
  ↳ test_dat$med_fam_income))
pred_deg3 <- predict(model_3, newdata = data.frame(med_fam_income =
  ↳ test_dat$med_fam_income))
print(paste("The test R squared value for the degree 1 model is",
  ↳ cor(test_dat$default_rate, pred_deg1) ^ 2))
print(paste("The test R squared value for the degree 2 model is",
  ↳ cor(test_dat$default_rate, pred_deg2) ^ 2))
```

```

print(paste("The test R squared value for the degree 3 model is",
  ↪cor(test_dat$default_rate, pred_deg3) ^ 2))
print(paste("The test R value for the degree 1 model is",
  ↪cor(test_dat$default_rate, pred_deg1)))
print(paste("The test R value for the degree 2 model is",
  ↪cor(test_dat$default_rate, pred_deg2)))
print(paste("The test R value for the degree 3 model is",
  ↪cor(test_dat$default_rate, pred_deg3)))

```

```

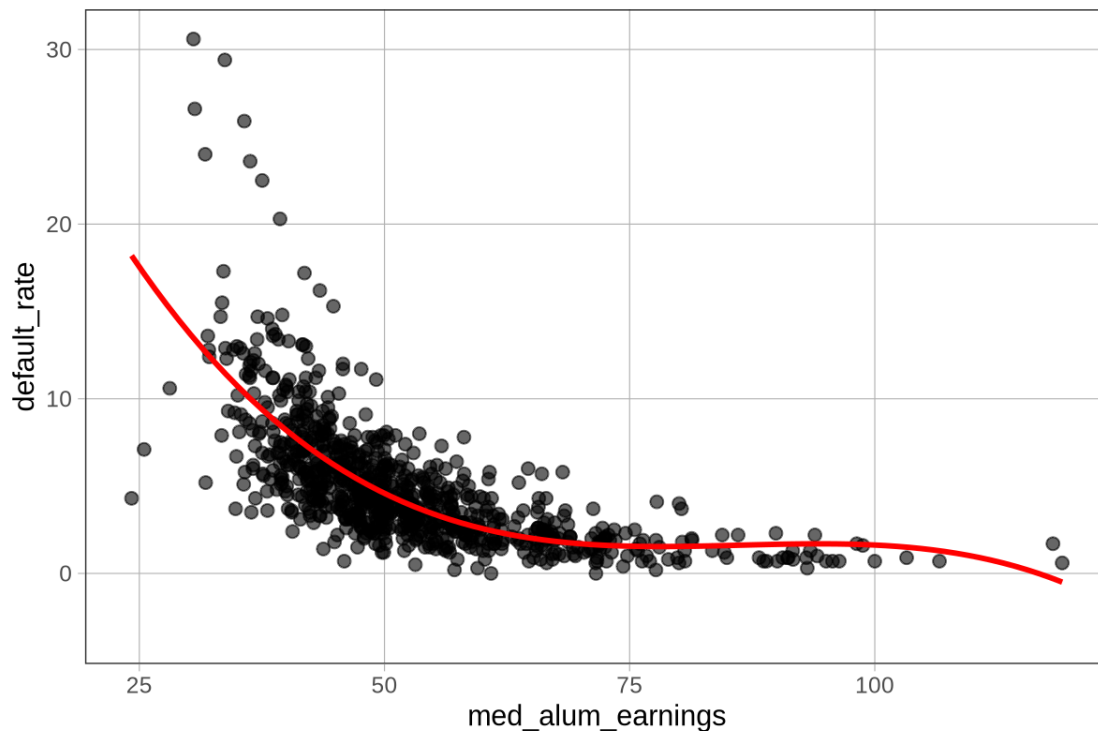
[1] "The test R squared value for the degree 1 model is 0.324789552646084"
[1] "The test R squared value for the degree 2 model is 0.366632791010403"
[1] "The test R squared value for the degree 3 model is 0.380295770600981"
[1] "The test R value for the degree 1 model is 0.569903108121094"
[1] "The test R value for the degree 2 model is 0.605502098270851"
[1] "The test R value for the degree 3 model is 0.616681255269674"

```

```

[11]: gf_point(default_rate ~ med_alum_earnings, data = sample_dat) %>% gf_lm(formula_
  ↪= y ~ poly(x,3), color="red")

```



```

[25]: model_1 <- lm(default_rate ~ poly(med_alum_earnings, 1), data = sample_dat)
model_2 <- lm(default_rate ~ poly(med_alum_earnings, 2), data = sample_dat)
model_3 <- lm(default_rate ~ poly(med_alum_earnings, 3), data = sample_dat)

```

```

pred_deg1 <- predict(model_1, newdata = data.frame(med_alum_earnings =
  ↪test_dat$med_alum_earnings))
pred_deg2 <- predict(model_2, newdata = data.frame(med_alum_earnings =
  ↪test_dat$med_alum_earnings))
pred_deg3 <- predict(model_3, newdata = data.frame(med_alum_earnings =
  ↪test_dat$med_alum_earnings))
print(paste("The test R squared value for the degree 1 model is",
  ↪cor(test_dat$default_rate, pred_deg1) ^ 2))
print(paste("The test R squared value for the degree 2 model is",
  ↪cor(test_dat$default_rate, pred_deg2) ^ 2))
print(paste("The test R squared value for the degree 3 model is",
  ↪cor(test_dat$default_rate, pred_deg3) ^ 2))
print(paste("The test R value for the degree 1 model is",
  ↪cor(test_dat$default_rate, pred_deg1)))
print(paste("The test R value for the degree 2 model is",
  ↪cor(test_dat$default_rate, pred_deg2)))
print(paste("The test R value for the degree 3 model is",
  ↪cor(test_dat$default_rate, pred_deg3)))

```

```

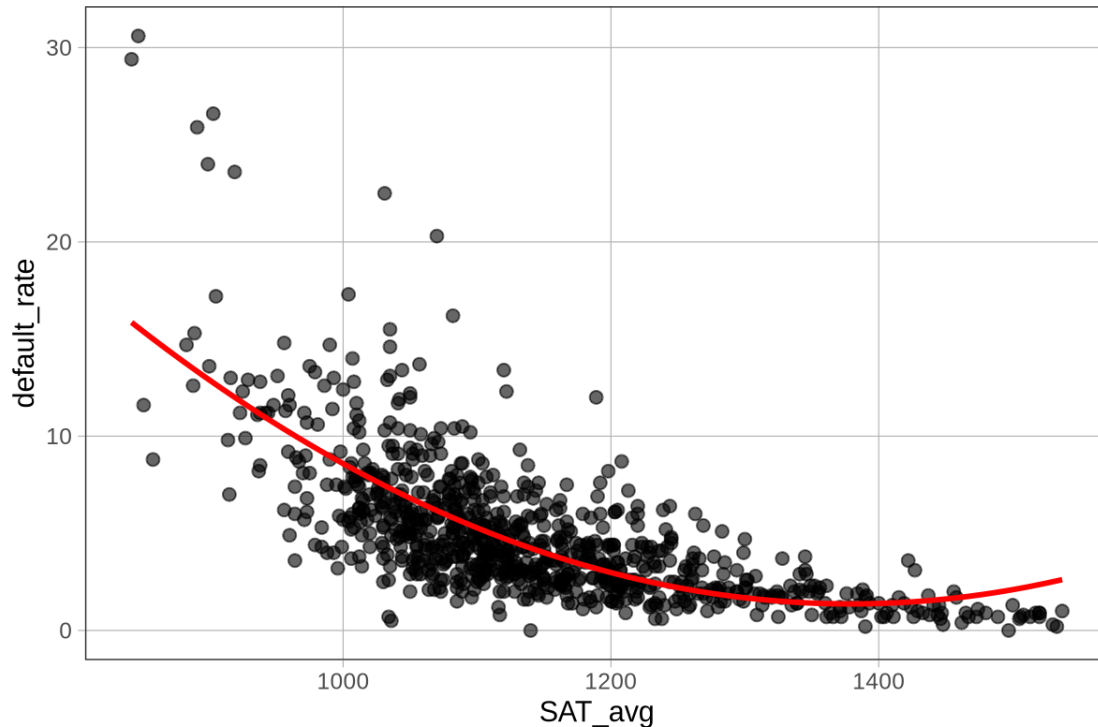
[1] "The test R squared value for the degree 1 model is 0.388138606283838"
[1] "The test R squared value for the degree 2 model is 0.508267621649189"
[1] "The test R squared value for the degree 3 model is 0.540190939875559"
[1] "The test R value for the degree 1 model is 0.623007709650401"
[1] "The test R value for the degree 2 model is 0.712928903642705"
[1] "The test R value for the degree 3 model is 0.7349768294821"

```

```

[269]: gf_point(default_rate ~ SAT_avg, data = sample_dat) %>% gf_lm(formula = y ~
  ↪poly(x,2), color="red")

```



```
[251]: model_1 <- lm(default_rate ~ poly(median_debt, 1), data = sample_dat)
model_2 <- lm(default_rate ~ poly(median_debt, 3), data = sample_dat)
model_3 <- lm(default_rate ~ poly(log(median_debt), 7), data = sample_dat)
pred_deg1 <- predict(model_1, newdata = data.frame(median_debt =
  ↳ test_dat$median_debt))
pred_deg2 <- predict(model_2, newdata = data.frame(median_debt =
  ↳ test_dat$median_debt))
pred_deg3 <- predict(model_3, newdata = data.frame(median_debt =
  ↳ test_dat$median_debt))
print(paste("The test R squared value for the degree 1 model is",
  ↳ cor(test_dat$default_rate, pred_deg1) ^ 2))
print(paste("The test R squared value for the degree 2 model is",
  ↳ cor(test_dat$default_rate, pred_deg2) ^ 2))
print(paste("The test R squared value for the degree 3 model is",
  ↳ cor(test_dat$default_rate, pred_deg3) ^ 2))
print(paste("The test R value for the degree 1 model is",
  ↳ cor(test_dat$default_rate, pred_deg1)))
print(paste("The test R value for the degree 2 model is",
  ↳ cor(test_dat$default_rate, pred_deg2)))
print(paste("The test R value for the degree 3 model is",
  ↳ cor(test_dat$default_rate, pred_deg3)))
```

```
[1] "The test R squared value for the degree 1 model is 0.0483338850153296"
```



```
[1] "The test R squared value for the degree 2 model is 0.0510691108391861"
[1] "The test R squared value for the degree 3 model is 0.0601127538246698"
[1] "The test R value for the degree 1 model is 0.219849687321428"
[1] "The test R value for the degree 2 model is 0.225984757979794"
[1] "The test R value for the degree 3 model is 0.245179024030747"
```

Check yourself: The R^2 for the degree 5 model should be about 0.6165

3.3 - Compare the R^2 estimates for each model. Which models did well? Which models did poorly? Why do you think this is?

Double-click this cell to type your answer here: The initial four models did well, with R^2 greater than 0.5, but the model for the 12 degree polynomial did poorly, as it had an R^2 value of 0.176, which is extremely low compared to the others. This could be because high degree polynomial models, which have more variables, have a higher likelihood of overfitting.

3.4 - In machine learning, the central goal is to build our models so as to avoid “underfitting” and “overfitting” our models to the training data. What do you think these terms mean? Which of our models were underfit? Which do you think were overfit? Explain.

Double-click this cell to type your answer here: When a model is underfitting, it is too simple and cannot capture the underlying trend in the data. When a model is overfitting, it is too complex and captures every little detail about the training data. The 12 degree polynomial model could be underfitting, as it has a low R^2 value. If the value for R^2 for the other four polynomial models were estimated using the training set, then compared to the results from using the testing set, it could be determined whether the models are overfitting or underfitting depending on whether the training R^2 is significantly higher than the testing R^2 or not.

Recall that we built our polynomial models here with just one predictor: x (`SAT_avg`). Yet, those models could end up being quite complex...

$$\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$$

Now, imagine that we wanted to bring in multiple predictors ($x_1 = \text{SAT_avg}$, $x_2 = \text{net_tuition}$, $x_3 = \text{grad_rate}$) for multiple regression. Plus, imagine that we decided to add in some polynomial terms for each of these predictors. We could end up with a model that looks ever more complicated, with literally hundreds of terms...

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_1^3 + \beta_4 x_2 + \beta_5 x_2^2 + \beta_6 x_2^3 + \beta_7 x_3 + \beta_8 x_3^2 + \dots$$

3.5 - Is it always good to add more predictors and add more polynomial terms to your model? Explain why or why not.

Double-click this cell to type your answer here: It isn't good to add more predictors and more polynomial terms to the model because it could result in overfitting due to the model having not learned the general relationship between the variables. This can also lead to the model having a high R^2 value on the training data.

0.1.6 4.0 - In-class prediction competition

Now you have all the tools you need to build very powerful prediction models! This means that it's time for a friendly competition :)

The code below takes the full dataset and splits it into larger train and test datasets. 80% of the colleges will go into the train dataset. 20% will go into the test dataset. Because we all are setting the same seed (2024), everyone will get the exact same train and test sets:

```
[5]: ## Run but do not edit this code

# set training data to be 80% of all colleges
train_size <- floor(0.8 * nrow(dat))

## sample row indeces
set.seed(2024)
train_ind <- sample(seq_len(nrow(dat)), size = train_size)

train <- dat[train_ind, ]
test <- dat[-train_ind, ]
```

```
[6]: dim(train)
```

```
1. 842 2. 26
```

```
[7]: dim(test)
```

```
1. 211 2. 26
```

Now it's time to compete!

Goal: Create the most accurate prediction model of colleges' default rates.

Evaluation: Whichever student has the highest R^2 on the test set wins.

Guidelines Save your best model as an object called `my_model`. You are only allowed to fit models on the train set (not on the test set). You may use as many predictors and as many polynomial terms as you'd like. Just be warned: Don't fall into the trap of overfitting! Choose only the most important variables and keep your models simple, so that you can generalize well to the test set. Periodically test your model on the test set and then make adjustments as necessary.

Go!

```
[14]: ## Your code goes here
# Replace 'grad_rate + poly(SAT_avg,3)' with your own combo of variables and
# poly terms
my_model <- lm(default_rate ~ poly(grad_rate,3) +
               poly(SAT_avg,2)+
               poly(net_tuition,2)+
               poly(pctPELL,4)+
               poly(pct_fed_loan,3)+
               poly(median_debt,1)+
```

```

poly(log(net_price),4)+
poly(log(avg_cost),3)+
poly(med_fam_income,4)+
poly(med_alum_earnings,4)+
poly(pct_firstgen,3)+
poly(log(enrollment),2)+highest_degree+hbcu+enrollment, data =
  ↪train)
train_pred = predict(my_model, newdata = train)
print(paste("The train R^2 value is:",cor(train$default_rate, train_pred)^2))
test_predictions = predict(my_model, newdata = test)
print(paste("The test R^2 value was: ", cor(test$default_rate,
  ↪test_predictions) ^ 2))

```

```

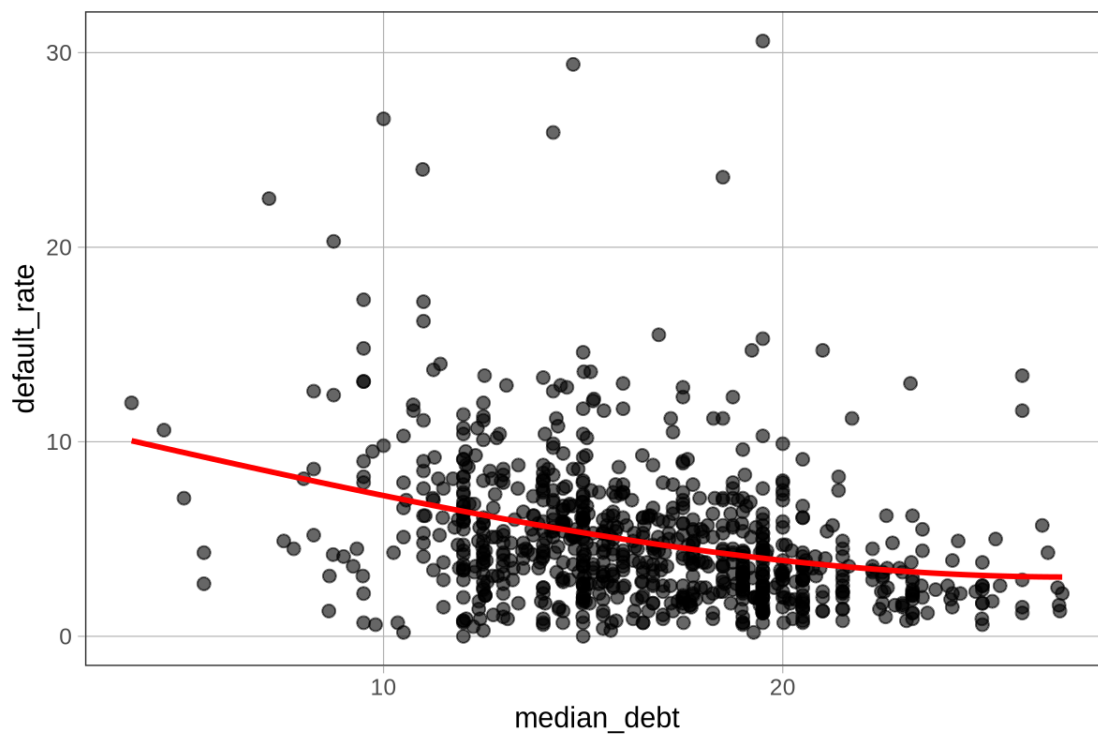
[1] "The train R^2 value is: 0.746041602218253"
[1] "The test R^2 value was: 0.758909619767055"

```

```

[146]: gf_point(default_rate ~ median_debt, data = sample_dat) %>% gf_lm(formula = y ~
  ↪poly(x,3), color="red")

```



```

[8]: str(dat)

```

```

'data.frame':  1053 obs. of  26 variables:
 $ OPEID      : int  100200 105200 105500 100500 105100 831000
100900 101200 100300 101900 ...

```

```

$ name          : chr "Alabama A & M University" "University of
Alabama at Birmingham" "University of Alabama in Huntsville" "Alabama State
University" ...
$ city          : chr "Normal" "Birmingham" "Huntsville" "Montgomery"
...
$ state         : chr "AL" "AL" "AL" "AL" ...
$ region        : chr "South" "South" "South" "South" ...
$ median_debt   : num 15.2 15.1 14 17.5 17.7 ...
$ default_rate  : num 12.1 4.8 4.7 12.8 4 8.2 2.6 4.4 9.9 10 ...
$ highest_degree : chr "Graduate" "Graduate" "Graduate" "Graduate" ...
$ ownership     : chr "Public" "Public" "Public" "Public" ...
$ locale        : chr "Small City" "Small City" "Small City" "Small
City" ...
$ hbcu          : chr "Yes" "No" "No" "Yes" ...
$ admit_rate    : num 89.7 80.6 77.1 98.9 80.4 ...
$ SAT_avg       : int 959 1245 1300 938 1262 1061 1302 1202 1068 1101
...
$ online_only   : chr "No" "No" "No" "No" ...
$ enrollment    : int 5090 13549 7825 3603 30610 4301 24368 1129 1834
917 ...
$ net_price     : num 15.5 16.5 17.2 19.5 20.9 ...
$ avg_cost      : num 23.4 25.5 24.9 21.9 30 ...
$ net_tuition   : num 8.1 11.99 8.28 9.3 14.71 ...
$ ed_spending_per_student : num 4.84 14.69 8.32 9.58 9.65 ...
$ avg_faculty_salary : num 7.6 11.38 9.7 7.19 10.35 ...
$ pctPELL       : num 71 34 24 73.7 17.2 ...
$ pct_fed_loan  : num 75 46.9 38.5 78 36.4 ...
$ grad_rate     : num 28.7 61.2 57.1 31.8 72.1 ...
$ pct_firstgen  : num 36.6 34.1 31 34.3 22.6 ...
$ med_fam_income : num 23.6 34.5 44.8 22.1 66.7 ...
$ med_alum_earnings : num 36.3 47 54.4 32.1 52.8 ...

```

```

[212]: for (i in 1:12) {
  model_x <- lm(default_rate ~ poly(log(net_tuition), i), data = sample_dat)
  pred_deg1 <- predict(model_x, newdata = data.frame(net_tuition =
↪test_dat$net_tuition))
  print(paste("The test R squared value for the model is :", i, " : ",
↪cor(test_dat$default_rate, pred_deg1) ^ 2))
  #print(paste("The test R value for the model is          :", i, " : ",
↪cor(test_dat$default_rate, pred_deg1)))
}

```

```

[1] "The test R squared value for the model is : 1 : 0.0606423544150001"
[1] "The test R squared value for the model is : 2 : 0.0989057217872761"
[1] "The test R squared value for the model is : 3 : 0.0994854151076478"
[1] "The test R squared value for the model is : 4 : 0.102223918877205"
[1] "The test R squared value for the model is : 5 : 0.0954917183096561"

```

```
[1] "The test R squared value for the model is : 6 : 0.0933385362107325"
[1] "The test R squared value for the model is : 7 : 0.0928327107735141"
[1] "The test R squared value for the model is : 8 : 0.0931167117217493"
[1] "The test R squared value for the model is : 9 : 0.0908184025253966"
[1] "The test R squared value for the model is : 10 : 0.092370360186314"
[1] "The test R squared value for the model is : 11 : 0.0827093001037942"
[1] "The test R squared value for the model is : 12 : 0.0801164325306628"
```

```
[9]: train_pred = predict(my_model, newdata = train)
      cor(train$default_rate, train_pred)^2
```

```
0.746041602218253
```

```
[8]: # run this code to get the R^2 value on the test set from your model
      test_predictions = predict(my_model, newdata = test)
      print(paste("The test R^2 value was: ", cor(test$default_rate,
      ↪test_predictions) ^ 2))
```

```
[1] "The test R^2 value was: 0.758909619767055"
```

0.1.7 5.0 - NATIONWIDE prediction competition

Competition: We’re hosting a *nationwide* competition to see which student can build the best model for predicting student loan default rates at different colleges. Here’s an article about last year’s winners.

Evaluation: Across the country, all students are using the same train and test sets as you did in the prior exercise to fit and evaluate their models. Your goal: Build a model that gives the best predictions on this test set. The student models that produce the highest R^2 value on the test set will be announced as champions!

Submission Process (due June 7, 2024 at 11:59pm CT): 1. Print and have a parent/guardian sign the media release form. This form gives permission to feature you and publish your results, in the event that you’re a finalist! Take a picture or scan the signed form and submit it during as a part of Step #2 (below). 2. Submit this google form (note: you’ll have to log into a google account), which allows you to upload your media release form, model, and notebook. This counts as your final submission.

Notes to avoid disqualification: - Do not change the seed (2024) in the code block that splits the data into the train and test sets. Using the common seed of 2024 will ensure everyone across the country has the exact same train/test split. - Make sure your model is fit using the `train` data. In other words, it should look like: `my_model <- lm(default_rate ~ ..., data = train)`. - Make sure you find the R^2 value on the `test` data, using the provided code. - There are ways to “cheat” on this competition by looking directly at the test set data values and designing your model to predict those values exactly (or approximately). However, based on the design of your model (which we’ll see when you share your notebook), it’s pretty easy for us to tell if you’ve done this. So, don’t do it! Your submission will be discarded.

0.1.8 Feedback (Required)

Please take 2 minutes to fill out this anonymous notebook feedback form, so we can continue improving this notebook for future years!