

432 Class 04 Slides

thomaseLove.github.io/432

2022-01-20

Today's Agenda

- Data Load from .Rds built in Class 3
- Fitting models with `lm`
 - Incorporating an interaction between factors
 - Incorporating polynomial terms
 - Incorporating restricted cubic splines
- Evaluating results in a testing sample with `yardstick`

Setup

```
knitr::opts_chunk$set(comment = NA)
options(width = 60)

library(here); library(magrittr)
library(janitor); library(knitr)
library(patchwork); library(broom)
library(rsample); library(yardstick)

library(rms)                ## new today: from Frank Harrell

library(tidyverse)

theme_set(theme_bw())
options(dplyr.summarise.inform = FALSE)
```

From Class 3

We developed the week2 data and performed a simple imputation for it (into week2im) in Class 3. Here, we'll read in those saved results, and then split into testing and training samples, as we did in Class 3.

```
week2 <- readRDS(here("data", "week2.Rds"))

week2im <- readRDS(here("data", "week2im.Rds"))

set.seed(432)
week2im_split <- initial_split(week2im, prop = 3/4)

train_w2im <- training(week2im_split)
test_w2im <- testing(week2im_split)
```

Codebook for useful week2 variables

- 894 subjects in Cleveland-Elyria with `bmi` and no history of diabetes

Variable	Description
<code>bmi</code>	(outcome) Body-Mass index in kg/m^2 .
<code>inc_imp</code>	income (imputed from grouped values) in \$
<code>fruit_day</code>	average fruit servings consumed per day
<code>drinks_wk</code>	average alcoholic drinks consumed per week
<code>female</code>	sex: 1 = female, 0 = male
<code>exerany</code>	any exercise in the past month: 1 = yes, 0 = no
<code>genhealth</code>	self-reported overall health (5 levels)
<code>race_eth</code>	race and Hispanic/Latinx ethnicity (5 levels)

- plus `ID`, `SEQNO`, `hx_diabetes` (all 0), `MMSA` (all Cleveland-Elyria)
- See Chapter 2 of the Course Notes for details on the variables

Class 03 and Class 04

In Class 03, we fit two models to predict bmi, using exerany and health, one with an interaction term and one without.

```
m_1 <- lm(bmi ~ exerany + health, data = train_w2im)
m_1int <- lm(bmi ~ exerany * health, data = train_w2im)
```

Adding in the covariate fruit_day...

```
m_2 <- lm(bmi ~ fruit_day + exerany + health,
          data = train_w2im)
m_2int <- lm(bmi ~ fruit_day + exerany * health,
            data = train_w2im)
```

Compare fits in the training sample?

mod	r.sq	adj.r.sq	sigma	df	df.res	AIC	BIC
m_1	0.089	0.082	6.12	5	664	4335.9	4367.5
m_2	0.098	0.090	6.09	6	663	4331.2	4367.3
m_1int	0.126	0.114	6.01	9	660	4315.8	4365.4
m_2int	0.138	0.125	5.97	10	659	4309.1	4363.2

- m_1 = no fruit_day
- m_2 = fruit_day included
- int = exerany*health interaction included

But the training sample cannot judge between models accurately. Our models have already *seen* that data. So fairer comparisons, consider the (held out) testing sample...

Model predictions of bmi in the test sample

We'll use `augment` from the `broom` package...

```
m1_test_aug <- augment(m_1, newdata = test_w2im)
m1int_test_aug <- augment(m_1int, newdata = test_w2im)
m2_test_aug <- augment(m_2, newdata = test_w2im)
m2int_test_aug <- augment(m_2int, newdata = test_w2im)
```

This adds fitted values (predictions) and residuals (errors) ...

```
m1_test_aug %>% select(ID, bmi, .fitted, .resid) %>%
  slice(1:2) %>% kable()
```

ID	bmi	.fitted	.resid
4	26.51	28.85006	-2.340059
5	24.25	28.85006	-4.600059

What does the `yardstick` package do?

For each subject in the testing set, we will need:

- `estimate` = model's prediction of that subject's `bmi`
- `truth` = the `bmi` value observed for that subject

Calculate a summary of the predictions across the n test subjects, such as:

- R^2 = square of the correlation between `truth` and `estimate`
- `mae` = mean absolute error ...

$$mae = \frac{1}{n} \sum |truth - estimate|$$

- `rmse` = root mean squared error ...

$$rmse = \sqrt{\frac{1}{n} \sum (truth - estimate)^2}$$

Testing Results (using R^2)

We can use the `yardstick` package and its `rsq()` function.

```
testing_r2 <- bind_rows(  
  rsq(m1_test_aug, truth = bmi, estimate = .fitted),  
  rsq(m1int_test_aug, truth = bmi, estimate = .fitted),  
  rsq(m2_test_aug, truth = bmi, estimate = .fitted),  
  rsq(m2int_test_aug, truth = bmi, estimate = .fitted)) %>%  
  mutate(model = c("m_1", "m_1int", "m_2", "m_2int"))  
testing_r2 %>% kable(dig = 4)
```

.metric	.estimator	.estimate	model
rsq	standard	0.0712	m_1
rsq	standard	0.0395	m_1int
rsq	standard	0.0647	m_2
rsq	standard	0.0361	m_2int

Mean Absolute Error?

Consider the mean absolute prediction error ...

```
testing_mae <- bind_rows(  
  mae(m1_test_aug, truth = bmi, estimate = .fitted),  
  mae(m1int_test_aug, truth = bmi, estimate = .fitted),  
  mae(m2_test_aug, truth = bmi, estimate = .fitted),  
  mae(m2int_test_aug, truth = bmi, estimate = .fitted)) %>%  
  mutate(model = c("m_1", "m_1int", "m_2", "m_2int"))  
testing_mae %>% kable(dig = 2)
```

.metric	.estimator	.estimate	model
mae	standard	4.43	m_1
mae	standard	4.63	m_1int
mae	standard	4.48	m_2
mae	standard	4.71	m_2int

Root Mean Squared Error?

How about the square root of the mean squared prediction error, or RMSE?

```
testing_rmse <- bind_rows(  
  rmse(m1_test_aug, truth = bmi, estimate = .fitted),  
  rmse(m1int_test_aug, truth = bmi, estimate = .fitted),  
  rmse(m2_test_aug, truth = bmi, estimate = .fitted),  
  rmse(m2int_test_aug, truth = bmi, estimate = .fitted)) %>%  
  mutate(model = c("m_1", "m_1int", "m_2", "m_2int"))  
testing_rmse %>% kable(digits = 3)
```

.metric	.estimator	.estimate	model
rmse	standard	5.736	m_1
rmse	standard	6.033	m_1int
rmse	standard	5.778	m_2
rmse	standard	6.091	m_2int

Other Summaries for Numerical Predictions

Within the `yardstick` package, there are several other summaries, including:

- `rsq_trad()` = defines R^2 using sums of squares.
 - The `rsq()` measure we showed a few slides ago is a squared correlation coefficient and is guaranteed to fall in $(0, 1)$.
- `mape()` = mean absolute percentage error
- `mpe()` = mean percentage error
- `huber_loss()` = Huber loss (often used in robust regression), which is less sensitive to outliers than `rmse()`.
- `ccc()` = concordance correlation coefficient, which attempts to measure both consistency/correlation (like `rsq()`) and accuracy (like `rmse()`).

See the `yardstick` home page for more details.

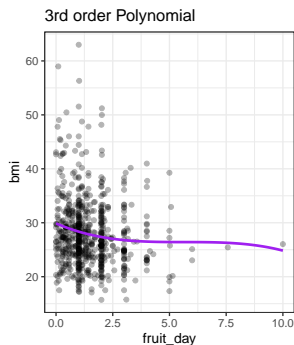
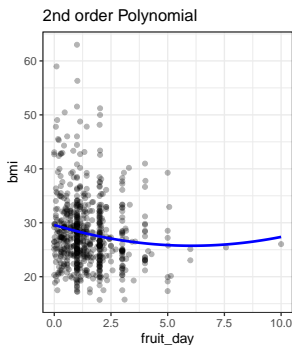
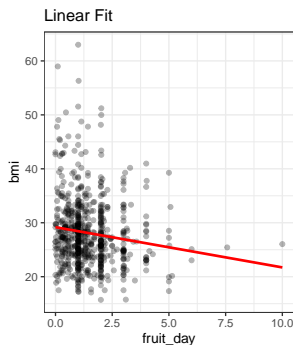
Incorporating a non-linear term for fruit_day

Suppose we wanted to include a polynomial term for fruit_day:

```
lm(bmi ~ fruit_day, data = train_w2im)
```

```
lm(bmi ~ poly(fruit_day, 2), data = train_w2im)
```

```
lm(bmi ~ poly(fruit_day, 3), data = train_w2im)
```



Polynomial Regression

A polynomial in the variable x of degree D is a linear combination of the powers of x up to D .

For example:

- Linear: $y = \beta_0 + \beta_1 x$
- Quadratic: $y = \beta_0 + \beta_1 x + \beta_2 x^2$
- Cubic: $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$
- Quartic: $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 x^4$
- Quintic: $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 x^4 + \beta_5 x^5$

Fitting such a model creates a **polynomial regression**.

Raw Polynomials vs. Orthogonal Polynomials

Predict bmi using fruit_day with a polynomial of degree 2.

```
(temp1 <- lm(bmi ~ fruit_day + I(fruit_day^2),  
             data = train_w2im))
```

Call:

```
lm(formula = bmi ~ fruit_day + I(fruit_day^2), data = train_w2im)
```

Coefficients:

(Intercept)	fruit_day	I(fruit_day^2)
29.5925	-1.2733	0.1051

This uses raw polynomials. Predicted bmi for fruit_day = 2 is

$$\begin{aligned}\text{bmi} &= 29.5925 - 1.2733 (\text{fruit_day}) + 0.1051 (\text{fruit_day}^2) \\ &= 29.5925 - 1.2733 (2) + 0.1051 (4) \\ &= 27.466\end{aligned}$$

Does the raw polynomial match our expectations?

```
temp1 <- lm(bmi ~ fruit_day + I(fruit_day^2),  
            data = train_w2im)  
  
augment(temp1, newdata = tibble(fruit_day = 2)) %>%  
  kable(digits = 3)
```

fruit_day	.fitted
2	27.466

and this matches our “by hand” calculation. But it turns out most regression models use *orthogonal* rather than raw polynomials. . .

Fitting an Orthogonal Polynomial

Predict bmi using fruit_day with an *orthogonal* polynomial of degree 2.

```
(temp2 <- lm(bmi ~ poly(fruit_day,2), data = train_w2im))
```

Call:

```
lm(formula = bmi ~ poly(fruit_day, 2), data = train_w2im)
```

Coefficients:

```
(Intercept)  poly(fruit_day, 2)1  
28.089          -21.636
```

```
poly(fruit_day, 2)2  
8.009
```

This looks very different from our previous version of the model.

- What happens when we make a prediction, though?

Prediction in the Orthogonal Polynomial Model

Remember that in our raw polynomial model, our “by hand” and “using R” calculations both concluded that the predicted `bmi` for a subject with `fruit_day = 2` was 27.466.

Now, what happens with the orthogonal polynomial model `temp2` we just fit?

```
augment(temp2, newdata = data.frame(fruit_day = 2)) %>%  
  kable(digits = 3)
```

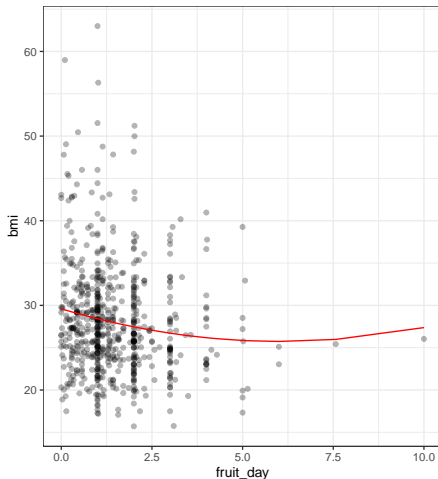
fruit_day		.fitted
2		27.466

- No change in the prediction.

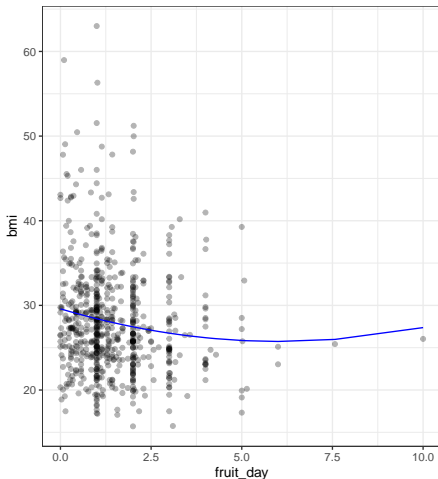
Fits of raw vs orthogonal polynomials

Comparing Two Methods of Fitting a Quadratic Polynomial

temp1: Raw fit, degree 2



temp2: Orthogonal fit, degree 2



- The two models are, in fact, identical.

Why do we use orthogonal polynomials?

- The main reason is to avoid having to include powers of our predictor that are highly collinear.
- Variance Inflation Factor assesses collinearity...

```
vif(temp1)          ## from rms package
```

```
fruit_day I(fruit_day^2)  
4.652178      4.652178
```

- Orthogonal polynomial terms are uncorrelated with one another, easing the process of identifying which terms add value to our model.

```
vif(temp2)
```

```
poly(fruit_day, 2)1 poly(fruit_day, 2)2  
1                      1
```

Why orthogonal rather than raw polynomials?

The tradeoff is that the raw polynomial is a lot easier to explain in terms of a single equation in the simplest case.

Actually, we'll usually avoid polynomials in our practical work, and instead use splines, which are more flexible and require less maintenance, but at the cost of pretty much requiring you to focus on visualizing their predictions rather than their equations.

Adding a Second Order Polynomial to our Models

```
m_3 <- lm(bmi ~ poly(fruit_day,2) + exerany + health,  
          data = train_w2im)
```

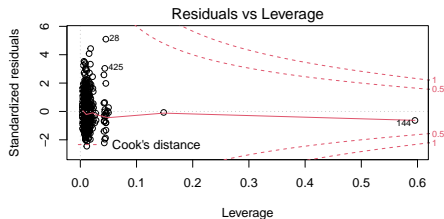
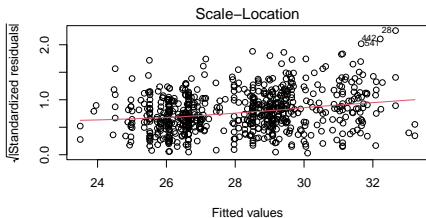
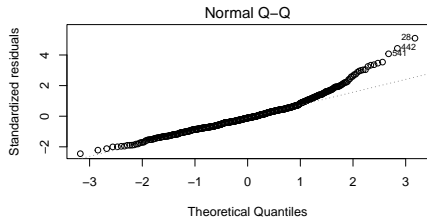
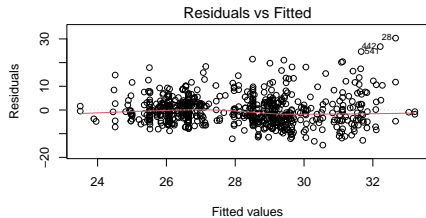
- Comparison to other models without the interaction...

mod	r.sq	adj.r.sq	sigma	df	df.res	AIC	BIC
m_1	0.0889	0.0821	6.12	5	664	4335.9	4367.5
m_2	0.0980	0.0898	6.09	6	663	4331.2	4367.3
m_3	0.0980	0.0885	6.09	7	662	4333.2	4373.8

Tidied summary of `m_3` coefficients

term	est	se	t	p	conf.low	conf.high
(Intercept)	27.87	0.74	37.54	0.000	26.65	29.10
poly(fruit_day, 2)1	-15.88	6.16	-2.58	0.010	-26.04	-5.73
poly(fruit_day, 2)2	1.08	6.24	0.17	0.863	-9.20	11.35
exerany	-2.03	0.56	-3.63	0.000	-2.95	-1.11
healthVG	0.55	0.70	0.79	0.430	-0.60	1.71
healthG	3.00	0.72	4.15	0.000	1.81	4.19
healthF	3.55	0.91	3.92	0.000	2.06	5.04
healthP	4.54	1.36	3.33	0.001	2.29	6.78

m_3 Residual Plots



Add in the interaction

```
m_3int <- lm(bmi ~ poly(fruit_day,2) + exerany * health,  
             data = train_w2im)
```

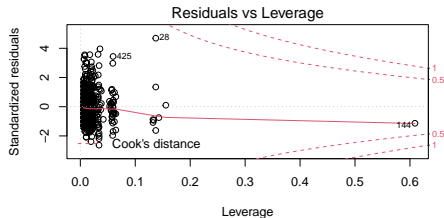
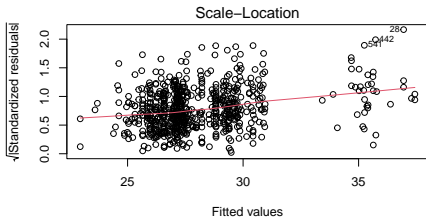
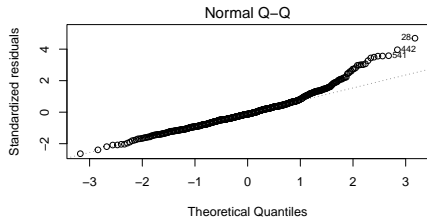
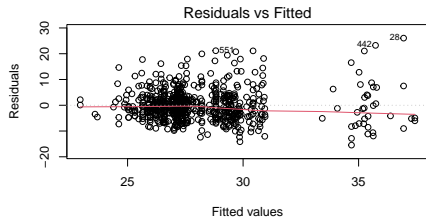
- Comparison to other models with the interaction...

mod	r.sq	adj.r.sq	sigma	df	df.res	AIC	BIC
m_1int	0.1264	0.1145	6.01	9	660	4315.8	4365.4
m_2int	0.1376	0.1245	5.97	10	659	4309.1	4363.2
m_3int	0.1377	0.1233	5.98	11	658	4311.1	4369.6

Tidied summary of `m_3int` coefficients

term	est	se	t	p	conf.low	conf.high
(Intercept)	27.40	1.41	19.43	0.000	25.07	29.72
poly(fruit_day, 2)1	-17.73	6.06	-2.93	0.004	-27.71	-7.75
poly(fruit_day, 2)2	-1.66	6.29	-0.26	0.792	-12.02	8.70
exerany	-1.46	1.55	-0.94	0.346	-4.00	1.09
healthVG	-0.66	1.63	-0.40	0.686	-3.34	2.02
healthG	2.75	1.61	1.71	0.089	0.09	5.41
healthF	7.58	1.77	4.28	0.000	4.66	10.50
healthP	9.27	2.61	3.55	0.000	4.97	13.56
exerany:healthVG	1.59	1.79	0.88	0.377	-1.37	4.54
exerany:healthG	0.42	1.80	0.23	0.814	-2.54	3.38
exerany:healthF	-6.40	2.06	-3.10	0.002	-9.80	-3.00
exerany:healthP	-6.71	3.05	-2.20	0.028	-11.74	-1.68

m_3int Residual Plots



How do models `m_3` and `m_3int` do in testing?

```
m3_test_aug <- augment(m_3, newdata = test_w2im)
m3int_test_aug <- augment(m_3int, newdata = test_w2im)

testing_r2 <- bind_rows(
  rsq(m1_test_aug, truth = bmi, estimate = .fitted),
  rsq(m2_test_aug, truth = bmi, estimate = .fitted),
  rsq(m3_test_aug, truth = bmi, estimate = .fitted),
  rsq(m1int_test_aug, truth = bmi, estimate = .fitted),
  rsq(m2int_test_aug, truth = bmi, estimate = .fitted),
  rsq(m3int_test_aug, truth = bmi, estimate = .fitted)) %>%
  mutate(model = c("m_1", "m_2", "m_3", "m_1int",
                    "m_2int", "m_3int"))
```

- I've hidden my calculations for RMSE and MAE here.

Results comparing all six models (testing)

```
bind_cols(testing_r2 %>% select(model, rsquare = .estimate),  
          testing_rmse %>% select(rmse = .estimate),  
          testing_mae %>% select(mae = .estimate)) %>%  
kable(digits = c(0, 4, 3, 3))
```

model	rsquare	rmse	mae
m_1	0.0712	5.736	4.432
m_2	0.0647	5.778	4.480
m_3	0.0651	5.776	4.480
m_1int	0.0395	6.033	4.628
m_2int	0.0361	6.091	4.710
m_3int	0.0354	6.098	4.711

- Did the polynomial term in m_3 and m_3int improve our predictions?

Splines

- A **linear spline** is a continuous function formed by connecting points (called **knots** of the spline) by line segments.
- A **restricted cubic spline** is a way to build highly complicated curves into a regression equation in a fairly easily structured way.
- A restricted cubic spline is a series of polynomial functions joined together at the knots.
 - Such a spline gives us a way to flexibly account for non-linearity without over-fitting the model.
 - Restricted cubic splines can fit many different types of non-linearities.
 - Specifying the number of knots is all you need to do in R to get a reasonable result from a restricted cubic spline.

The most common choices are 3, 4, or 5 knots.

- 3 Knots, 2 degrees of freedom, allows the curve to “bend” once.
- 4 Knots, 3 degrees of freedom, lets the curve “bend” twice.
- 5 Knots, 4 degrees of freedom, lets the curve “bend” three times.

A simulated data set

```
set.seed(4322021)

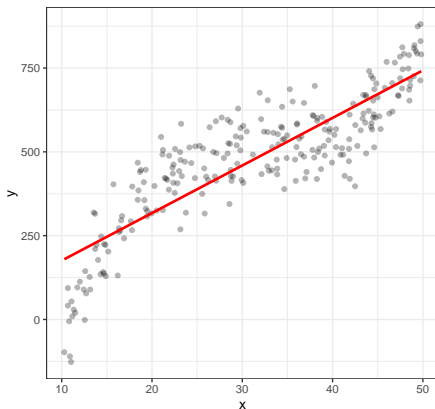
sim_data <- tibble(
  x = runif(250, min = 10, max = 50),
  y = 3*(x-30) - 0.3*(x-30)^2 + 0.05*(x-30)^3 +
    rnorm(250, mean = 500, sd = 70)
)

head(sim_data, 2)

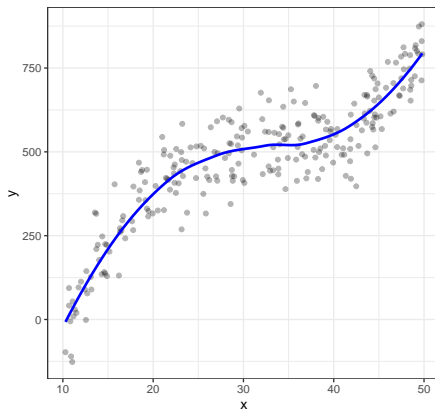
# A tibble: 2 x 2
      x      y
<dbl> <dbl>
1  42.5  397.
2  35.9  414.
```


The `sim_data`, plotted.

With Linear Fit



With Loess Smooth

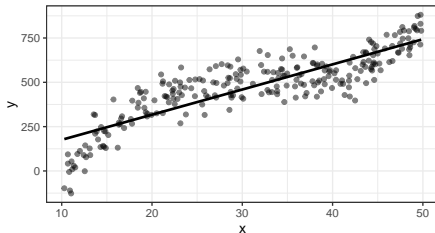


Fitting Restricted Cubic Splines with `lm` and `rcs`

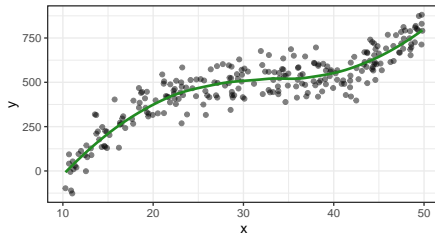
```
sim_linear <- lm(y ~ x, data = sim_data)
sim_poly2  <- lm(y ~ poly(x, 2), data = sim_data)
sim_poly3  <- lm(y ~ poly(x, 3), data = sim_data)
sim_rcs3   <- lm(y ~ rcs(x, 3), data = sim_data)
sim_rcs4   <- lm(y ~ rcs(x, 4), data = sim_data)
sim_rcs5   <- lm(y ~ rcs(x, 5), data = sim_data)
```

Looking at the Polynomial Fits

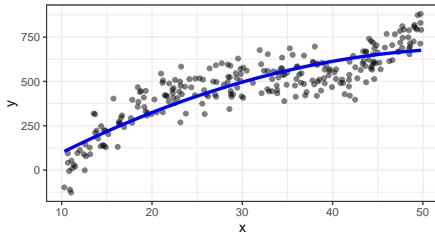
Linear Fit



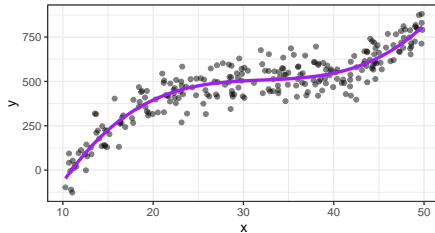
Loess Smooth



Quadratic Polynomial

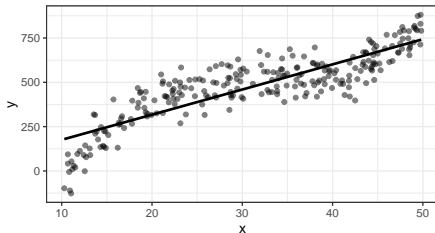


Cubic Polynomial

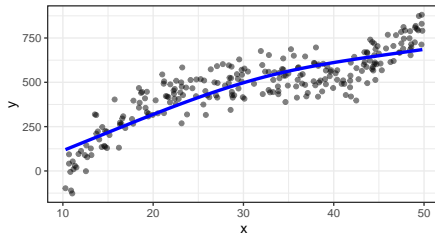


Looking at the Restricted Cubic Spline Fits

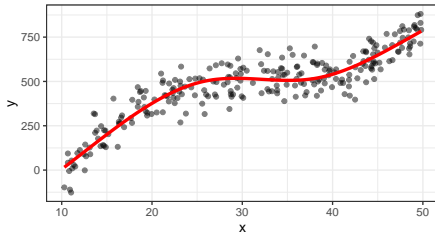
Linear Fit



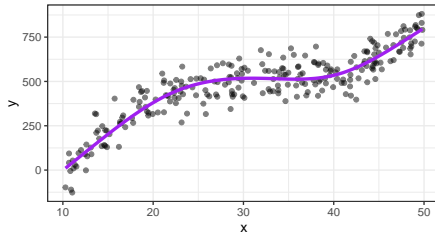
RCS with 3 knots



RCS with 4 knots



RCS with 5 knots



Fitting Restricted Cubic Splines with `lm` and `rcs`

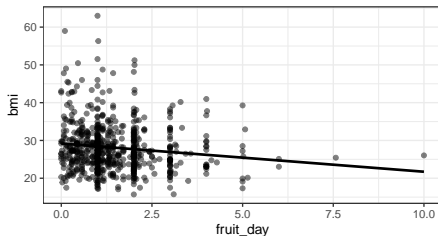
For most applications, three to five knots strike a nice balance between complicating the model needlessly and fitting data pleasingly. Let's consider a restricted cubic spline model for `bmi` based on `fruit_day` again, but now with:

- in `temp3`, 3 knots, and
- in `temp4`, 4 knots,

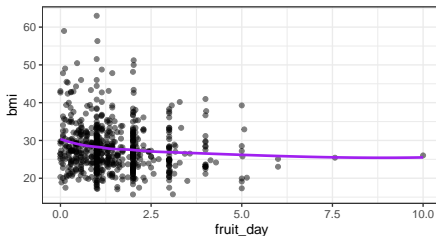
```
temp3 <- lm(bmi ~ rcs(fruit_day, 3), data = train_w2im)
temp4 <- lm(bmi ~ rcs(fruit_day, 4), data = train_w2im)
```

Spline models for bmi and fruit_day

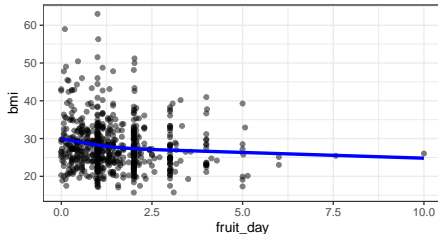
Linear Fit



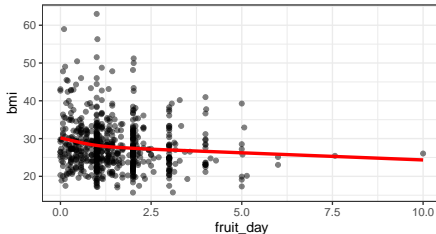
Loess Smooth



RCS, 3 knots



RCS, 4 knots



Let's try an RCS with 4 knots

```
m_4 <- lm(bmi ~ rcs(fruit_day, 4) + exerany + health,  
          data = train_w2im)
```

```
m_4int <- lm(bmi ~ rcs(fruit_day, 4) + exerany * health,  
            data = train_w2im)
```

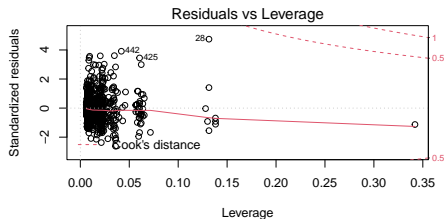
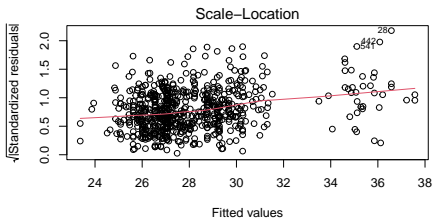
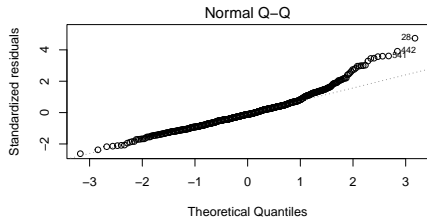
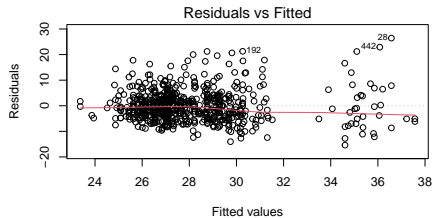
Comparing 4 models including the `exerany*health` interaction...

mod	fruit	r.sq	adj.r.sq	sigma	df	AIC	BIC
m_1int	not in	0.1264	0.1145	6.007	9	4315.8	4365.4
m_2int	linear	0.1376	0.1245	5.973	10	4309.1	4363.2
m_3int	poly(2)	0.1377	0.1233	5.977	11	4311.1	4369.6
m_4int	rcs(4)	0.1379	0.1222	5.981	12	4312.9	4376.0

Tidied summary of `m_4int` coefficients

term	est	se	t	p	lo90	hi90
(Intercept)	28.57	1.54	18.50	0.000	26.02	31.11
rcs(fruit_day, 4)fruit_day	-1.22	1.27	-0.96	0.337	-3.30	0.87
rcs(fruit_day, 4)fruit_day'	2.33	5.82	0.40	0.689	-7.26	11.92
rcs(fruit_day, 4)fruit_day''	-6.00	16.58	-0.36	0.717	-33.32	21.31
exerany	-1.34	1.55	-0.87	0.387	-3.90	1.21
healthVG	-0.64	1.63	-0.39	0.696	-3.32	2.04
healthG	2.78	1.61	1.72	0.086	0.12	5.44
healthF	7.64	1.77	4.30	0.000	4.71	10.56
healthP	9.09	2.56	3.54	0.000	4.86	13.31
exerany:healthVG	1.56	1.80	0.87	0.385	-1.40	4.52
exerany:healthG	0.36	1.80	0.20	0.841	-2.60	3.32
exerany:healthF	-6.47	2.07	-3.13	0.002	-9.87	-3.06
exerany:healthP	-6.55	3.02	-2.17	0.031	-11.53	-1.57

m_4int Residual Plots



How do models `m_4` and `m_4int` do in testing?

model	rsquare	rmse	mae
<code>m_1</code>	0.0712	5.736	4.432
<code>m_2</code>	0.0647	5.778	4.480
<code>m_3</code>	0.0651	5.776	4.480
<code>m_4</code>	0.0682	5.771	4.474
<code>m_1int</code>	0.0395	6.033	4.628
<code>m_2int</code>	0.0361	6.091	4.710
<code>m_3int</code>	0.0354	6.098	4.711
<code>m_4int</code>	0.0393	6.099	4.707

I'll note that there's a fair amount of very repetitive code in the R Markdown file to create that table.

- What are our conclusions?

Next Week

- What if we want to build models for a binary outcome, rather than a quantitative one?
- Lab 1 due Monday at 9 PM.
- Tuesday's class will be a Zoom session, Thursday in person.