

ISLR Chapter 4 Exercises

Ryan Heslin

11/17/2020

Lab 5: Cross-Validation and Bootstrapping

2. We will now derive the probability that a given observation is part of a bootstrap sample. Suppose that we obtain a bootstrap sample from a set of n observations
 - a. $\frac{n-1}{n}$. Let $j = 1$. There are $n - 1$ other values the bootstrap observation could take.
 - b. Again $\frac{n-1}{n}$, because each observation has an equal chance of appearing in each draw.
 - c. This is a case of a repeated event with unvarying probability. The expression I gave above can be rewritten as $1 - \frac{1}{n}$. This is simply raised to the n th power, since it is repeated once for each observation.

```
library(rlang)
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.0 --

## v ggplot2 3.3.2      v purrr 0.3.4
## v tibble 3.0.1       v dplyr 1.0.0
## v tidyr 1.1.0        v stringr 1.4.0
## v readr 1.4.0        v forcats 0.5.0

## -- Conflicts ----- tidyverse_conflicts() --
## x purrr::%@%()      masks rlang::%@%()
## x purrr::as_function() masks rlang::as_function()
## x dplyr::filter()   masks stats::filter()
## x purrr::flatten()  masks rlang::flatten()
## x purrr::flatten_chr() masks rlang::flatten_chr()
## x purrr::flatten_dbl() masks rlang::flatten_dbl()
## x purrr::flatten_int() masks rlang::flatten_int()
## x purrr::flatten_lgl() masks rlang::flatten_lgl()
## x purrr::flatten_raw() masks rlang::flatten_raw()
## x purrr::invoke()   masks rlang::invoke()
## x dplyr::lag()       masks stats::lag()
## x purrr::list_along() masks rlang::list_along()
## x purrr::modify()    masks rlang::modify()
## x purrr::prepend()   masks rlang::prepend()
## x purrr::splice()    masks rlang::splice()
```

```
library(boot)
p_in <- expr(1 -(1-1/n)^n)

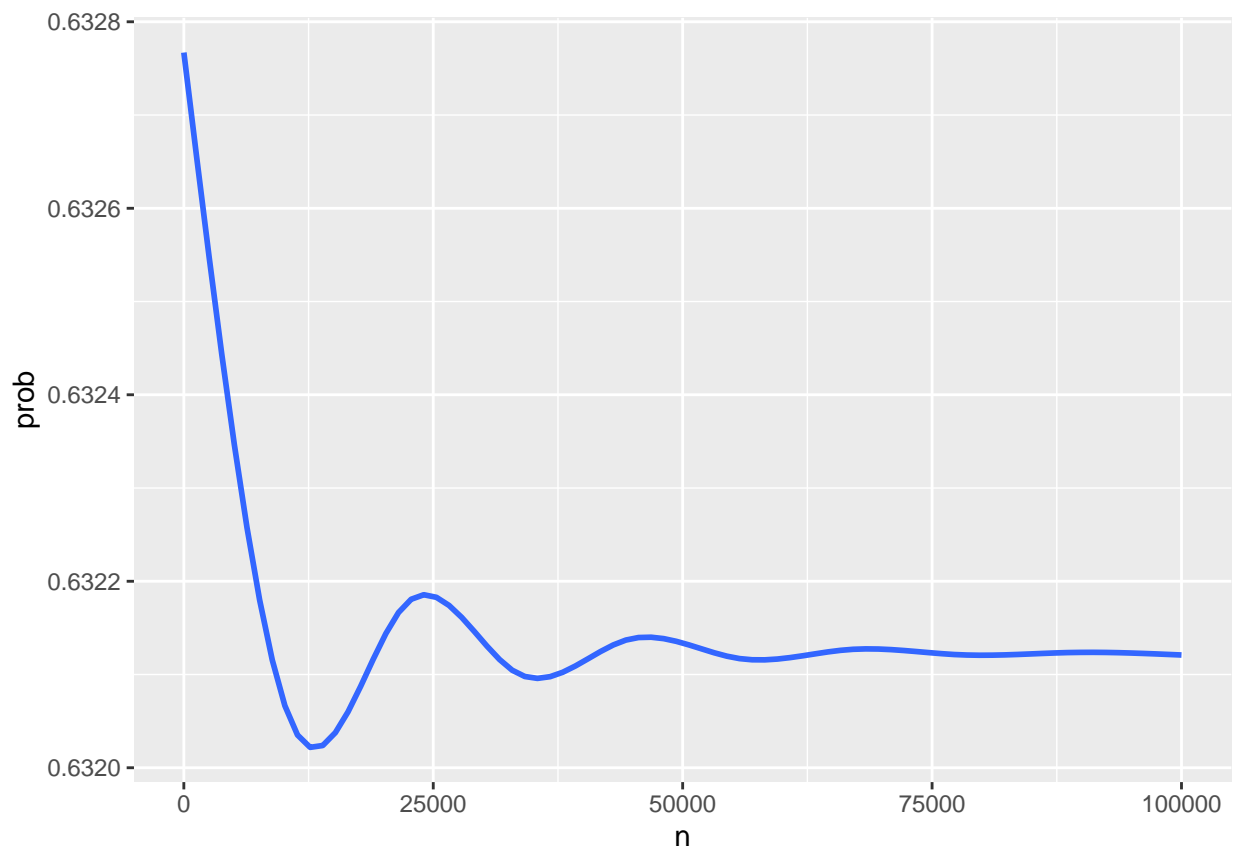
eval(p_in, env = list(n= c(5, 10, 10000)))
```

```
## [1] 0.6723200 0.6513216 0.6321390
```

```
dat <- 1:100000

tibble(n = 1:100000, prob = rlang::eval_tidy(p_in)) %>%
  ggplot(aes(x = n, y = prob)) +
  geom_smooth(method = "gam", se = F)
```

```
## 'geom_smooth()' using formula 'y ~ s(x, bs = "cs")'
```



Notice the amusing way the smoother overfits.

```
store=rep(NA, 10000)
for(i in 1:10000){
  store[i]=sum(sample(1:100, rep=TRUE)==4)>0
}
mean(store)
```

```
## [1] 0.6329
```

```
store <- replicate(10000, sum(sample(1:100, rep=TRUE)==4)>0)
```

The proportion is about 60%. From the expression given above, we know

$$P(j \text{ in bootstrap}) = (1 - \frac{1}{n})^n$$

If we subtract this from 1 to get the odds of j appearing in the bootstrap, we get roughly the observed mean.

3. We now review k-fold cross-validation. (a) Explain how k-fold cross-validation is implemented. (b) What are the advantages and disadvantages of k-fold cross-validation relative to: i. The validation set approach? ii. LOOCV?
 - a. The data are divided into K equally sized groups, usually about five or ten. One is used as a validation set, and the remainder for training. After each iteration, a training fold swaps places with the validation fold.

Compared to validation sets, k-folding preserves much more data for training and therefore has less bias. Compared to LOOCV, it is less variant because each fitted model has $n - \frac{n}{K}$ observations in common with the others (the one overlapping fold) in common with the others. In LOOCV, the models differ by just one row and are thus highly correlated.

4. Suppose that we use some statistical learning method to make a prediction for the response Y for a particular value of the predictor X . Carefully describe how we might estimate the standard deviation of our prediction.

We use an extension of the typical formula with a LOOCV sample. We sum the difference of each predicted value from the mean of the every other predicted value, then divide this by $n - 1$ to get the standard error.

5. In Chapter 4, we used logistic regression to predict the probability of default using income and balance on the Default data set. We will now estimate the test error of this logistic regression model using the validation set approach. Do not forget to set a random seed before beginning your analysis.

I fit a model and test it using a validation set for four iterations, the results are pretty similar, as we'd expect. Adding a student dummy variable makes little difference.

```
library(ISLR)
set.seed(1)
default <- Default
train_inds <- replicate(sample(seq(1, nrow(default)), 5000, replace = F), n = 4) %>% as.data.frame()

map(train_inds, function(train_inds){
  train <- default[train_inds,]
  valid <- default[-train_inds,]
  mod <- glm(data = train, default ~ income + balance, family = binomial)

  predict(mod, valid, type = "response") %>%
    round %>% factor(labels = c("No", "Yes")) %>%
    table(., valid$default)
}) %>% set_names(1:4)
```

```
## $'1'
##
```

```
## .      No  Yes
##  No  4824  108
##  Yes   19   49
##
## $'2'
##
## .      No  Yes
##  No  4813  112
##  Yes   25   50
##
## $'3'
##
## .      No  Yes
##  No  4827  103
##  Yes   19   51
##
## $'4'
##
## .      No  Yes
##  No  4824  106
##  Yes   16   54
```

```
default <- default %>% mutate(student_dummy = as.numeric(student)-1)
mod2 <- glm(data =default[train_inds[,1],], default ~ income+ balance + student_dummy, family= binomial)
predict(mod2, default[-train_inds[,1],], type = "response") %>%
  round %>% factor(labels = c("No", "Yes")) %>%
  table
```

```
## .
##  No  Yes
## 4937  63
```

```
(predict(mod2, default[-train_inds[,1],], type = "response") %>%
  round() %>% factor(labels = c("No", "Yes")) != default[-train_inds[,1],$default)
```

```
## [1] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [13] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [25] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [37] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [49] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [61] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE FALSE
## [73] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [85] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [97] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [109] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [121] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [133] FALSE FALSE FALSE FALSE FALSE TRUE FALSE FALSE FALSE FALSE FALSE FALSE
## [145] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [157] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [169] FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE FALSE FALSE FALSE TRUE
## [181] FALSE FALSE FALSE FALSE TRUE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [193] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [205] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
```

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

```
## [4753] FALSE FALSE FALSE FALSE FALSE TRUE FALSE FALSE FALSE FALSE FALSE FALSE
## [4765] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE FALSE FALSE
## [4777] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [4789] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [4801] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [4813] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [4825] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [4837] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [4849] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [4861] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [4873] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [4885] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE FALSE
## [4897] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE FALSE
## [4909] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [4921] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [4933] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [4945] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [4957] FALSE FALSE FALSE TRUE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [4969] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE FALSE FALSE
## [4981] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## [4993] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
```

This test error seems good, but it's not much better than a naive classifier that never predicted default which would have an error rate of:

```
mean(default$default=="Yes")
```

```
## [1] 0.0333
```

Now we are asked to compute SE's using both bootstraps and the glm function.

The text notes that bootstrap SE estimates are more accurate than the standard formulas because they do not assume the model is accurate when estimating the error variance and do not assume the x_i are fixed.

The bootstrap reports much lower standard errors for each term, interestingly. All the errors look normally distributed.

```
#map_dfr(list(mod1, mod2), broom::tidy)

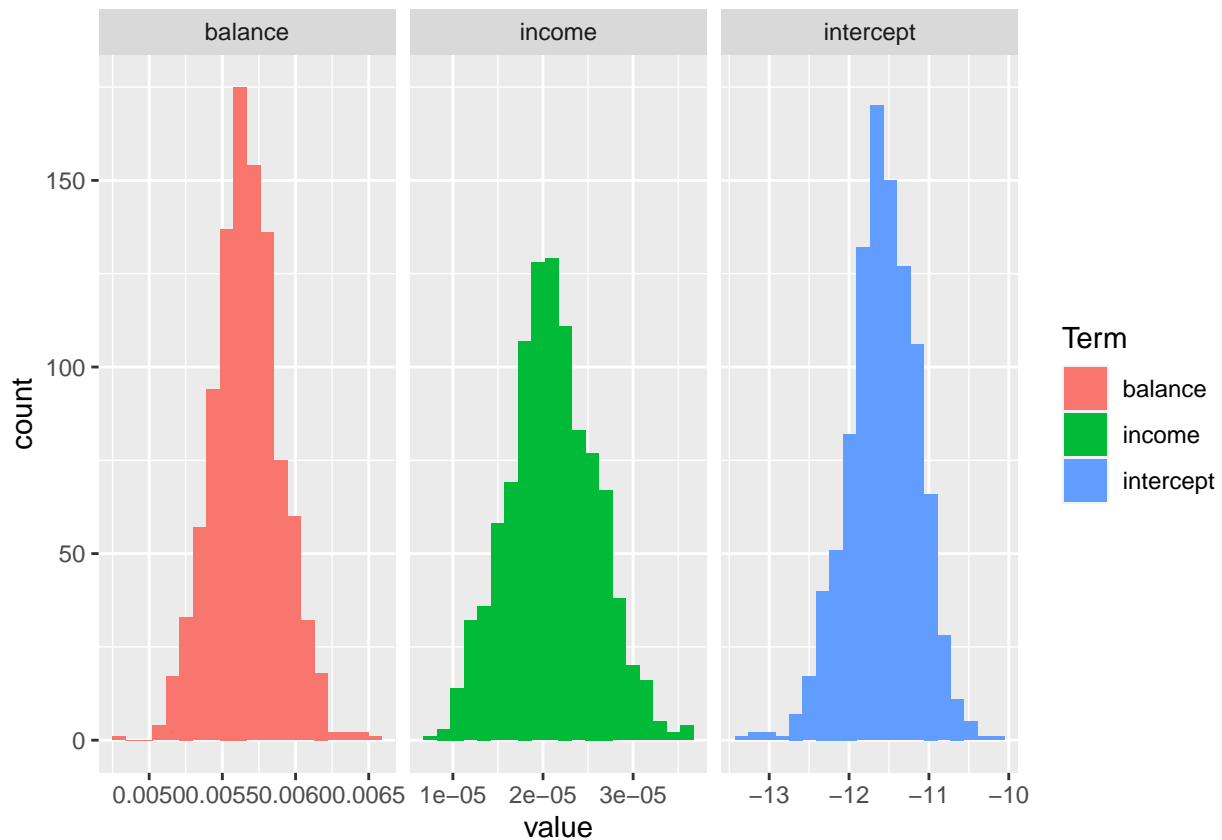
get_coefs <- function(df, inds){
  glm(data = df[inds,], default ~ income + balance, family = binomial) %>% coef()
}
get_coefs(default, sample(nrow(default), nrow(default), replace = TRUE))
```

```
## (Intercept)      income      balance
## -1.139432e+01  2.333588e-05  5.505038e-03
```

```
estimate <- boot(default, get_coefs, 1000)

estimate$t %>% as.data.frame() %>%
  set_names(c("intercept", "income", "balance")) %>%
  pivot_longer(everything(), names_to = "Term") %>%
```

```
ggplot(aes(x = value, fill = Term)) +
  geom_histogram(bins = 20) +
  facet_wrap(. ~ Term, scales = "free_x")
```



7. We are now asked to use LOOCV to validate a logistic model. I hold out the first observation and verify it was correctly classified

```
weekly <- Weekly
mod <- glm(data = weekly, Direction ~ Lag1 + Lag2, family = binomial)
mod_train <- glm(data = weekly[-1,], Direction ~ Lag1 + Lag2, family = binomial)

(predict(mod_train, weekly, mode = "response") %>% {. [1]} %>% round()) == (as.numeric(weekly$Direction) - 1)

##      1
## TRUE
```

Next we iterate LOOCV over the dataset. The model does not improve much over chance.

```
error <- vector(length = nrow(weekly))
for(i in seq(1, nrow(weekly))) {
  dat <- weekly[-i,]
  mod <- glm(data = dat, Direction ~ Lag1 + Lag2, family = binomial)
  error[i] <- (predict(mod, weekly, mode = "response") %>% {. [i]} %>% round()) == (as.numeric(weekly$Direction) - 1)[i]
}
```

```
}
mean(error)
```

```
## [1] 0.4517906
```

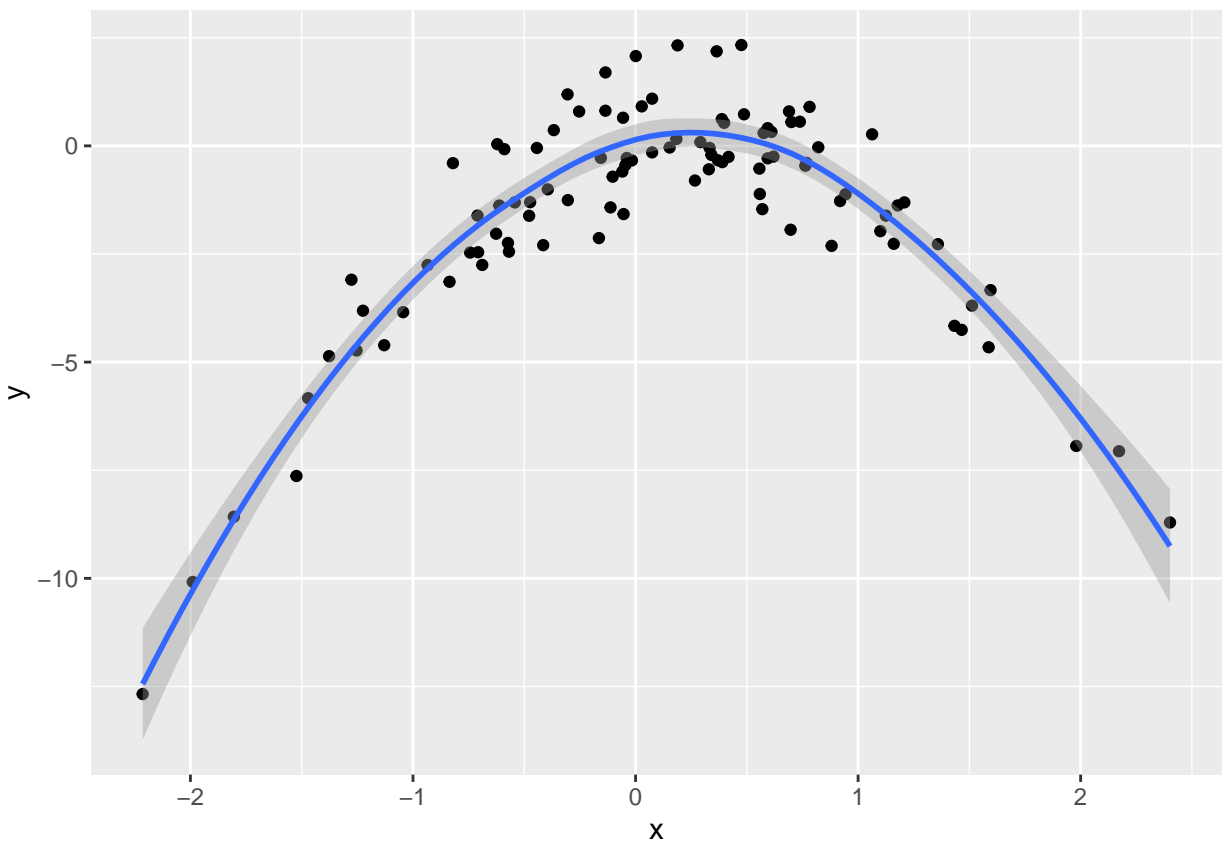
8. Next we cross-validate simulated data. The relationship is quadratic, obviously.

```
set.seed(1)
x=rnorm(100)
y=x-2*x^2+rnorm(100)

dat <- data.frame(x, y)

dat %>% ggplot(aes(x, y))+
  geom_point() +
  geom_smooth()
```

```
## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'
```



The true function here is:

$$y = 2x - 2x^2 + \epsilon$$

$$y = 0 + \epsilon$$

I had to look up the answer for this one. `rnorm`

After some aggressive programming, I find that the linear model's MSE is terrible, but the others are comparable. The random seed is irrelevant, since it plays no role in fitting the models.

The standard errors for the polynomial models are actually higher, but much smaller relative to the coefficients. I don't know how polynomial terms affect the SE formula, for simple regression square root of the ratio of RSS to TSS divided by degrees of freedom. I completely forgot to use the secret shortcut formula, which does work for polynomial regression.

```
set.seed(3454354)

forms <- append(expr(y~x), map(2:4, ~substitute(y ~poly(x, n), list(n = as.numeric(.x))))) %>%
  set_names(c("linear", "quad", "cubic", "quartic"))

map(forms, ~lm(data = dat, formula = .x))

## $linear
##
## Call:
## lm(formula = .x, data = dat)
##
## Coefficients:
## (Intercept)          x
##      -1.6254       0.6925
##
##
## $quad
##
## Call:
## lm(formula = .x, data = dat)
##
## Coefficients:
## (Intercept) poly(x, 2)1 poly(x, 2)2
##      -1.550       6.189      -23.948
##
##
## $cubic
##
## Call:
## lm(formula = .x, data = dat)
##
## Coefficients:
## (Intercept) poly(x, 3)1 poly(x, 3)2 poly(x, 3)3
##      -1.5500       6.1888      -23.9483       0.2641
##
##
## $quartic
##
## Call:
## lm(formula = .x, data = dat)
##
## Coefficients:
## (Intercept) poly(x, 4)1 poly(x, 4)2 poly(x, 4)3 poly(x, 4)4
##      -1.5500       6.1888      -23.9483       0.2641       1.2571
```

```

get_loocv <- function(df, form){

  resp <- form[[2]]
  out <- map2_dbl(list(form), 1:nrow(df), function(form, exclude){
    mod <- lm(data = df[-exclude], formula = form)
    error <- (predict(mod, df)[exclude] -df[exclude,as.character(resp)])^2
    error
  })
  out
}

map(forms, ~get_loocv(df = dat, .x)) %>% map_dbl(mean)

```

```

##      linear      quad      cubic    quartic
## 6.6255041 0.8902911 0.8895936 0.8737907

```

```

map(forms, ~lm(data = dat, formula = .x)) %>% map(summary)

```

```

## $linear
##
## Call:
## lm(formula = .x, data = dat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9.5161 -0.6800  0.6812  1.5491  3.8183
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -1.6254     0.2619  -6.205 1.31e-08 ***
## x              0.6925     0.2909   2.380  0.0192 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.6 on 98 degrees of freedom
## Multiple R-squared:  0.05465,    Adjusted R-squared:  0.045
## F-statistic: 5.665 on 1 and 98 DF,  p-value: 0.01924
##
## $quad
##
## Call:
## lm(formula = .x, data = dat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.9650 -0.6254 -0.1288  0.5803  2.2700
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -1.5500     0.0958  -16.18 < 2e-16 ***
## poly(x, 2)1   6.1888     0.9580   6.46 4.18e-09 ***

```



```

## poly(x, 2)2 -23.9483      0.9580  -25.00  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.958 on 97 degrees of freedom
## Multiple R-squared:  0.873, Adjusted R-squared:  0.8704
## F-statistic: 333.3 on 2 and 97 DF,  p-value: < 2.2e-16
##
##
## $cubic
##
## Call:
## lm(formula = .x, data = dat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.9765 -0.6302 -0.1227  0.5545  2.2843
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -1.55002    0.09626  -16.102  < 2e-16 ***
## poly(x, 3)1    6.18883    0.96263   6.429 4.97e-09 ***
## poly(x, 3)2  -23.94830    0.96263  -24.878  < 2e-16 ***
## poly(x, 3)3    0.26411    0.96263   0.274   0.784
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.9626 on 96 degrees of freedom
## Multiple R-squared:  0.8731, Adjusted R-squared:  0.8691
## F-statistic: 220.1 on 3 and 96 DF,  p-value: < 2.2e-16
##
##
## $quartic
##
## Call:
## lm(formula = .x, data = dat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.0550 -0.6212 -0.1567  0.5952  2.2267
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -1.55002    0.09591  -16.162  < 2e-16 ***
## poly(x, 4)1    6.18883    0.95905   6.453 4.59e-09 ***
## poly(x, 4)2  -23.94830    0.95905  -24.971  < 2e-16 ***
## poly(x, 4)3    0.26411    0.95905   0.275   0.784
## poly(x, 4)4    1.25710    0.95905   1.311   0.193
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.9591 on 95 degrees of freedom
## Multiple R-squared:  0.8753, Adjusted R-squared:  0.8701
## F-statistic: 166.7 on 4 and 95 DF,  p-value: < 2.2e-16

```

9. The last problem asks us to find a bootstrap estimate of a population mean. the bootstrap gives lower SE's than a difference of means test.

```
library(MASS)

##
## Attaching package: 'MASS'

## The following object is masked from 'package:dplyr':
##
##      select

boston <- Boston

mu_hat <- mean(boston$medv)
mu_hat

## [1] 22.53281

se_hat <- with(boston, sd(medv)/sqrt(length(medv)))
se_hat

## [1] 0.4088611

boot_mean <- function(vec, inds)
  mean(vec[inds])
booted <- boot(boston$medv, boot_mean, 1000)
booted

##
## ORDINARY NONPARAMETRIC BOOTSTRAP
##
##
## Call:
## boot(data = boston$medv, statistic = boot_mean, R = 1000)
##
##
## Bootstrap Statistics :
##      original      bias    std. error
## t1* 22.53281 0.005100198  0.3967026

t.test(Boston$medv) %>% broom::glance()

## Warning: '...' is not empty.
##
## We detected these problematic arguments:
## * 'needs_dots'
##
## These dots only exist to allow future extensions and should be empty.
## Did you misspecify an argument?
```

```
## # A tibble: 1 x 8
##   estimate statistic    p.value parameter conf.low conf.high method  alternative
##   <dbl>      <dbl>      <dbl>      <dbl>      <dbl>      <dbl> <chr>      <chr>
## 1      22.5      55.1 9.37e-216        505      21.7      23.3 One Sam~ two.sided
```

Now the confidence interval of the bootstrap. I can't figure out how to extract it, but it looks about the same as the sample CI.

```
boot_med <- function(vec, inds) median(vec[inds])
boot_quant <- function(n){
  force(n)
  function(vec, inds){
    quantile(vec[inds], probs = n)
  }
}

median <- boot(boston$medv, boot_med, 1000)
median
```

```
##
## ORDINARY NONPARAMETRIC BOOTSTRAP
##
##
## Call:
## boot(data = boston$medv, statistic = boot_med, R = 1000)
##
##
## Bootstrap Statistics :
##   original    bias    std. error
## t1*      21.2 -0.00565   0.3762474
```

```
boot(boston$medv, boot_quant(.1), 1000)
```

```
##
## ORDINARY NONPARAMETRIC BOOTSTRAP
##
##
## Call:
## boot(data = boston$medv, statistic = boot_quant(0.1), R = 1000)
##
##
## Bootstrap Statistics :
##   original    bias    std. error
## t1*      12.75 -0.0043   0.5022437
```

```
c(quantile(boston$medv, .1) - 2*se_hat, quantile(boston$medv, .1) + 2* se_hat)
```

```
##      10%      10%
## 11.93228 13.56772
```