

STAT346: Statistical Data Science I

Final: Thursday, Dec 15, 2022, 02:00–03:15 p.m.

Instructions

1. This exam covers course materials (Chapters 11–18).
2. You may use any books or online resources you want during this examination, but you may not communicate with any person other than your examiner or your TAs.
3. You are required to use the RStudio IDE for this exam.
4. You should work on the provided exam template. When you finalize your exam, you should submit your paper in pdf as well as its .rmd source file. They should have the following name:
 - stat346_final_yourID.pdf
 - stat346_final_yourID.rmd
5. You should submit your paper no later than 3:20 p.m. After that, there will be a deduction for the late submission (2 points per 1 minute). Still you have to submit your paper by 3:30 p.m.

Problem Set #0 (5 Points)

Run the following code, and show the result.

```
rm(list = ls())  
ls()
```

```
## character(0)
```

Problem Set #1 (20 Points)

We use the Teams, Batting and Salaries data by calling `library(Lahman)`:

```
library(Lahman)
library(broom)
data(Teams)
data(Batting)
data(Salaries)
```

For (a)-(b), use the Teams data.

- (a) [5 point] We use data from 1961 to 2001, and use BB, `singles`, `doubles`, `triples`, and HR (per game) as explanatory variables to predict R (run per game). Please complete the following program first, and provide the fitted regression formula.

```
fit = Teams %>% filter(yearID %in% 1961:2001) %>%
  mutate(BB = BB/G,
         singles = (H-X2B-X3B-HR)/G,
         doubles = X2B/G,
         triples = X3B/G,
         HR = HR/G,
         R = R/G) %>%
  lm(R ~ BB + singles + doubles + triples + HR, data = .)
tidy(fit, conf.int = T)
```

```
## # A tibble: 6 x 7
##   term          estimate std.error statistic    p.value conf.low conf.high
##   <chr>          <dbl>     <dbl>     <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)   -2.77      0.0862    -32.1 4.76e-157  -2.94    -2.60
## 2 BB            0.371     0.0117     31.6 1.87e-153   0.348     0.394
## 3 singles       0.519     0.0127     40.8 8.67e-217   0.494     0.544
## 4 doubles       0.771     0.0226     34.1 8.44e-171   0.727     0.816
## 5 triples       1.24      0.0768     16.1 2.12e- 52   1.09      1.39
## 6 HR            1.44      0.0243     59.3 0           1.40      1.49
```

Answer: $\hat{R} = -2.77 + 0.37BB + 0.52singles + 0.77doubles + 1.24triples + 1.44HR$

- (b) [10 points] Suppose we obtained the average number of team plate appearances per game as 38.74. Then we compute the per-plate-appearance rates for players on data from 1998-2001. Then we filter out players with less than 200 plate appearances per year and compute player-specific predicted runs in the object `R_hat`. Then make a ggplot of `R_hat` versus `R`. Interpret the prediction result using `R` and `R_hat` for the first player `abreubo01`.

```
pa_per_game = 38.74
players = Batting %>%
  filter(yearID %in% 1998:2001) %>%
```

```

group_by(playerID) %>%
mutate(PA = BB+AB) %>%
summarize(G=sum(PA)/pa_per_game,
          BB = sum(BB)/G, singles = sum(H-X2B-X3B-HR)/G,
          doubles = sum(X2B)/G, triples = sum(X3B)/G,
          HR = sum(HR)/G, AVG = sum(H)/sum(AB),
          PA = sum(PA), R = sum(R)/G) %>%
filter(PA >= 800) %>%
select(-G) %>%
mutate(R_hat = predict(fit, newdata = .))
head(players)

```

```

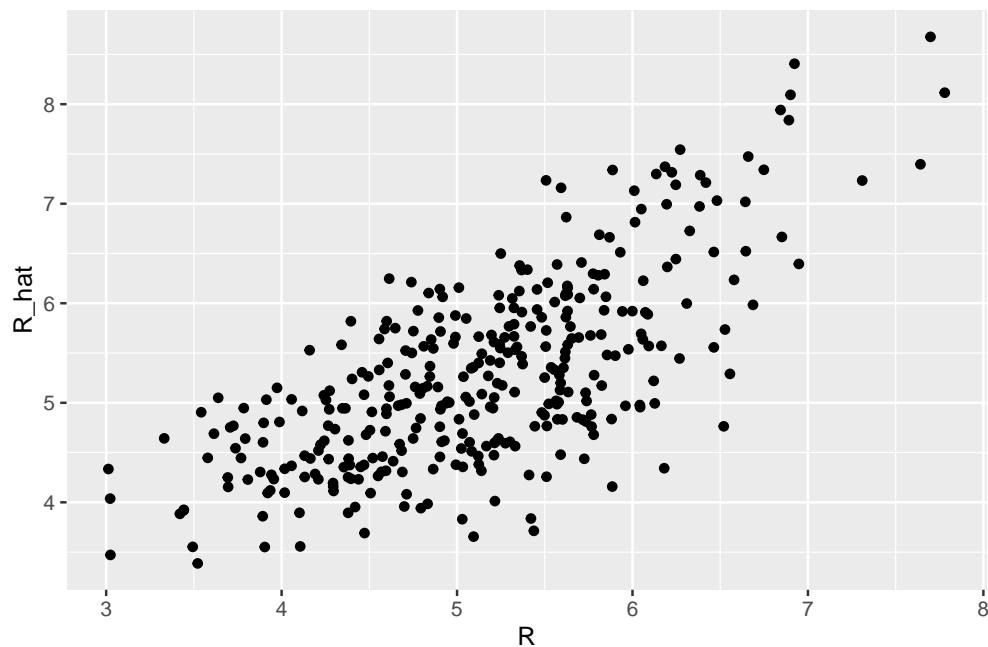
## # A tibble: 6 x 10
##   playerID      BB singles doubles triples   HR   AVG   PA    R R_hat
##   <chr>      <dbl>  <dbl>  <dbl>  <dbl> <dbl> <dbl> <int> <dbl> <dbl>
## 1 abreubo01  5.93    6.12    2.29  0.461  1.38  0.313  2606  6.05  6.95
## 2 agbaybe01  4.50    6.29    1.86  0.219  1.28  0.282  1060  4.75  5.72
## 3 alfoned01  4.62    6.42    2.04  0.0780 1.34  0.290  2482  6.09  5.89
## 4 alicelu01  3.80    6.96    1.71  0.388  0.492 0.273  1497  5.51  4.77
## 5 alomaro01  4.43    7.12    2.20  0.264  1.13  0.312  2638  6.58  6.24
## 6 alomasa02  1.65    6.74    2.08  0.165  0.760 0.263  1172  4.13  4.25

```

```

ggplot(aes(x = R, y = R_hat), data = players) +
geom_point()

```



Answer: For the first player abreubo01, Real value of R is 6.05 and Predicted value of R (R_{hat}) is 6.95.

Predicted value of R is overestimated (by 0.9) to real value of R.

- (c) [5 points] Select three variables `yearID`, `playerID` and `salary` in the data `Salaries`. Then join `players` for the year 2002 with this data for the year 2002 by `playerID`. Here we keep the rows in `players` data. Show the first 6 observations.

```
Salaries1 <- Salaries %>%
  filter(yearID == 2002) %>%
  select(yearID, playerID, salary)

left_join(players, Salaries1, by = "playerID") %>%
  head(6)
```

```
## # A tibble: 6 x 12
##   playerID      BB singles doubles triples   HR   AVG   PA    R R_hat yearID
##   <chr>      <dbl>   <dbl>   <dbl>   <dbl> <dbl> <dbl> <int> <dbl> <dbl>   <int>
## 1 abreu01    5.93     6.12     2.29   0.461   1.38  0.313  2606  6.05  6.95   2002
## 2 agbaybe01  4.50     6.29     1.86   0.219   1.28  0.282  1060  4.75  5.72   2002
## 3 alfoned01  4.62     6.42     2.04   0.0780  1.34  0.290  2482  6.09  5.89   2002
## 4 alicelu01  3.80     6.96     1.71   0.388   0.492  0.273  1497  5.51  4.77   2002
## 5 alomaro01  4.43     7.12     2.20   0.264   1.13  0.312  2638  6.58  6.24   2002
## 6 alomasa02  1.65     6.74     2.08   0.165   0.760  0.263  1172  4.13  4.25   2002
## # ... with 1 more variable: salary <int>
```

Problem Set #2 (40 Points)

In every problems (a)–(e), set the seed number as 2022 at the beginning.

- (a) [5 points] Generate the data consisting of zeroes and ones where we obtain 1 with probability 0.3. The size of the data is 10^4 and we treat this data as our population. Save this population data into the object name `vote`. Check that the true proportion p is 0.295.

```
set.seed(2022)
vote <- sample(c(0, 1), size = 10^4, replace = TRUE, prob = c(0.7, 0.3))
p <- mean(vote)
p
```

```
## [1] 0.295
```

- (b) [5 points] Take a random sample of size 100 (without replacement) from the `vote`, and calculate the sample proportion.

```
set.seed(2022)
prop <- sample(vote, size = 100)
mean(prop)
```

```
## [1] 0.25
```

Answer: The sample proportion is 0.25.

- (c) [10 points] Run a Monte Carlo simulation to confirm that a 90% confidence interval includes p about 90% of the time. For this, repeat (1000 times) the random sampling from the population data `vote` where the sample size is 100 and calculate the sample mean and check if the true proportion is included or not in the approximate 90% confidence interval based on the sample data.

```
set.seed(2022)
n = 100 #sample size
B = 1000 #iteration number
inside = replicate(B, {
  x = sample(c(0,1), size = n, replace = TRUE,
             prob = c(1-p, p)) #random sample from vote
  x_hat = mean(x) # compute the average
  se_hat = sqrt(x_hat * (1 - x_hat) / n) # compute the standard error
  between(p, x_hat + qnorm(0.05) * se_hat,
          x_hat + qnorm(0.95) * se_hat) # use `between` function
})
mean(inside)
```

```
## [1] 0.888
```

Answer: 90% confidence interval includes p about 88.8% of the time

- (d) [10 points] Use a Monte Carlo simulation to learn the distribution of the sample mean of size 100. Obtain 90% confidence interval with the iteration number as 10^4 .

```
set.seed(2022)
B = 10^4
vote_p = replicate(B, {
  X = sample(vote, size = 100) # random sample from vote
  mean(X) # compute the average
})
quantile(vote_p, c(0.05, 0.95)) # compute the 90% confidence interval
```

```
##    5%  95%
## 0.22 0.37
```

Answer: [0.22, 0.37]

- (e) [10 points] Set the seed number as 2022 and obtain 90% confidence interval for the true proportion using bootstrap. Use the bootstrap iteration number as 10^4 . Compare this with 90% confidence interval based on the CLT.

```
set.seed(2022)
X <- sample(vote, size = n)
B = 10^4
vote_p_star = replicate(B, {
  X_star = sample(X, size = n, replace = TRUE)
  mean(X_star)
})
quantile(vote_p_star, c(0.05, 0.95))
```

```
##    5%  95%
## 0.18 0.32
```

```
mean(X) + qnorm(0.95) * sd(X) / sqrt(100) * c(-1, 1)
```

```
## [1] 0.1784169 0.3215831
```

Answer: 90% confidence interval for the true proportion using bootstrap is [0.18, 0.32].

90% confidence interval by CLT is [0.1784169, 0.3215831].

a confidence interval constructed with the bootstrap is much closer to one constructed with the theoretical distribution.

Problem Set #3 (15 Points)

We consider UCBAmissions data set, which contains data on applicants to graduate school at Berkeley for the six largest departments in 1973 classified by admission and gender.

```
data(UCBAmissions)
dat = as.data.frame(UCBAmissions)[c("Dept", "Admit", "Gender", "Freq")]
```

- (a) [5 points] If we think of an observation as a Dept, then this data `dat` is not tidy. Use the `pivot_wider` function to wrangle into tidy shape: one row for each major with variable names as `Admitted_Male`, `Rejected_Male`, `Admitted_Female`, and `Rejected_Female`. Save this data set as an object `dat1` (as shown below).

```
dat1 <- dat %>%
  pivot_wider(names_from = c(Admit, Gender), values_from = Freq)
dat1
```

```
## # A tibble: 6 x 5
##   Dept Admitted_Male Rejected_Male Admitted_Female Rejected_Female
##   <fct>      <dbl>      <dbl>          <dbl>          <dbl>
## 1 A           512         313             89             19
## 2 B           353         207             17              8
## 3 C           120         205            202            391
## 4 D           138         279            131            244
## 5 E            53         138             94            299
## 6 F            22         351             24            317
```

- (b) [5 points] Now use `pivot_longer` to wrangle the above `dat1` so that we obtain the following `dat2` (head part is shown).

```
dat2 <- dat1 %>%
  pivot_longer(-Dept, names_to = 'Admit_Gender', values_to = 'Freq')
head(dat2)
```

```
## # A tibble: 6 x 3
##   Dept Admit_Gender  Freq
##   <fct> <chr>      <dbl>
## 1 A     Admitted_Male  512
## 2 A     Rejected_Male  313
## 3 A     Admitted_Female  89
## 4 A     Rejected_Female  19
## 5 B     Admitted_Male  353
## 6 B     Rejected_Male  207
```

- (c) [5 points] Note that the above `dat2` has a variable `Admit_Gender` containing information from the two variables `Admit` and `Gender`. Use `separate` function to obtain the same data as `dat` as shown below (only the head part is shown and column types may be different from those in 'dat').

```
dat2 %>%  
  separate(Admit_Gender, c('Admit', 'Gender')) %>%  
  head()
```

```
## # A tibble: 6 x 4  
##   Dept Admit   Gender Freq  
##   <fct> <chr>   <chr> <dbl>  
## 1 A     Admitted Male    512  
## 2 A     Rejected Male    313  
## 3 A     Admitted Female    89  
## 4 A     Rejected Female    19  
## 5 B     Admitted Male   353  
## 6 B     Rejected Male   207
```


Problem Set #4 (20 Points)

For this problem, we use actual polls from the 2016 election. We will use all the national polls.

```
library(dslabs)
library(lubridate)
data("polls_us_election_2016")
polls = polls_us_election_2016 %>% filter(state == "U.S.")
```

- (a) [5 points] Sort `startdate` and show the first three dates. Consider the variable `enddate`, and find the proportion of October 2016.

```
polls %>%
  arrange(startdate) %>%
  head(3)
```

```
##   state startdate   enddate      pollster grade samplesize population
## 1  U.S. 2015-11-13 2015-11-16 Morning Consult <NA>      2001         rv
## 2  U.S. 2015-11-15 2015-12-02   Marist College    A       2360         rv
## 3  U.S. 2015-11-16 2015-11-17 Public Policy Polling B+      1360          v
##   rawpoll_clinton rawpoll_trump rawpoll_johnson rawpoll_mcmullin
## 1              44             43              NA              NA
## 2              52             41              NA              NA
## 3              45             44              NA              NA
##   adjpoll_clinton adjpoll_trump adjpoll_johnson adjpoll_mcmullin
## 1      45.14983      44.84425              NA              NA
## 2      50.48969      42.09936              NA              NA
## 3      43.56539      42.89190              NA              NA
```

Answer: 2015-11-13, 2015-11-15, 2015-11-16

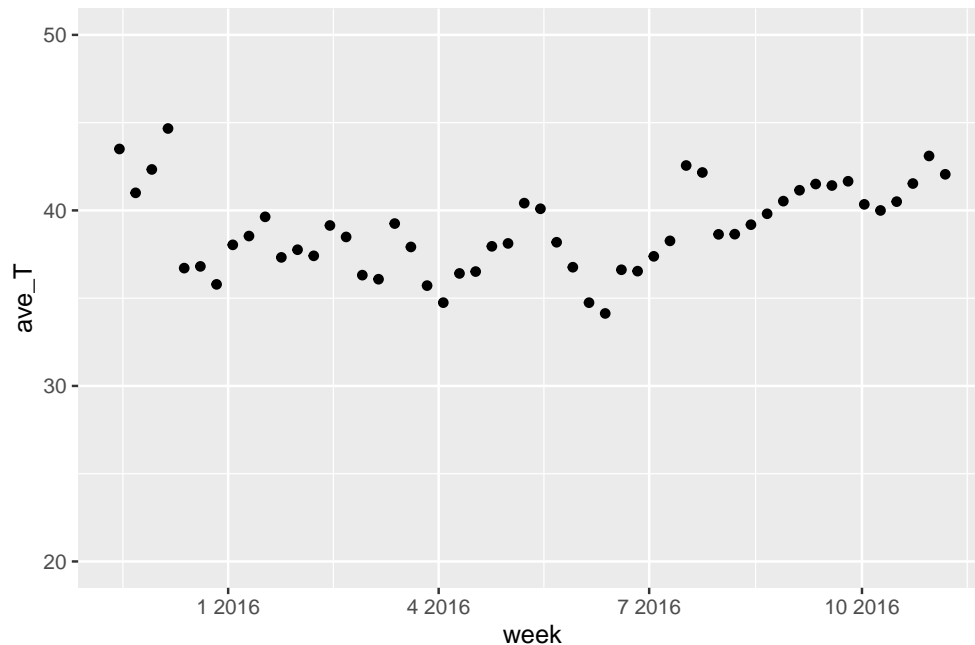
```
polls %>%
  summarize(proportion =
    sum(year(enddate) == 2016 & month(enddate) == 10)
    / length(enddate))
```

```
##   proportion
## 1 0.2106691
```

Answer: 21.06691%

- (b) [5 points] Group all the polls by week of the year using `round_date` function (with `startdate` variable), and compute the average `rawpoll_trump` per week (saving this as `ave_T`), and draw a scatter plot using `qplot` function.

```
polls %>%
  mutate(week = round_date(startdate, 'week')) %>%
  group_by(week) %>%
  summarize(ave_T = mean(rawpoll_trump)) %>%
  qplot(week, ave_T, data=., ylim = c(20,50))
```



- (c) [10 points] Provide the loess fit for `rawpoll_trump` using `span = 0.2` and `degree = 1`. In order to use the loess function, we make a new variable `day` and fit a smooth line of `rawpoll_trump` versus `day`. Provide the following plot.

```
span = 0.2
polls$day = as.numeric(polls$startdate)
fit = loess(rawpoll_trump ~ day, degree = 1,
            span = span, data = polls)
polls %>% mutate(smooth = fit$fitted) %>%
  ggplot(aes(startdate, rawpoll_trump)) +
  geom_point(size = 3, alpha = 0.5, color = "grey") +
  geom_line(aes(startdate, smooth), color = "red")
```

