

Homework 1: review and coding questions

Solution Set

This assignment is meant to serve multiple objectives:

- You will gain familiarity with R, RStudio, R Markdown, and GitHub
- You will perform at least one iteration of the courses's data science inspired workflow
- You will gain experience with typing mathematics
- You will learn some `dplyr` and `ggplot2` basics

STAT 528 is a collaborative course environment, especially for assignments that involve coding, modeling, and/or data analysis. You are encouraged to ask for help from other students. Coding and data science work flow can be very tedious. Having someone else look over your work or answering a basic question can save you a lot of time. However, direct copying is not accepting. All final work must be your own.

Mathematical review questions

Problem 1: Prove that the Binomial distribution arises as a sum of n iid Bernoulli trials each with success probability p .

Answer to P1: The moment generating function of the sum of n iid Bernoulli random variables X_1, \dots, X_n is

$$\mathbb{E} \left(e^{t(X_1 + \dots + X_n)} \right) = \mathbb{E} \left(e^{tX_1} \right) \dots \mathbb{E} \left(e^{tX_n} \right) = \left(\mathbb{E} \left(e^{tX_1} \right) \right)^n = (pe^t + 1 - p)^n,$$

which is equal to the moment generating function of a Binomial random variable.

Problem 2: Let $l(\theta)$ denote a twice continuously differentiable log likelihood corresponding to an iid sample under density f_θ where n is the sample size. The score function is defined as

$$u(\theta) = \frac{\partial l(\theta)}{\partial \theta},$$

and the Fisher information matrix is defined as

$$I(\theta) = -\mathbb{E} \left(\frac{\partial^2 l(\theta)}{\partial \theta^2} \right),$$

where the expectation is over the assumed distribution for the data when the parameter value is θ . Prove that

$$\mathbb{E}(u(\theta)) = 0 \quad \text{and} \quad \text{Var}(u(\theta)) = I(\theta).$$

i)

The score function is,

$$u(\theta) = \frac{\partial l(\theta)}{\partial \theta} = \frac{\partial}{\partial \theta} \log \left[\prod_{i=1}^n f_{\theta}(x_i) \right] = \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f_{\theta}(x_i) = \sum_{i=1}^n \frac{\frac{\partial}{\partial \theta} f_{\theta}(x_i)}{f_{\theta}(x_i)}$$

So we have,

$$\begin{aligned} E(u(\theta)) &= E\left(\frac{\partial l(\theta)}{\partial \theta}\right) \\ &= \sum \int_{\chi} \frac{\frac{\partial}{\partial \theta} f_{\theta}(x)}{f_{\theta}(x)} f_{\theta}(x) dx \\ &= \sum \int_{\chi} \frac{\partial}{\partial \theta} f_{\theta}(x) dx \\ &= \sum \frac{\partial}{\partial \theta} \int_{\chi} f_{\theta}(x) dx \\ &= \sum \frac{\partial}{\partial \theta} 1 \\ &= 0 \end{aligned}$$

ii)

By the results we just proved in i),

$$\begin{aligned} 0 &= E(u(\theta)) \\ 0 &= \frac{\partial}{\partial \theta} E(u(\theta)) \\ &= \frac{\partial}{\partial \theta} E\left(\sum \frac{\partial}{\partial \theta} \log f_{\theta}(x)\right) \\ &= \frac{\partial}{\partial \theta} \int_{\chi} \sum \frac{\partial \log f_{\theta}(x)}{\partial \theta} f_{\theta}(x) dx \\ &= \sum \int_{\chi} \frac{\partial}{\partial \theta} \left[\frac{\partial \log f_{\theta}(x)}{\partial \theta} f_{\theta}(x) \right] dx \\ &= \sum \int_{\chi} \left[\frac{\partial^2 \log f_{\theta}(x)}{\partial \theta^2} f_{\theta}(x) + \frac{\partial \log f_{\theta}(x)}{\partial \theta} \frac{\partial f_{\theta}(x)}{\partial \theta} \right] dx \\ &= \sum \int_{\chi} \frac{\partial^2 \log f_{\theta}(x)}{\partial \theta^2} f_{\theta}(x) dx + \sum \int_{\chi} \frac{\frac{\partial}{\partial \theta} f_{\theta}(x)}{f_{\theta}(x)} \frac{\frac{\partial}{\partial \theta} f_{\theta}(x)}{f_{\theta}(x)} f_{\theta}(x) dx \end{aligned}$$

Rearrange the above equality to get,

$$\begin{aligned} - \sum \int_{\chi} \frac{\partial^2 \log f_{\theta}(x)}{\partial \theta^2} f_{\theta}(x) dx &= \sum \int_{\chi} \frac{\frac{\partial}{\partial \theta} f_{\theta}(x)}{f_{\theta}(x)} \frac{\frac{\partial}{\partial \theta} f_{\theta}(x)}{f_{\theta}(x)} f_{\theta}(x) dx \\ &= E\left(\left(\frac{\partial l(\theta)}{\partial \theta}\right)^2\right) \\ I(\theta) &= E[u(\theta)]^2 \end{aligned}$$

Finally, note that $E[u(\theta)] = 0$, hence we arrive at,

$$I(\theta) = E[u(\theta)]^2 = E[u(\theta) - E(u(\theta))]^2 = \text{Var}(u(\theta))$$

Coding questions

Problem 3: The data we will use to accomplish this task will come from Lahman’s Baseball Database. Thankfully, there is an R package, `Lahman`, that makes importing this data into R very easy. If you have not done so previously, install this package using:

```
install.packages("Lahman")
```

While there many metrics that could be used to determine who is the “best” baseball player, because we are focusing on `batters`, we will use the `on-base plus slugging (OPS)` statistic. This statistic measures both a batter’s ability to “get on base” and “hit for power.”

- [YouTube: Moneyball, “He Gets on Base”](#)

Additionally, our definition of “best” will be based on a player’s career statistics, but an alternative argument could be made based on single season efforts.

After loading the Lahman package, you will have access to several data frames containing historical baseball data from 1871 - 2022. You will need to interact with the following data frames:

- `Schools`
- `CollegePlaying`
- `Batting`
- `People`

You should spend some time exploring these datasets and reading the relevant documentation.

Create a tibble named `illini_mlb_batters` that contains the following elements, in this order:

- `playerID`
- `nameFirst`
- `nameLast`
- `birthYear`
- `G`
- `AB`
- `R`
- `H`
- `X2B`
- `X3B`
- `HR`
- `RBI`
- `SB`
- `CS`
- `BB`
- `SO`
- `IBB`
- `HBP`
- `SH`
- `SF`
- `GIDP`
- `PA`
- `TB`
- `BA`

- OBP
- SLG
- OPS

The rows of the tibble should be sorted from highest OPS to lowest OPS. Each row should represent the career statistics for the player with ID playerID. Only include players that had at least one at-bat and one plate appearance. Except for PA, TB, AVG, OBP, SLG, and OPS, the (sometimes season-level) variables listed can be found in one of the four data frames listed above. The remaining values can be calculated as follows:

- $PA = AB + BB + HBP + SH + SF$
- $TB = H + X2B + 2 * X3B + 3 * HR$
- $BA = H / AB$
- $OBP = (H + BB + HBP) / (PA - SH)$
- $SLG = TB / AB$
- $OPS = OBP + SLG$

Round any rate statistics to three decimals places, as is customary in baseball.

Answer to P3:

```
library(Lahman)
library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.1 --

## v ggplot2 3.4.3      v purrr 1.0.2
## v tibble 3.2.1       v dplyr 1.1.3
## v tidyr 1.3.0        v stringr 1.5.0
## v readr 1.4.0        v forcats 0.5.1

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()

as_tibble(People)

## # A tibble: 20,676 x 26
##   playerID birthYear birthMonth birthDay birthCountry birthState birthCity
##   <chr>      <int>      <int>    <int> <chr>          <chr>    <chr>
## 1 aardsda01 1981         12      27 USA          CO      Denver
## 2 aaronha01 1934          2       5 USA          AL      Mobile
## 3 aaronto01 1939          8       5 USA          AL      Mobile
## 4 aasedo01 1954          9       8 USA          CA      Orange
## 5 abadan01 1972          8      25 USA          FL      Palm Beach
## 6 abadfe01 1985         12      17 D.R.        La Romana La Romana
## 7 abadijo01 1850         11       4 USA          PA      Philadelphia
## 8 abbated01 1877          4      15 USA          PA      Latrobe
## 9 abbeybe01 1869         11      11 USA          VT      Essex
## 10 abbeych01 1866         10      14 USA          NE      Falls City
## # i 20,666 more rows
## # i 19 more variables: deathYear <int>, deathMonth <int>, deathDay <int>,
```

```
## # deathCountry <chr>, deathState <chr>, deathCity <chr>, nameFirst <chr>,
## # nameLast <chr>, nameGiven <chr>, weight <int>, height <int>, bats <fct>,
## # throws <fct>, debut <chr>, finalGame <chr>, retroID <chr>, bbrefID <chr>,
## # deathDate <date>, birthDate <date>
```

```
as_tibble(Batting)
```

```
## # A tibble: 112,184 x 22
##   playerID yearID stint teamID lgID      G    AB    R    H   X2B   X3B   HR
##   <chr>      <int> <int> <fct> <fct> <int> <int> <int> <int> <int> <int> <int>
## 1 abercda01  1871     1 TR0   NA      1     4     0     0     0     0     0
## 2 addybo01   1871     1 RC1   NA     25    118    30    32     6     0     0
## 3 allisar01  1871     1 CL1   NA     29    137    28    40     4     5     0
## 4 allisdo01  1871     1 WS3   NA     27    133    28    44    10     2     2
## 5 ansonca01  1871     1 RC1   NA     25    120    29    39    11     3     0
## 6 armstbo01  1871     1 FW1   NA     12     49     9    11     2     1     0
## 7 barkeal01  1871     1 RC1   NA      1     4     0     1     0     0     0
## 8 barnero01  1871     1 BS1   NA     31    157    66    63    10     9     0
## 9 barrebi01  1871     1 FW1   NA      1     5     1     1     1     0     0
## 10 barrofr01 1871     1 BS1   NA     18     86    13    13     2     1     0
## # i 112,174 more rows
## # i 10 more variables: RBI <int>, SB <int>, CS <int>, BB <int>, SO <int>,
## # IBB <int>, HBP <int>, SH <int>, SF <int>, GIDP <int>
```

```
as_tibble(Schools)
```

```
## # A tibble: 1,207 x 5
##   schoolID name_full      city      state country
##   <chr>      <chr>      <chr>      <chr> <chr>
## 1 abilchrist Abilene Christian University Abilene TX USA
## 2 adelphi Adelphi University Garden City NY USA
## 3 adrianmi Adrian College Adrian MI USA
## 4 akron University of Akron Akron OH USA
## 5 alabama University of Alabama Tuscaloosa AL USA
## 6 alabamaam Alabama A&M University Normal AL USA
## 7 alabamast Alabama State University Montgomery AL USA
## 8 albanyst Albany State University Albany GA USA
## 9 albertsnid Albertson College Caldwell ID USA
## 10 albevil Bevill State Community College Sumiton AL USA
## # i 1,197 more rows
```

```
as_tibble(CollegePlaying)
```

```
## # A tibble: 17,350 x 3
##   playerID schoolID yearID
##   <chr>      <chr>      <int>
## 1 aardsda01 pennst      2001
## 2 aardsda01 rice        2002
## 3 aardsda01 rice        2003
## 4 abadan01 gamiddl      1992
## 5 abadan01 gamiddl      1993
## 6 abbeybe01 vermont      1889
```

```
## 7 abbeybe01 vermont 1890
## 8 abbeybe01 vermont 1891
## 9 abbeybe01 vermont 1892
## 10 abbotje01 kentucky 1991
## # i 17,340 more rows
```

```
Schools %>%
  filter(city == "Champaign")
```

```
##   schoolID          name_full      city state country
## 1 illinois University of Illinois at Urbana-Champaign Champaign IL USA
## 2 ilparkl      Parkland College Champaign IL USA
```

```
foo = People %>%
  select(playerID, nameFirst, nameLast, birthYear)

illiniIDs = CollegePlaying %>%
  filter(schoolID == "illinois") %>%
  pull(playerID) %>%
  unique()

illini_mlb_batters = Batting %>%
  filter(playerID %in% illiniIDs) %>%
  select(playerID, G:GIDP) %>%
  mutate(across(G:GIDP, ~replace_na(.x,0))) %>%
  group_by(playerID) %>%
  summarise(across(G:GIDP,sum)) %>%
  mutate(PA = AB + BB + HBP + SH + SF,
         TB = H + X2B + 2 * X3B + 3 * HR,
         BA = H / AB,
         OBP = (H + BB + HBP) / (PA - SH),
         SLG = TB / AB,
         OPS = OBP + SLG) %>%
  left_join(foo, by = "playerID") %>%
  select(playerID, nameFirst, nameLast, birthYear, everything()) %>%
  mutate(across(BA:OPS, ~round(.x, 3))) %>%
  arrange(desc(OPS)) %>%
  filter(AB >= 1) %>%
  filter(PA >= 1)

print.data.frame(head(illini_mlb_batters))
```

```
##   playerID nameFirst nameLast birthYear    G  AB  R    H X2B X3B  HR RBI SB
## 1 boudrlo01      Lou Boudreau    1917 1646 6029 861 1779 385  66  68 789 51
## 2 eversho01      Hoot  Evers    1921 1142 3801 556 1055 187  41  98 565 45
## 3 halleto01      Tom  Haller    1937 1294 3935 461 1011 153  31 134 504 14
## 4 spiezsc01     Scott Spiezio    1972 1274 3899 517  996 225  27 119 549 33
## 5 mcurha01     Harry McCurdy    1899  543 1157 148  326  71  12   9 148 12
## 6 fletcda01    Darrin Fletcher    1966 1245 3902 377 1048 214   8 124 583  2
##   CS  BB  SO  IBB HBP  SH SF GIDP  PA  TB  BA  OBP  SLG  OPS
## 1 50 796 309   0 34 164   0 155 7023 2500 0.295 0.380 0.415 0.795
## 2 36 415 420   0 27  65   2 116 4310 1618 0.278 0.353 0.426 0.778
```

```
## 3 30 477 593 96 35 35 37 60 4519 1628 0.257 0.340 0.414 0.753
## 4 23 412 594 35 35 25 41 77 4412 1632 0.255 0.329 0.419 0.747
## 5 9 129 108 0 3 25 0 0 1314 448 0.282 0.355 0.387 0.743
## 6 6 255 399 31 49 13 51 122 4270 1650 0.269 0.318 0.423 0.740
```

Problem 4: The data we will use to accomplish this task will come from the **Teams** data frame in Lahman’s Baseball Database. In this problem we will visualize the [Pythagorean Theorem of Baseball](#). This “Theorem” states that winning percentage is given by the following nonlinear equation:

$$WP = \frac{R^2}{R^2 + RA^2}$$

where

- WP is winning percentage
- R is total runs scored by a baseball team
- RA is total runs allowed by a baseball team

For this problem, plot the estimated number of wins as predicted by the Pythagorean equation and actual wins (denoted W). The estimated number of wins as predicted by the Pythagorean equation

$$162 * \frac{R^2}{R^2 + RA^2}.$$

Provide a line of best fit. Restrict attention to the 1990 season and beyond. Note that there are two shortened seasons that need to be treated separately from the remaining seasons. These seasons are 1994 and 2020. The [1994 season](#) was cut short because of a labor strike. The [2020 season](#) was cut short due to COVID.

Answer to P4:

```
library(ggplot2)
as_tibble(Teams)
```

```
## # A tibble: 3,015 x 48
##   yearID lgID teamID franchID divID Rank   G Ghome   W   L DivWin WCWin
##   <int> <fct> <fct>   <fct>   <chr> <int> <int> <int> <int> <chr> <chr>
## 1  1871 NA   BS1     BNA     <NA>    3    31    NA   20   10 <NA> <NA>
## 2  1871 NA   CH1     CNA     <NA>    2    28    NA   19    9 <NA> <NA>
## 3  1871 NA   CL1     CFC     <NA>    8    29    NA   10   19 <NA> <NA>
## 4  1871 NA   FW1     KEK     <NA>    7    19    NA    7   12 <NA> <NA>
## 5  1871 NA   NY2     NNA     <NA>    5    33    NA   16   17 <NA> <NA>
## 6  1871 NA   PH1     PNA     <NA>    1    28    NA   21    7 <NA> <NA>
## 7  1871 NA   RC1     ROK     <NA>    9    25    NA    4   21 <NA> <NA>
## 8  1871 NA   TRO     TRO     <NA>    6    29    NA   13   15 <NA> <NA>
## 9  1871 NA   WS3     OLY     <NA>    4    32    NA   15   15 <NA> <NA>
## 10 1872 NA   BL1     BLC     <NA>    2    58    NA   35   19 <NA> <NA>
## # i 3,005 more rows
## # i 36 more variables: LgWin <chr>, WSWin <chr>, R <int>, AB <int>, H <int>,
## #   X2B <int>, X3B <int>, HR <int>, BB <int>, SO <int>, SB <int>, CS <int>,
## #   HBP <int>, SF <int>, RA <int>, ER <int>, ERA <dbl>, CG <int>, SHO <int>,
## #   SV <int>, IPouts <int>, HA <int>, HRA <int>, BBA <int>, SOA <int>, E <int>,
## #   DP <int>, FP <dbl>, name <chr>, park <chr>, attendance <int>, BPF <int>,
## #   PPF <int>, teamIDBR <chr>, teamIDlahman45 <chr>, teamIDretro <chr>
```

```
Teams %>%
  filter(yearID >= 1990) %>%
  select(yearID, W, R, RA) %>%
  mutate(WP = 162*R^2/(R^2+RA^2))%>%
  mutate(seasons = case_when(
    yearID == 1994 ~ "labor strike",
    yearID == 2020 ~ "COVID",
    .default = "normal"
  ))%>%
  ggplot() +
  aes(x = W, y = WP, color = seasons) +
  geom_point() +
  geom_smooth(method = "lm") +
  labs(title = "Linear Relationship between Pythagorean wins and Observed wins",
    x = "Observed wins", y = "Pythagorean wins") +
  theme_minimal()
```

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

