

Homework 2

Sam Rosenberg

Due October 4 at 11:59pm

1 Instructions

Setup. Pull the latest version of this assignment from Github and set your working directory to `stat-961-fall-2021/homework/homework-2`. Consult the [getting started guide](#) if you need to brush up on R, LaTeX, or Git.

Collaboration. The collaboration policy is as stated on the Syllabus:

“Students are permitted to work together on homework assignments, but solutions must be written up and submitted individually. Students must disclose any sources of assistance they received; furthermore, they are prohibited from verbatim copying from any source and from consulting solutions to problems that may be available online and/or from past iterations of the course.”

In accordance with this policy,

Please list anyone you discussed this homework with:

Please list what external references you consulted (e.g. articles, books, or websites):

Writeup. Use this document as a starting point for your writeup, adding your solutions between `\begin{sol}` and `\end{sol}`. See the [preparing reports guide](#) for guidance on compilation, creation of figures and tables, and presentation quality. Show all the code you wrote to produce your numerical results, and include complete derivations typeset in LaTeX for the mathematical questions.

Programming. The `tidyverse` paradigm for data manipulation (`dplyr`) and plotting (`ggplot2`) are strongly encouraged, but points will not be deducted for using base R.

```
library(tidyverse)
```

Grading. Each sub-part of each problem will be worth 3 points: 0 points for no solution or completely wrong solution; 1 point for some progress; 2 points for a mostly correct solution; 3 points for a complete and correct solution modulo small flaws. The presentation quality of the solution for each problem (as exemplified by the guidelines in Section 3 of the [preparing reports guide](#)) will be evaluated out of an additional 3 points.

Submission. Compile your writeup to PDF and submit to [Gradescope](#).

Problem 1. Likelihood inference in linear regression.

Let's consider the usual linear regression setup. Given a full-rank $n \times p$ model matrix \mathbf{X} , a coefficient vector $\boldsymbol{\beta} \in \mathbb{R}^p$, and a noise variance $\sigma^2 > 0$, suppose

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(0, \sigma^2 \mathbf{I}_n). \quad (1)$$

The goal of this problem is to connect linear regression inference with classical likelihood-based inference (below is a quick refresher).

- For the sake of simplicity, let's start by assuming σ^2 is known. Under the fixed-design model, why does the linear regression model (1) not fit into the classical inferential setup (2)? Write the linear model in as close a form as possible to (2).
- Continue assuming that σ^2 is known. Why does the Fisher information (4) not immediately make sense for the linear regression model? Propose and compute an analog to this quantity, and using this quantity exhibit a result analogous to the asymptotic normality (3).
- Now assume that neither $\boldsymbol{\beta}$ nor σ^2 is known. Derive the maximum likelihood estimates for $(\boldsymbol{\beta}, \sigma^2)$. How do these compare to the estimates $(\hat{\boldsymbol{\beta}}, \hat{\sigma}^2)$ discussed in class?
- Continuing to assume that neither $\boldsymbol{\beta}$ nor σ^2 is known, consider the null hypothesis $H_0 : \boldsymbol{\beta}_S = \mathbf{0}$ for some $S \subseteq \{1, \dots, p\}$. Write this hypothesis in the form (5), and derive the likelihood ratio test for this hypothesis. Discuss the connection of this test with the F -test.

Refresher on likelihood inference. In classical likelihood inference, we have observations

$$y_i \stackrel{\text{i.i.d.}}{\sim} p_{\boldsymbol{\theta}}, \quad i = 1, \dots, n \quad (2)$$

from some model parameterized by a vector $\boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^d$. Under regularity conditions, the maximum likelihood estimate $\hat{\boldsymbol{\theta}}_n$ is known to converge to a normal distribution centered at its true value:

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}) \xrightarrow{d} N(0, \mathbf{I}(\boldsymbol{\theta})^{-1}), \quad (3)$$

where

$$\mathbf{I}(\boldsymbol{\theta}) \equiv -\mathbb{E}_{\boldsymbol{\theta}} \left[\frac{\partial^2}{\partial \boldsymbol{\theta}^2} \log p_{\boldsymbol{\theta}}(y) \right] \quad (4)$$

is the Fisher information matrix. Furthermore, an optimal test of the null hypothesis

$$H_0 : \boldsymbol{\theta} \in \Theta_0 \quad \text{versus} \quad H_1 : \boldsymbol{\theta} \in \Theta_1 \quad (5)$$

for some $\Theta_0 \subseteq \Theta_1 \subseteq \Theta$ is the likelihood ratio test based on the test statistic

$$\Lambda = \frac{\max_{\boldsymbol{\theta} \in \Theta_1} \prod_{i=1}^n p_{\boldsymbol{\theta}}(y_i)}{\max_{\boldsymbol{\theta} \in \Theta_0} \prod_{i=1}^n p_{\boldsymbol{\theta}}(y_i)}. \quad (6)$$

Under H_0 , we have the convergence

$$2 \log \Lambda \xrightarrow{d} \chi_k^2, \quad \text{where} \quad k \equiv \dim(\Theta_1) - \dim(\Theta_0). \quad (7)$$

Solution 1.

(a) We can rewrite (1) as follows:

$$\begin{aligned} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} &= \begin{pmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{np} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{x}_1 \cdot \boldsymbol{\beta} + \epsilon_1 \\ \vdots \\ \mathbf{x}_n \cdot \boldsymbol{\beta} + \epsilon_n \end{pmatrix}, \end{aligned}$$

where the covariance matrix of $\boldsymbol{\epsilon}$ is $\sigma^2 \mathbf{I}_n$. Note that this means the ϵ_i are all uncorrelated, but the fact that they are multivariate normal and uncorrelated implies that the ϵ_i are independent standard normal random variables. So, $y_i \sim \mathbf{x}_i \cdot \boldsymbol{\beta} + \epsilon_i$, where $\epsilon_i \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma^2)$.

Though each y_i is parametrized by $\boldsymbol{\beta} \in \mathbb{R}^p$, it is also a function of \mathbf{x}_i , which is not necessarily the same for each i since the \mathbf{x}_i are regarded as fixed. As a result, the fixed-design model does not fit into the classical inferential setup, since the y_i need not be identically distributed.

- (b) The Fisher information does not make sense in this context because the y_i each have different PDFs, since they are not i.i.d. as noted in (a). We have previously seen that $\hat{\beta}_j \sim N(\beta_j, \sigma^2/s_j^2)$, so $\sqrt{n}(\hat{\beta}_j - \beta_j) \sim N(0, \sigma^2/(s_j/n)^2)$. So, an analog for the Fisher information would be $(s_j/n)^2/\sigma^2$.
- (c) Define our parameter vector of interest to be $\boldsymbol{\theta} := (\boldsymbol{\beta}, \sigma^2)$. We saw in Unit 2 that $\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$. So, we know that the joint density (equivalent to the likelihood function) for the fixed-design model is

$$\begin{aligned} p_{\boldsymbol{\theta}}(\mathbf{y}) &= \frac{1}{\sqrt{(2\pi)^n |\sigma^2 \mathbf{I}_n|}} \exp\left[-\frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\sigma^2 \mathbf{I}_n)^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right] \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left[-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2\right]. \end{aligned}$$

Because $\log(x)$ is a monotone increasing function, any maximum of $p_{\boldsymbol{\theta}}(\mathbf{y})$ is also a maximum of $\ell_{\boldsymbol{\theta}}(\mathbf{y}) := \log[p_{\boldsymbol{\theta}}(\mathbf{y})]$ and vice versa.

We can simplify as follows:

$$\begin{aligned} \ell_{\boldsymbol{\theta}}(\mathbf{y}) &= \log[p_{\boldsymbol{\theta}}(\mathbf{y})] \\ &= -\frac{1}{2}[n \log(2\pi\sigma^2) + \frac{1}{\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2]. \end{aligned}$$

So,

$$\frac{\partial}{\partial \boldsymbol{\beta}} \ell_{\boldsymbol{\theta}}(\mathbf{y}) = \frac{1}{\sigma^2} \mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

Under the full rank assumption we receive the same maximum likelihood estimator for $\boldsymbol{\beta}$ by setting the partials equal to 0 and solving for $\boldsymbol{\beta}$; that is, $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$.

We also have that

$$\frac{\partial}{\partial \sigma} \ell_{\boldsymbol{\theta}}(\mathbf{y}) = \frac{1}{\sigma^3} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 - \frac{n}{\sigma}.$$

Again setting the partial equal to 0 and solving for σ^2 while substituting our estimator $\hat{\boldsymbol{\beta}}$, we see that the maximum likelihood estimator is $\hat{\sigma}^2 = \frac{1}{n} \|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2 = \frac{1}{n} \|(\mathbf{I} - \mathbf{H})\mathbf{y}\|^2 = \frac{1}{n} \|\hat{\boldsymbol{\epsilon}}\|^2$. Thus, our estimator for $\boldsymbol{\beta}$ remains the same as the least squares estimator, while the estimator for ϵ differs from our unbiased estimator by a factor of $\frac{n}{n-p}$.

- (d) Define $\Theta := (\sigma^2, \beta_1, \dots, \beta_p) \in \mathbb{R}^{p+1}$. We also take $\Theta_1 := \Theta$ and $\Theta_0 := \{\theta = (\sigma^2, \beta_1, \dots, \beta_p) \in \mathbb{R}^{p+1} | \beta_S = 0\}$. Then we have that we are testing

$$H_0 : \boldsymbol{\theta} \in \Theta_0 \text{ versus } H_1 : \boldsymbol{\theta} \in \Theta_1.$$

Define $\ell_i(\mathbf{y}) := \max_{\boldsymbol{\theta} \in \Theta_i} p_{\boldsymbol{\theta}}(\mathbf{y})$, where $p_{\boldsymbol{\theta}}(\mathbf{y})$ is the density for the linear model; i.e.

$$p_{\boldsymbol{\theta}}(\mathbf{y}) := \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left[-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2\right].$$

We know from (c) that $(\hat{\sigma}_1^2, \hat{\beta}^1) := \arg \max_{\boldsymbol{\theta} \in \Theta} p_{\boldsymbol{\theta}}(\mathbf{y}) = (\frac{1}{n} \|\mathbf{y} - \mathbf{X}\hat{\beta}^1\|^2, (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y})$. Also,

$$\begin{aligned} \arg \max_{\boldsymbol{\theta} \in \Theta_0} p_{\boldsymbol{\theta}}(\mathbf{y}) &= \arg \max_{\boldsymbol{\theta} \in \Theta: \beta_S = 0} \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left[-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2\right] \\ &= \arg \max_{\boldsymbol{\theta} \in \Theta: \beta_S = 0} \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left[-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}_{-S}\boldsymbol{\beta}_{-S}\|^2\right]. \end{aligned}$$

But this is exactly the MLE for the model given by the partial model matrix \mathbf{X}_{-S} . So, $(\hat{\sigma}_0^2, \hat{\beta}^0) := \arg \max_{\boldsymbol{\theta} \in \Theta_0} p_{\boldsymbol{\theta}}(\mathbf{y}) = (\frac{1}{n} \|\mathbf{y} - \mathbf{X}_{-S}\hat{\beta}_{-S}\|^2, (\mathbf{X}_{-S}^T \mathbf{X}_{-S})^{-1} \mathbf{X}_{-S}^T \mathbf{y})$.

We then have

$$\begin{aligned} \Lambda &:= \frac{\ell_1(\mathbf{y})}{\ell_0(\mathbf{y})} \\ &= \frac{(2\pi\hat{\sigma}_1^2)^{-n/2} \exp\left[-\frac{1}{2\hat{\sigma}_1^2} \|\mathbf{y} - \mathbf{X}\hat{\beta}^1\|^2\right]}{(2\pi\hat{\sigma}_0^2)^{-n/2} \exp\left[-\frac{1}{2\hat{\sigma}_0^2} \|\mathbf{y} - \mathbf{X}\hat{\beta}^0\|^2\right]} \\ &= \left(\frac{\hat{\sigma}_0^2}{\hat{\sigma}_1^2}\right)^{n/2} \exp\left[\frac{1}{2} \left(\frac{1}{\hat{\sigma}_0^2} \|\mathbf{y} - \mathbf{X}\hat{\beta}^0\|^2 - \frac{1}{\hat{\sigma}_1^2} \|\mathbf{y} - \mathbf{X}\hat{\beta}^1\|^2\right)\right]. \end{aligned}$$

So,

$$\begin{aligned} 2 \log \Lambda &= n \log\left(\frac{\hat{\sigma}_0^2}{\hat{\sigma}_1^2}\right) + \left(\frac{1}{\hat{\sigma}_0^2} \|\mathbf{y} - \mathbf{X}\hat{\beta}^0\|^2 - \frac{1}{\hat{\sigma}_1^2} \|\mathbf{y} - \mathbf{X}\hat{\beta}^1\|^2\right) \\ &= n \log\left(\frac{\|\mathbf{y} - \mathbf{X}_{-S}\hat{\beta}_{-S}\|^2}{\|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2}\right) + n \left(\frac{\|\mathbf{y} - \mathbf{X}\hat{\beta}^0\|^2}{\|\mathbf{y} - \mathbf{X}_{-S}\hat{\beta}_{-S}\|^2} - \frac{\|\mathbf{y} - \mathbf{X}\hat{\beta}^1\|^2}{\|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2}\right) \\ &= n \log\left(\frac{\|\mathbf{y} - \mathbf{X}_{-S}\hat{\beta}_{-S}\|^2}{\|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2}\right) \\ &= n \log\left(\frac{\|(\mathbf{I} - \mathbf{H}_{-S})\mathbf{y}\|^2}{\|(\mathbf{I} - \mathbf{H})\mathbf{y}\|^2}\right), \end{aligned}$$

since $\|\mathbf{y} - \mathbf{X}\hat{\beta}^0\|^2 = \|\mathbf{y} - \mathbf{X}_{-S}\hat{\beta}_{-S}\|^2$ and $\|\mathbf{y} - \mathbf{X}\hat{\beta}^1\|^2 = \|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2$. Then

$$\Lambda^{2/n} = \frac{\|(\mathbf{I} - \mathbf{H}_{-S})\mathbf{y}\|^2}{\|(\mathbf{I} - \mathbf{H})\mathbf{y}\|^2},$$

so

$$\frac{n-p}{|S|} (\Lambda^{2/n} - 1) = \frac{\|(\mathbf{I} - \mathbf{H}_{-S})\mathbf{y}\|^2 - \|(\mathbf{I} - \mathbf{H})\mathbf{y}\|^2}{\|(\mathbf{I} - \mathbf{H})\mathbf{y}\|^2} = F.$$

Thus, we have that

Problem 2. Relationships among t -tests, F -tests, and R^2 .

Consider the linear regression model (1), such that $\mathbf{x}_{*,0} = \mathbf{1}_n$ is an intercept term (note that there are only $p - 1$ other predictors, for a total of p).

- Relate the R^2 of the linear regression to the F -statistic for a certain hypothesis test. What is the corresponding null hypothesis? What is the null distribution of the F -statistic? Are R^2 and F positively or negative related, and why does this make sense?
- Use the relationship found in part (a) to simulate the null distribution of the R^2 by repeatedly sampling from an F distribution (via `rf`). Fix $n = 100$ and try $p \in \{2, 25, 50, 75, 99\}$. Comment on these null distributions, how they change as a function of p , and why.
- Consider the null hypothesis $H_0 : \beta_j = 0$, which can be tested using either a t -test or an F -test. Write down the corresponding t and F statistics, and prove that the latter is the square of the former.
- Now suppose we are interested in testing the null hypothesis $H_0 : \beta_{-0} = \mathbf{0}$. One way of going about this is to start with the usual test statistic $t(\mathbf{c})$ for the null hypothesis $H_0 : \mathbf{c}^T \beta_{-0} = 0$, and then maximize over all $\mathbf{c} \in \mathbb{R}^{p-1}$:

$$t_{\max} \equiv \max_{\mathbf{c} \in \mathbb{R}^{p-1}} t(\mathbf{c}). \quad (8)$$

What is the null distribution of t_{\max}^2 ? What F -statistic is t_{\max}^2 equivalent to? How does the null distribution of t_{\max}^2 compare to that of $t(\mathbf{c})^2$?

Solution 2.

- Recall that

$$R^2 = \frac{\text{SSR}}{\text{SST}} = \frac{\|\mathbf{X}\hat{\beta} - \bar{y}\mathbf{1}_n\|^2}{\|\mathbf{y} - \bar{y}\mathbf{1}_n\|^2} = \frac{\|(\mathbf{H} - \mathbf{H}_1)\mathbf{y}\|^2}{\|(\mathbf{I} - \mathbf{H}_1)\mathbf{y}\|^2}.$$

Then we also have

$$1 - R^2 = \frac{\|(\mathbf{I} - \mathbf{H}_1)\mathbf{y}\|^2 - \|(\mathbf{H} - \mathbf{H}_1)\mathbf{y}\|^2}{\|(\mathbf{I} - \mathbf{H}_1)\mathbf{y}\|^2} = \frac{\|(\mathbf{I} - \mathbf{H})\mathbf{y}\|^2}{\|(\mathbf{I} - \mathbf{H}_1)\mathbf{y}\|^2}.$$

So, $\left(\frac{n-p}{p-1}\right) \frac{R^2}{1-R^2} = \frac{\|(\mathbf{H} - \mathbf{H}_1)\mathbf{y}\|^2/(p-1)}{\|(\mathbf{I} - \mathbf{H})\mathbf{y}\|^2/(n-p)} =: F$. But this is exactly the F -statistic for the hypothesis test with null hypothesis $H_0 : \beta_1 = \dots = \beta_{p-1} = 0$. We can also invert this relationship and find that for $c := \frac{p-1}{n-p} F$, $R^2 = \frac{c}{c+1}$.

Under the null distribution, the F -statistic is F distributed with $p - 1$ and $n - p$ degrees of freedom. Note that R^2 and F are positively related which makes sense - a higher R^2 indicates that the full model explains more of the variance in the observed data, which in turn suggests that we do not have $\beta_1 = \dots = \beta_{p-1} = 0$.

- # Simulation parameters*

```
n <- 100
p_list <- c(2, 25, 50, 75, 99)
```

```
# Run simulation for different values of p
for(p in p_list){
```

```

# Sample F-statistics
f <- rf(n = n, df1 = p-1, df2 = n-p)

# Compute R^2 from F-statistics
c <- (p-1)/(n-p)*F
R_2 <- c/(c+1)

# Plot?
}

```

- (c) Recall that $t = \frac{\hat{\beta}_j}{\hat{\sigma}/s_j}$, so $t^2 = \frac{\hat{\beta}_j^2 s_j^2}{\hat{\sigma}^2}$ and that $F = \frac{(\|\mathbf{X}\hat{\beta} - \mathbf{X}_{-j}\hat{\beta}_{-j}\|^2)/1}{(\|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2)/(n-p)}$. But, $\hat{\sigma}^2 = \frac{1}{n-p}\|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2$, so

$$t^2 = \frac{\hat{\beta}_j^2 s_j^2}{\|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2/(n-p)}$$

and thus it suffices to show that $\hat{\beta}_j^2 s_j^2 = \|\mathbf{X}\hat{\beta} - \mathbf{X}_{-j}\hat{\beta}_{-j}\|^2 = \|(\mathbf{H} - \mathbf{H}_{-j})\mathbf{y}\|^2 = \|\mathbf{H}_j\mathbf{y}\|^2 = \|\hat{\beta}_{j|-j}\mathbf{x}_{*j}^\perp\|^2$. Note that $s_j^2 = \|\mathbf{x}_{*j}^\perp\|^2$, so this equivalent to proving that $\|\hat{\beta}_{j|-j}\mathbf{x}_{*j}^\perp\|^2 = \hat{\beta}_j^2 s_j^2 = \hat{\beta}_j^2 \|\mathbf{x}_{*j}^\perp\|^2 = \|\hat{\beta}_j\mathbf{x}_{*j}^\perp\|^2$.

- (d)

Problem 3. Case study: Violent crime.

The `Statewide_crime.dat` file under `stat-961-fall-2021/data` contains information on the number of violent crimes and murders for each U.S. state in a given year, as well as three socioeconomic indicators: percent living in metropolitan areas, high school graduation rate, and poverty rate.

```
crime_data = read_tsv("../data/Statewide_crime.dat")
print(crime_data, n = 5)
```

```
## # A tibble: 51 x 6
##   STATE Violent Murder Metro HighSchool Poverty
##   <chr>   <dbl>   <dbl> <dbl>   <dbl>   <dbl>
## 1 AK      593      6  65.6    90.2     8
## 2 AL      430      7  55.4    82.4    13.7
## 3 AR      456      6  52.5    79.2    12.1
## 4 AZ      513      8  88.2    84.4    11.9
## 5 CA      579      7  94.4    81.3    10.5
## # ... with 46 more rows
```

The goal of this problem is to study the relationship between the three socioeconomic indicators and the per capita violent crime rate.

- These data contain the total number of violent crimes per state, but it is more meaningful to model violent crime rate per capita. To this end, go online to find a table of current populations for each state. Augment `crime_data` with a new variable called `Pop` with this population information (see `dplyr::left_join`) and create a new variable called `CrimeRate` defined as `CrimeRate = Violent/Pop` (see `dplyr::mutate`).
- Explore the variation and covariation among the variables `CrimeRate`, `Metro`, `HighSchool`, `Poverty` with the help of visualizations and summary statistics.
- Construct linear model based hypothesis tests and confidence intervals associated with the relationship between `CrimeRate` and the three socioeconomic variables, printing and/or plotting your results. Discuss the results in technical terms.
- Discuss your interpretation of the results from part (c) in language that a policymaker could comprehend, including any caveats or limitations of the analysis. Comment on what other data you might want to gather for a more sophisticated analysis of violent crime.

Solution 3.