

Homework 2

Anirban Chatterjee

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) The linear regression in (1) is equivalent to,

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i \quad (8)$$

where $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$, $\mathbf{x}_i, 1 \leq i \leq n$ are the rows of \mathbf{X} and $\boldsymbol{\epsilon} = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)^T$. Under the fixed design set up, (8) is further equivalent to,

$$y_i \sim N(\mathbf{x}_i^T \boldsymbol{\beta}, \sigma^2), \quad 1 \leq i \leq n \quad (9)$$

where $y_i, 1 \leq i \leq n$ are independent. Define

$$p_{\boldsymbol{\beta}, i} \equiv N(\mathbf{x}_i^T \boldsymbol{\beta}, \sigma^2), \quad 1 \leq i \leq n \quad (10)$$

Then combining (9) and (10) we can see that the linear regression model (1) does not fit into the classical inferential set up (2) since it does not have the i.i.d. assumption, instead each y_i are drawn independently from the distribution $p_{\boldsymbol{\beta}, i}$ which (possibly) differs based upon i . The closest form that we can get is,

$$y_i \stackrel{\text{ind}}{\sim} p_{\boldsymbol{\beta}, i}, \quad 1 \leq i \leq n$$

where **ind** implies $\{y_i : 1 \leq i \leq n\}$ are generated independently.

(b) The Fisher Information Matrix given in (4) is defined for a single random variable y coming from an distribution $p_{\boldsymbol{\theta}}$ index by a (multidimensional) parameter $\boldsymbol{\theta}$. The definition works well in the classical inferential set up because of the presence of i.i.d. sample, whereby the expectation does not differ from sample to sample. But in the regression set up of (1), this definition will not make sense because each sample comes from (possibly) different distributions. As a result we consider the following,

$$\tilde{\mathbf{I}}(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\boldsymbol{\theta}} \left[\frac{\partial^2}{\partial \boldsymbol{\beta}^2} p_{\boldsymbol{\beta}, i}(y_i) \right] \quad (11)$$

where $p_{\boldsymbol{\beta}, i}$ is defined in (10). It is easy to observe that under i.i.d. set up (11) reduces to the usual Fisher Information from (4). Recalling the definition from (10) it can be easily seen that,

$$\frac{\partial^2}{\partial \boldsymbol{\beta}^2} p_{\boldsymbol{\beta}, i}(y_i) = \frac{1}{\sigma^2} \mathbf{x}_{*i} \mathbf{x}_{*i}^T, \quad 1 \leq i \leq p$$

where \mathbf{x}_{*i} denotes the i^{th} column of \mathbf{X} for all $1 \leq i \leq p$. Therefore,

$$\tilde{\mathbf{I}}(\boldsymbol{\beta}) = \frac{1}{n\sigma^2} \mathbf{X}^T \mathbf{X}$$

Recall from Unit 1,

$$\hat{\boldsymbol{\beta}} \sim N\left(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}\right) \equiv N\left(\boldsymbol{\beta}, \frac{1}{n} \tilde{\mathbf{I}}^{-1}(\boldsymbol{\beta})\right)$$

Implying,

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \sim N(\mathbf{0}, \tilde{\mathbf{I}}^{-1}(\boldsymbol{\beta}))$$

which is in the same flavor as (3).

Although the above statement is similar to (3), it does not talk about asymptotics. Consider the following reformulation

$$\sqrt{n}\tilde{\mathbf{I}}(\boldsymbol{\beta})^{-1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \sim N(0, \mathbf{I}).$$

This statement is valid for all $n \in \mathbb{N}$, in particular for $n \rightarrow \infty$.

(c) Under the given model the likelihood is given by,

$$\begin{aligned} L(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right) \end{aligned}$$

where $\mathbf{x}_i, 1 \leq i \leq n$ are the rows of \mathbf{X} . Then the log-likelihood is given by,

$$l(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) = \log L(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

Taking derivatives with respect to $\boldsymbol{\beta}$ and σ^2 and equating to 0 we have,

$$\begin{aligned} \frac{\partial l}{\partial \boldsymbol{\beta}} &= \mathbf{X}^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = 0 \\ \frac{\partial l}{\partial \sigma^2} &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = 0 \end{aligned} \quad (12)$$

Solving the system of equations in (12) we find the solutions to be,

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad \hat{\sigma}^2 = \frac{1}{n}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}})^2$$

where $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$ and $\mathbf{x}_i, 1 \leq i \leq n$ are the rows of the matrix \mathbf{X} . Notice that $\hat{\boldsymbol{\beta}}$ equals the estimate derived in class, where as $\hat{\sigma}^2$ is smaller than the corresponding estimate derived in class, and as a result $\hat{\sigma}^2$ is a biased estimate of σ^2 .

(d) The given hypothesis can be reformulated in the form of (5) as follows,

$$\mathbf{H}_0 : (\boldsymbol{\beta}, \sigma^2) \in \Theta_0 \text{ vs } \mathbf{H}_1 : (\boldsymbol{\beta}, \sigma^2) \in \Theta_1$$

where $\Theta_0 = \{\boldsymbol{\beta} \in \mathbb{R}^p : \beta_S = 0\} \times \mathbb{R}^+$ and $\Theta_1 = \mathbb{R}^p \times \mathbb{R}^+$. Then the likelihood ratio test statistic is given by,

$$\begin{aligned} \Lambda &= \frac{\max_{(\boldsymbol{\beta}, \sigma^2) \in \Theta_1} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2}(y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2\right]}{\max_{(\boldsymbol{\beta}, \sigma^2) \in \Theta_0} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2}(y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2\right]} \\ &= \frac{\max_{(\boldsymbol{\beta}, \sigma^2) \in \Theta_1} \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left[-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right]}{\max_{(\boldsymbol{\beta}, \sigma^2) \in \Theta_0} \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left[-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right]} \end{aligned} \quad (13)$$

By part (c) we know that for the numerator the optimal values of β and σ^2 are given by,

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad \hat{\sigma}^2 = \frac{1}{n} (\mathbf{y} - \mathbf{X} \hat{\beta})^T (\mathbf{y} - \mathbf{X} \hat{\beta}) \quad (14)$$

Considering \mathbf{H} to be the projection matrix onto $C(\mathbf{X})$ we have,

$$\hat{\sigma}^2 = \frac{1}{n} \mathbf{y}^T (\mathbf{I} - \mathbf{H})^T (\mathbf{I} - \mathbf{H}) \mathbf{y} = \frac{1}{n} \mathbf{y}^T (\mathbf{I} - \mathbf{H}) \mathbf{y}$$

Using the optimal values found in (14) the numerator from (13) becomes,

$$\frac{1}{(2\pi\hat{\sigma}^2)^{n/2}} \exp(-n/2)$$

Now observe that the denominator is equivalent to,

$$\max_{(\beta_{-S}, \sigma^2) \in \mathbb{R}^{p-|S|} \times \mathbb{R}^+} \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left[-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}_{-S} \beta_{-S})^T (\mathbf{y} - \mathbf{X}_{-S} \beta_{-S}) \right] \quad (15)$$

where \mathbf{X}_{-S} is the matrix that we get by removing the columns S from the matrix \mathbf{X} and β_S is the vector that we get by removing the coordinates S from the vector β (This equivalency is because by definition of Θ_0 , $\beta_S = \mathbf{0}$ and thus $\mathbf{X}\beta = \mathbf{X}_{-S}\beta_{-S}$). Similar computations as in part (c) gives the following optimal values for (15),

$$\hat{\beta}_{-S} = (\mathbf{X}_{-S}^T \mathbf{X}_{-S})^{-1} \mathbf{X}_{-S}^T \mathbf{y}, \quad \hat{\sigma}_0^2 = \frac{1}{n} (\mathbf{y} - \mathbf{X}_{-S} \hat{\beta}_{-S})^T (\mathbf{y} - \mathbf{X}_{-S} \hat{\beta}_{-S}) \quad (16)$$

Using the equivalency of the denominator and (15) and the optimal values obtained in (16) the denominator from (13) becomes,

$$\frac{1}{(2\pi\hat{\sigma}_0^2)^{n/2}} \exp(-n/2)$$

Implying,

$$\Lambda = \left(\frac{\hat{\sigma}_0}{\hat{\sigma}} \right)^{n/2} = \left(\frac{\|\mathbf{y} - \mathbf{X}_{-S} \hat{\beta}_{-S}\|^2}{\|\mathbf{y} - \mathbf{X} \hat{\beta}\|^2} \right)^{n/2} \quad (17)$$

Recall the F -statistic is given by,

$$F = \frac{(\|\mathbf{y} - \mathbf{X}_{-S} \hat{\beta}_{-S}\|^2 - \|\mathbf{y} - \mathbf{X} \hat{\beta}\|^2) / |S|}{\|\mathbf{y} - \mathbf{X} \hat{\beta}\|^2 / (n - p)} \quad (18)$$

Combining (17) and (18) we conclude that,

$$\Lambda = \left(1 + \frac{|S|}{n - p} F \right)^{n/2}$$

It can be easily seen that Λ is an increasing function of F and thus the likelihood ratio test is equivalent to the F test.

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 $\mathbf{c} \in \mathbb{R}^{p-1}$:

$$t_{\max} \equiv \max_{\mathbf{c} \in \mathbb{R}^{p-1}} t(\mathbf{c}). \quad (19)$$

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.

- For the linear regression in (1) recall the R^2 is given by,

$$R^2 = \frac{\text{SSR}}{\text{SST}} = \frac{\|\mathbf{X}\hat{\beta} - \bar{y}\mathbf{1}_n\|^2}{\|\mathbf{y} - \bar{y}\mathbf{1}_n\|^2} \quad (20)$$

Consider testing for any significant coefficients except the intercept. Then the null hypothesis is given by $H_0 : \beta_1 = \dots = \beta_{p-1} = 0$ and the corresponding F -statistic is given by,

$$F \equiv \frac{(\|\mathbf{y} - \bar{y}\mathbf{1}_n\|^2 - \|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2) / (p - 1)}{\|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2 / (n - p)} \quad (21)$$

Recall the ANOVA decomposition of the variation in \mathbf{y} ,

$$\|\mathbf{y} - \bar{y}\mathbf{1}_n\|^2 = \|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2 + \|\mathbf{X}\hat{\beta} - \bar{y}\mathbf{1}_n\|^2 \quad (22)$$

Combining (20), (21) and (22) it is easy to see that,

$$R^2 = \frac{F}{F + \frac{n-p}{p-1}} \quad (23)$$

Under H_0 the F -statistic defined in (21) has a $F_{p-1, n-p}$ distribution. From the relation (23) it can be easily seen that F and R^2 are positively related (this follows since $x/(x+c)$ is an increasing function of x when $c > 0$).

This relation does make sense. Notice that large values of F -statistic gives enough evidence to reject the null hypothesis, thereby showing that other $p - 1$ covariates have significant contribution towards explaining the variation of y , which is exactly what large values of R^2 indicates.

- (b) Recall that under null $F \sim F_{p-1, n-p}$. For the given values of $p \in \{2, 25, 50, 75, 99\}$ the chunk below simulates the null distribution of R^2 using the relation (23).

```
# Setting seed for reproducibility
set.seed(961)

# Choosing n and number of covariates (p)
n <- 100; p <- c(2, 25, 50, 75, 99)

# Function for generating R^2 from F
gen.R2 <- function(x, n, p){return (x/(x+(n-p)/(p-1)))}

# Sampling R^2
R2.samp <- c()
for (i in 1:length(p)){
  F.samp <- rf(n, p[i]-1, n-p[i])
  R2.samp <- cbind(R2.samp, gen.R2(F.samp, n, p[i]))
}
R2.samp <- as.data.frame(R2.samp)
colnames(R2.samp) <- c("p = 2", "p = 25", "p = 50", "p = 75", "p = 99")

# Histograms of R^2 for each value of p
R2.hist <- reshape2::melt(as.data.frame(R2.samp), id.vars = NULL) %>%
  mutate(variable = factor(variable)) %>%
  ggplot(aes(x=value, fill=variable)) +
    geom_histogram(binwidth=0.03)+
    labs(x = expression(R^2))+
    facet_grid(variable~.)

# Saving the histogram
ggsave(R2.hist, file = paste0("R2hist.png"), width = 5, height = 5)
```

From Fig 1 it is clear that the distribution of R^2 shifts towards 1 as p increases to n . Notice that as p increases to n we are producing better and better fits of the data (Think about the extreme case $p = n$, assuming the covariates are linearly independent we can have exact fit for the data, and in this case we would have explained all the variation in the data using our model). As a result, by definition, increasing p would increase the value of R^2 .

- (c) The t -statistics for the test is given by,

$$t_j = \frac{\hat{\beta}_j}{\sqrt{\frac{1}{n-p} \|\hat{\epsilon}\|^2 / s_j^2}} \quad (24)$$

where $\hat{\beta}_j$ is the least square estimate of the coefficient β_j from the regression problem in (1), $s_j^2 = \left[\left(\mathbf{X}^T \mathbf{X} \right)_{jj}^{-1} \right]^{-1}$ and $\hat{\epsilon}$ is the residual vector. The F -statistics is given by,

$$F = \frac{\|\mathbf{X}\hat{\beta} - \mathbf{X}_{*, -j}\hat{\beta}_{-j}\|^2}{\|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2 / (n-p)} = \frac{\|\mathbf{X}\hat{\beta} - \mathbf{X}_{*, -j}\hat{\beta}_{-j}\|^2}{\frac{1}{n-p} \|\hat{\epsilon}\|^2} \quad (25)$$

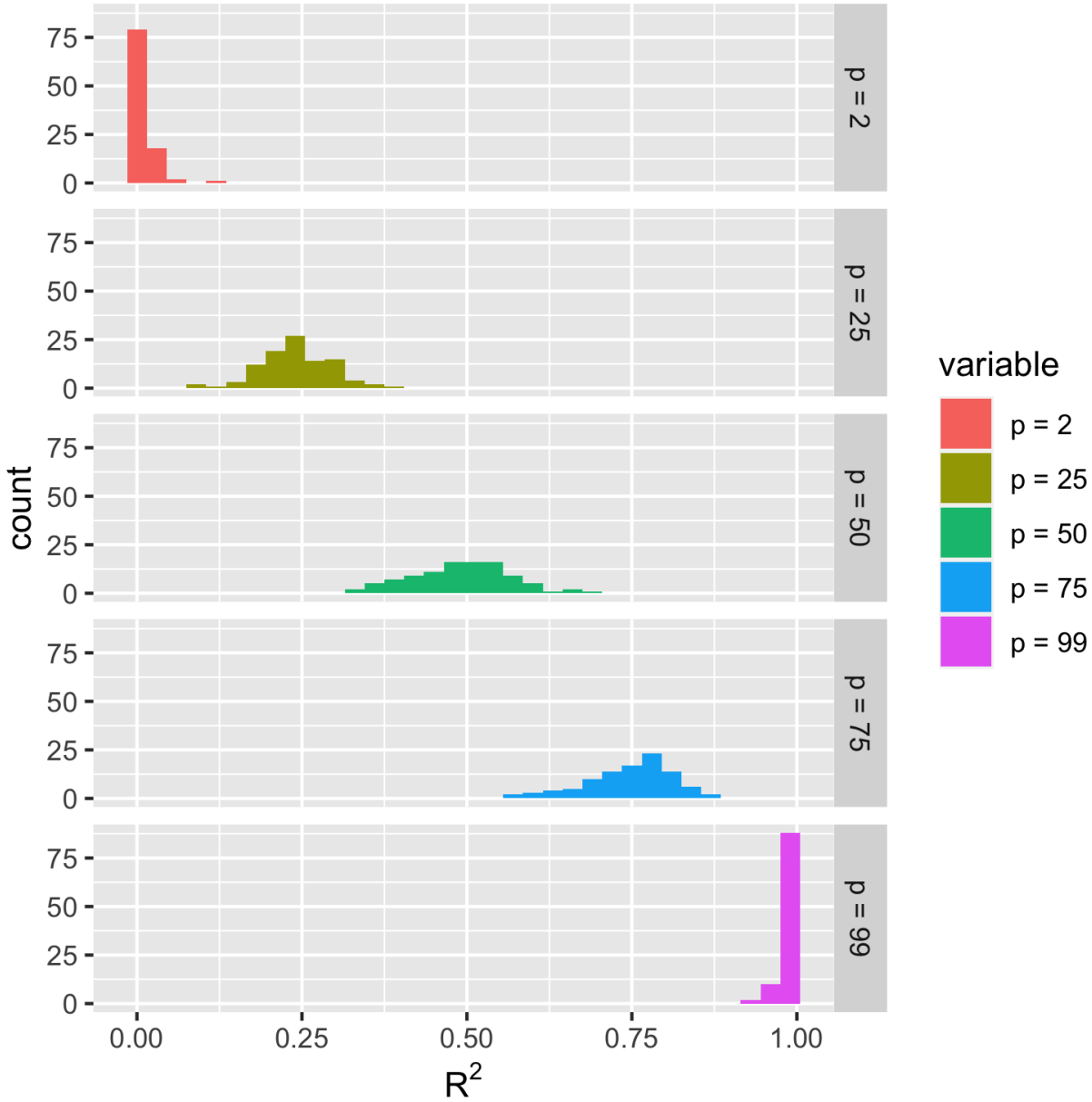


Figure 1: Histogram of R^2 for $p \in \{2, 25, 50, 75, 99\}$.

where $\hat{\beta}_{-j}$ is the least squares coefficients coming from the partial model,

$$\mathbf{y} = \mathbf{X}_{*, -j} \beta_{-j} + \epsilon$$

From Unit 1 recall that the regression problem (1) can be rewritten as,

$$\mathbf{y} = \mathbf{x}_{*j}^\perp \beta_j + \mathbf{X}_{*, -j} \beta'_{-j} + \epsilon$$

where \mathbf{x}_{*j}^\perp is the residual from regressing \mathbf{x}_{*j} on $\mathbf{X}_{*, -j}$ or in other words it is the projection of \mathbf{x}_{*j} onto $C(\mathbf{X}_{*, -j})^\perp$ (Note that the coefficient β_j remains unchanged in the two models). Using orthogonality, the least squares estimates of β_j and β'_{-j} can be found by regressing \mathbf{y}

on \mathbf{x}_{*j}^\perp and $\mathbf{X}_{*, -j}$ separately. Then by definition we must have,

$$\hat{\beta}_{-j} = \hat{\beta}'_{-j} \quad (26)$$

Consider the permutation matrix \mathbf{P} with rows $\mathbf{p}_i^T, 1 \leq i \leq p$ defined by,

$$\mathbf{p}_i = \begin{cases} \mathbf{e}_{i+1} & 1 \leq i \leq j-1 \\ \mathbf{e}_1 & i = j \\ \mathbf{e}_i & j+1 \leq i \leq p \end{cases}$$

where \mathbf{e}_k denotes the $p \times 1$ vector with 1 at position k and 0 at all other positions, for all $1 \leq k \leq p$. Define $\hat{\gamma}$ to be the least squares coefficient from regressing \mathbf{x}_{*j} on $\mathbf{X}_{*, -j}$. Then it is easy to observe that,

$$\mathbf{X} \left(\mathbf{P} \begin{bmatrix} 1 & \mathbf{0}^T \\ -\hat{\gamma} & \mathbf{I}_{p-1} \end{bmatrix} \right) = [\mathbf{x}_{*j}^\perp \quad \mathbf{X}_{*, -j}] =: \mathbf{X}'$$

Thus from Question 1 of Homework 1 we conclude that $C(\mathbf{X}) = C(\mathbf{X}')$. Thus considering the projection of \mathbf{y} on $C(\mathbf{X})$ and $C(\mathbf{X}')$ and using (26) we must have,

$$\mathbf{X}\hat{\beta} = \mathbf{x}_{*j}^\perp \hat{\beta}_j + \mathbf{X}_{*, -j} \hat{\beta}_{-j}$$

Then the F -statistics (from (25)) becomes,

$$F = \frac{\hat{\beta}_j^2 \|\mathbf{x}_{*j}^\perp\|^2}{\frac{1}{n-p} \|\hat{\epsilon}\|^2} \quad (27)$$

Recalling the variance of $\hat{\beta}_j$ (from unit 1), it can be seen that $s_j^2 = \|\mathbf{x}_{*j}^\perp\|^2$. Then from (24) and (27) we can conclude that,

$$t_j^2 = F$$

- (d) Observe that the hypothesis $\mathbf{H}_0 : \mathbf{c}^T \beta_{-0} = 0$ is equivalent to the hypothesis $\mathbf{H}_0 : \tilde{\mathbf{c}}^T \beta = 0$, where $\tilde{\mathbf{c}}^T = (0, \mathbf{c}^T)$. Then the usual test statistic $t(\mathbf{c})$ is given by,

$$t(\mathbf{c}) = \frac{\tilde{\mathbf{c}}^T \hat{\beta}}{\hat{\sigma} \sqrt{\tilde{\mathbf{c}}^T (\mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{c}}}} \quad (28)$$

Consider $\mathbf{e}_i, 1 \leq i \leq p$ to be the standard basis vectors of \mathbb{R}^p (also defined above in part (c)). Consider an $p \times p$ matrix $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_p]$ where

$$\mathbf{a}_j = \begin{cases} \mathbf{e}_1 & \text{if } j = 1 \\ \mathbf{e}_j - \bar{x}_{j-1} \mathbf{e}_1 & \text{if } 2 \leq j \leq p \end{cases}$$

and \bar{x}_j denotes the mean of column j of \mathbf{X} for all $1 \leq j \leq p-1$. Then

$$\mathbf{X}' = \mathbf{X}\mathbf{A} = [\mathbf{1}_n, \mathbf{x}_{*1} - \bar{x}_1 \mathbf{1}_n, \dots, \mathbf{x}_{*p-1} - \bar{x}_{p-1} \mathbf{1}_n]$$

where \mathbf{x}_{*j} denotes the j^{th} column of \mathbf{X} for all $1 \leq j \leq p-1$. It is easy to observe that

$$\mathbf{A}^{-1} = \begin{cases} \mathbf{e}_1 & \text{if } j = 1 \\ \mathbf{e}_j + \bar{x}_{j-1} \mathbf{e}_1 & \text{if } 2 \leq j \leq p \end{cases}$$

is the inverse of \mathbf{A} (The proof was done for $p = 3$ in HW 1, same proof would follow for the general case). Next consider the following linear model,

$$\mathbf{y} = \mathbf{X}'\boldsymbol{\beta}' + \boldsymbol{\epsilon} \quad (29)$$

By HW1 Problem 1(b) we know that $\hat{\boldsymbol{\beta}}' = \mathbf{A}^{-1}\hat{\boldsymbol{\beta}}$, implying,

$$\hat{\boldsymbol{\beta}}'_{-0} = \hat{\boldsymbol{\beta}}_{-0}$$

Recalling the expression (28) we have,

$$t(\mathbf{c}) = \frac{\tilde{\mathbf{c}}^T \hat{\boldsymbol{\beta}}'}{\hat{\sigma} \sqrt{\tilde{\mathbf{c}}^T (\mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{c}}}} \quad (30)$$

Using inversion of block matrices it can be easily seen that,

$$\tilde{\mathbf{c}}^T (\mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{c}} = \mathbf{c}^T (\mathbf{X}_{-0}^T \mathbf{X}_{-0} - n\bar{\mathbf{x}}\bar{\mathbf{x}}^T)^{-1} \mathbf{c} \quad (31)$$

where $\bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_{p-1})^T$. Define,

$$\mathbf{Z} = \mathbf{X}_{-0} - \mathbf{1}_n \bar{\mathbf{x}}^T$$

Observe that $C(\mathbf{Z})$ is orthogonal to $\mathbf{1}_n$. By definition $\mathbf{X}' = [\mathbf{1}_n, \mathbf{Z}]$ and $\mathbf{Z}^T \mathbf{Z} = \mathbf{X}_{-0}^T \mathbf{X}_{-0} - n\bar{\mathbf{x}}\bar{\mathbf{x}}^T$. Recall the model (29), then orthogonality of $\mathbf{1}_n$ and $C(\mathbf{Z})$ gives,

$$\hat{\boldsymbol{\beta}}'_{-0} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{y} \quad (32)$$

Combining (30), (31) and (32), $t(\mathbf{c})$ becomes,

$$t(\mathbf{c}) = \frac{\mathbf{c}^T (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{y}}{\hat{\sigma} \sqrt{\mathbf{c}^T (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{c}}}$$

Using Cauchy Schwarz Inequality we have,

$$|t(\mathbf{c})| = \frac{\left| \left((\mathbf{Z}^T \mathbf{Z})^{-1/2} \mathbf{c} \right)^T (\mathbf{Z}^T \mathbf{Z})^{-1/2} \mathbf{Z}^T \mathbf{y} \right|}{\hat{\sigma} \left\| (\mathbf{Z}^T \mathbf{Z})^{-1/2} \mathbf{c} \right\|_2} \leq \frac{1}{\hat{\sigma}} \left\| (\mathbf{Z}^T \mathbf{Z})^{-1/2} \mathbf{Z}^T \mathbf{y} \right\|_2$$

and equality would follow if and only if

$$t (\mathbf{Z}^T \mathbf{Z})^{-1/2} \mathbf{Z}^T \mathbf{y} = (\mathbf{Z}^T \mathbf{Z})^{-1/2} \mathbf{c} \text{ for some } t \in \mathbb{R} \setminus \{0\} \implies \mathbf{c} = t \mathbf{Z}^T \mathbf{y} \text{ for some } t \in \mathbb{R} \setminus \{0\}.$$

Then the maximum is attained and hence,

$$t_{\max}^2 = \frac{\left\| (\mathbf{Z}^T \mathbf{Z})^{-1/2} \mathbf{Z}^T \mathbf{y} \right\|^2}{\hat{\sigma}^2} = \frac{\mathbf{y}^T \mathbf{Z} (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{y}}{\hat{\sigma}^2} = \frac{\mathbf{y}^T \mathbf{P}_{\mathbf{Z}} \mathbf{y}}{\hat{\sigma}^2}$$

where \mathbf{P}_Z is the projection matrix onto $C(\mathbf{Z})$. Observe that under the null ($\mathbf{H}_0 : \beta_{-0} = \mathbf{0}$), $\mathbf{X}\beta = \mathbf{1}_n$. Then orthogonality of $\mathbf{1}_n$ and $C(\mathbf{Z})$ gives

$$\mathbf{P}_Z \mathbf{y} \sim N(0, \sigma^2 \mathbf{P}_Z).$$

Using Lemma 1.1 of Unit 2 along with $\mathbf{P}_Z^2 = \mathbf{P}_Z$ and $\text{trace}(\mathbf{P}_Z) = \dim(C(\mathbf{Z})) = p - 1$, we have

$$t_{\max}^2 = \frac{\mathbf{y}^T \mathbf{P}_Z \mathbf{y} / \sigma^2}{\hat{\sigma}^2 / \sigma^2} \sim (p - 1) F_{p-1, n-p} \text{ under null.}$$

Once again using the orthogonality,

$$\mathbf{P}_X = \mathbf{P}_Z + \mathbf{P}_{\mathbf{1}_n} \implies \mathbf{y}^T \mathbf{P}_Z \mathbf{y}^T = \|\mathbf{X}\hat{\beta} - \bar{y}\mathbf{1}_n\|^2$$

Thus t_{\max}^2 is equivalent to the F statistics used for testing $\mathbf{H}_0 : \beta_{-0} = \mathbf{0}$. Finally observe that under the null (any one of the two null hypothesis would suffice, since $\mathbf{H}_0 : \mathbf{c}^T \beta_{-0} = 0$ is an implication of $\mathbf{H}_0 : \beta_{-0} = \mathbf{0}$)

$$\tilde{\mathbf{c}}^T \hat{\beta} \sim N\left(0, \sigma^2 \tilde{\mathbf{c}}^T (\mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{c}}\right)$$

Then,

$$t(\mathbf{c}) \sim t_{n-p} \implies t(\mathbf{c})^2 \sim F_{1, n-p} \text{ under null.}$$

Considering the mean and variance of F -distribution it can be easily seen that the null distribution of t_{\max}^2 will be shifted towards the right from the null distribution of $t(\mathbf{c})^2$ and will have higher variance than the null distribution of $t(\mathbf{c})^2$. The comparison is shown in Fig 2.

```
n <- 1000; p <- 10

# Generating samples
x <- (p-1)*rf(n, p-1, n-p)
y <- rf(n, 1, n-p)

# X denotes t^2_max and Y denotes t(c)
xlab <- rep('X', n)
ylab <- rep('Y', n)

# Building a dataframe
df <- data.frame(value=c(x,y), lab=c(xlab,ylab), stringsAsFactors = F)

#Plot
dencomp <- ggplot(df, aes(x=value, fill=lab, color=lab, group=lab))+
  geom_histogram(aes(y = ..density..),
                 alpha = 0.4, position = position_dodge(),
                 binwidth = 0.5)+
  geom_line(aes(y = ..density..),
```

```
stat = 'density', show.legend = F)

# Saving the density comparison
ggsave(plot = dencomp, file = paste0("dencomp.png"), width = 5, height = 5)
```

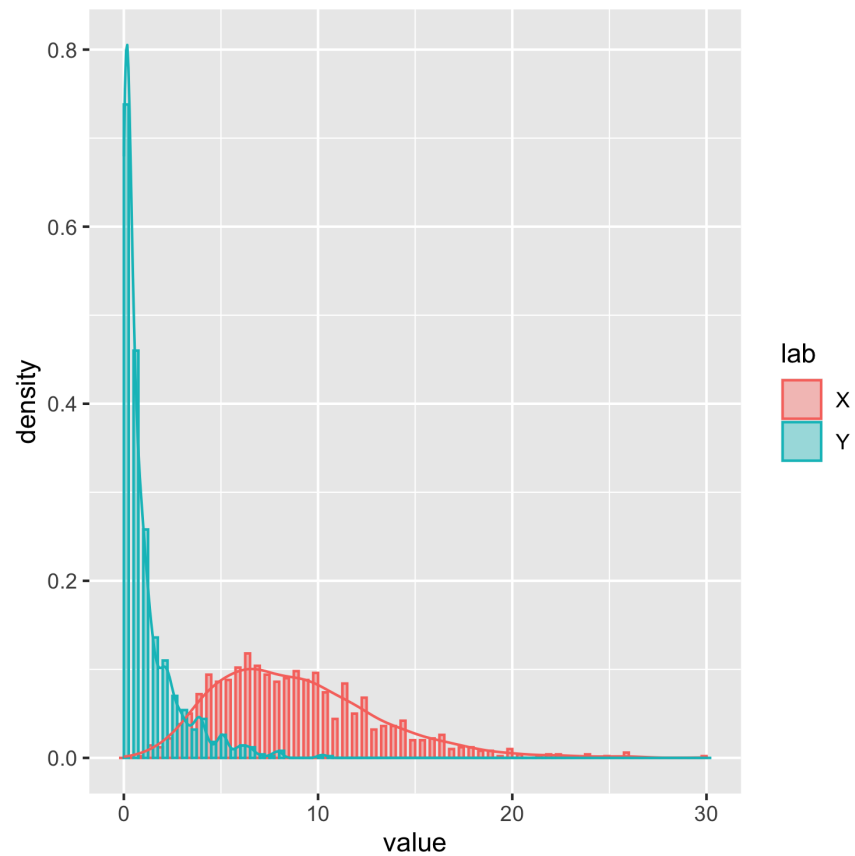


Figure 2: Comparison between $X \sim (p-1)F_{p-1, n-p}$ and $Y \sim F_{1, n-p}$.

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.

- I took the population data from [World Population Review](#). The below chunk creates the new variables `CrimeRate`.

```
# Reading population data
PopData <- read_csv("csvData.csv")[-31,2:3]

# Preprocessing population data

## Abbreviating state names
PopData$State <- state.abb[match(PopData$State,state.name)]

## Matching abbreviations to crime_data
```

```

PopData$State[49] <- "DC";PopData$State[32] <- "IO"

## Preparing for adding to crime_data
colnames(PopData) <- c("STATE","Pop")
PopData <- as.data.frame(PopData)

# Adding the population (pop) and crime rate (Violent/Pop) variables
crime_data <- crime_data %>%
  left_join(.,PopData) %>%
  mutate(CrimeRate = Violent/Pop) %>%
  print(.,n=5)

## Joining, by = "STATE"

## # A tibble: 51 x 8
##   STATE Violent Murder Metro HighSchool Poverty      Pop CrimeRate
##   <chr>   <dbl>   <dbl> <dbl>   <dbl>   <dbl>   <int>   <dbl>
## 1 AK      593      6  65.6     90.2     8    724357 0.000819
## 2 AL      430      7  55.4     82.4    13.7  4934193 0.0000871
## 3 AR      456      6  52.5     79.2    12.1  3033946 0.000150
## 4 AZ      513      8  88.2     84.4    11.9  7520103 0.0000682
## 5 CA      579      7  94.4     81.3    10.5  39613493 0.0000146
## # ... with 46 more rows

```

- (b) From this point onwards we will consider CrimeRate to be violent crimes per 10^4 person.

```

# Changing CrimeRate to per 10^4 person
crime_data <- crime_data %>%
  mutate(CrimeRate = CrimeRate*(10^4))

# Variation and Covariation Plots
cov.plot <- crime_data %>%
  GGally::ggpairs(
    columns = c(4,5,6,8),
    lower = list(continuous = 'smooth'))

## Registered S3 method overwritten by 'GGally':
##   method from
##   +.gg      ggplot2

ggsave(plot = cov.plot, filename = "covariance_plot.png",
        device = "png", width = 7, height = 6)

# Correlation plot of the variables
crime.slice <- crime_data[,c("Metro",
                             "HighSchool", "Poverty", "CrimeRate")]
corr.plot <- ggcorrplot::ggcorrplot(round(cor(crime.slice), 3))

ggsave(plot = corr.plot, filename = "corr_plot.png",

```

```

device = "png", width = 6, height = 6)

# Summary statistics of crime_data
crime_summary <- crime_data %>%
  vtable::st(out = "return",
    vars = c("Metro", "HighSchool", "Poverty", "CrimeRate"),
    digits = 2) %>%
  kableExtra::kable(format = "latex",
    row.names = NA, booktabs = T,
    digits = 2, align = rep("c", 4)) %>%
  kableExtra::save_kable("summary_stat.pdf")

```

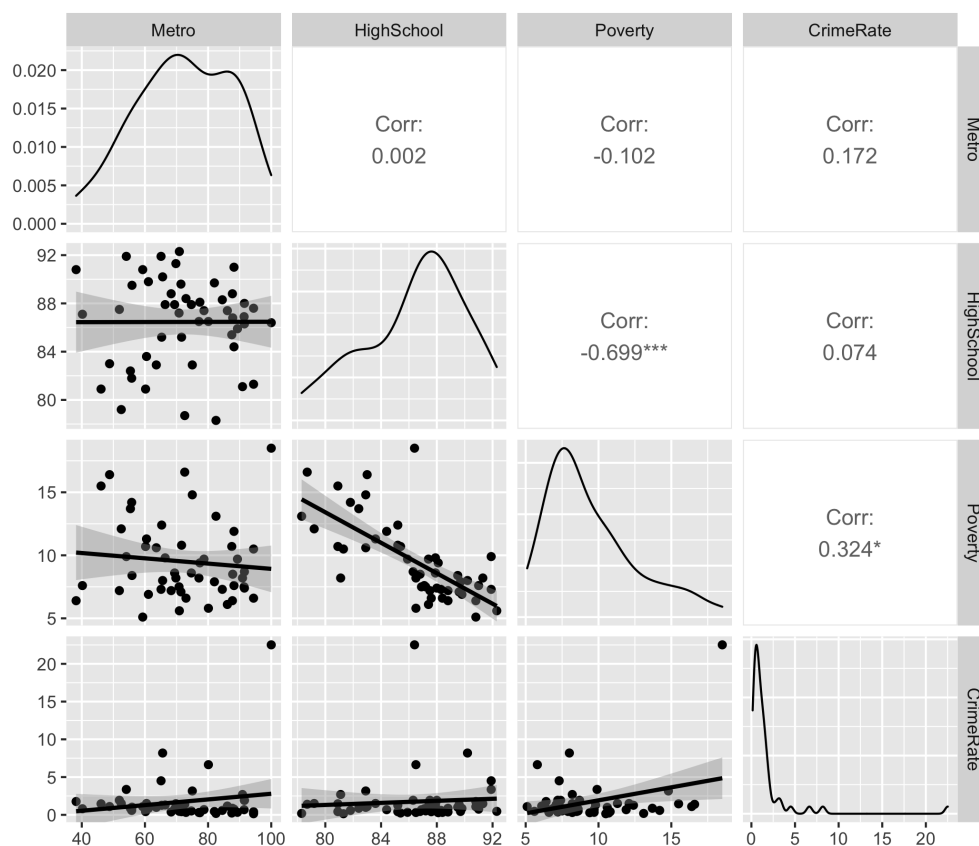


Figure 3: Scatterplots of all variables combinations among Metro, HighSchool, Poverty and CrimeRate with CrimeRate shown in scale of violent crimes per 10^4 people.

The variation and covariations can be observed from Fig 3 and 4 and the summary statistics are given in Table 5. From Fig 4 we can see that based on the data, poverty rate has a strong negative correlation with high school graduation rate, which in line with our intuition that in areas where high school graduation rates are high, poverty rate would come down. One surprising observation from Fig 3 and 4 is the positive correlation between all the socio-economic factors and crime rate. We will explore this further in part (c). From Table 5 we can observe that there is lot of variability in the percentage of population living in metropolitan areas. Also looking at the interquartile range and maximum values of CrimeRate, there

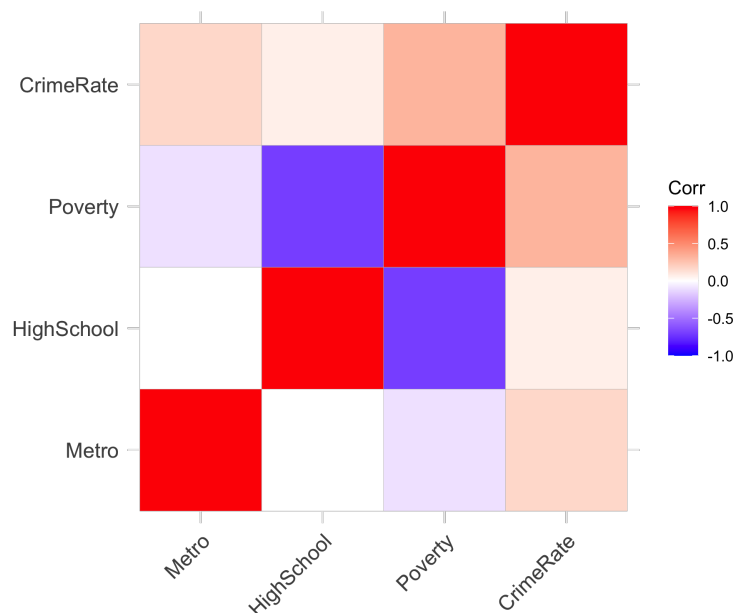


Figure 4: Correlation heatmap of **Metro**, **HighSchool**, **Poverty** and **CrimeRate** with **CrimeRate** shown in scale of violent crimes per 10^4 people.

Variable	N	Mean	Std. Dev.	Min	Pctl. 25	Pctl. 75	Max
Metro	51	72.25	15.28	38.2	60.8	86.8	100
HighSchool	51	86.46	3.62	78.3	84	88.8	92.3
Poverty	51	9.51	3.13	5.1	7.3	10.75	18.5
CrimeRate	51	1.75	3.33	0.15	0.45	1.49	22.52

Figure 5: Summary statistics of **crime_data** with **CrimeRate** shown in scale of violent crimes per 10^4 people.

appears to be presence of outlier in the data, which will be dealt with subsequently.

(c) We consider the following model,

$$\text{CrimeRate} = \beta_0 + \beta_1 \text{Metro} + \beta_2 \text{HighSchool} + \beta_3 \text{Poverty} + \epsilon \quad (33)$$

where $\epsilon \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}_{n \times n})$, variable names in the model indicates the n dimensional column vector of the same name from **crime_data** and **CrimeRate** has been modified to be in scale of violent crimes per 10^4 people. In order to ascertain that there is some socio-economic effect on the crime rate in the states we conduct the following hypothesis test,

$$\mathbf{H}_0 : \beta_{-0} = \mathbf{0} \text{ vs } \mathbf{H}_1 : \beta_{-0} \neq \mathbf{0} \quad (34)$$

In the following chunk we fit the linear model in (33).


```

# Fitting the linear model
crime.lm <- crime_data %>%
  lm(CrimeRate ~ Metro + HighSchool + Poverty,.)

# Summary of the fitted linear model
summary(crime.lm)

##
## Call:
## lm(formula = CrimeRate ~ Metro + HighSchool + Poverty, data = .)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.1346 -1.5118 -0.6766  0.6883 11.7901
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -59.52692    14.80527   -4.021 0.000209 ***
## Metro         0.05468     0.02607    2.097 0.041367 *
## HighSchool    0.57139     0.15318    3.730 0.000515 ***
## Poverty       0.83333     0.17793    4.683 2.43e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.788 on 47 degrees of freedom
## Multiple R-squared:  0.342, Adjusted R-squared:  0.3
## F-statistic: 8.144 on 3 and 47 DF,  p-value: 0.0001802

```

From the F-statistic given in the summary we can conclude that socio-economic effect on crime rate in states is significant at the 5% level. Based on the above conclusion we wish to look which of the socio-economic factors, namely **Metro**, **HighSchool** and **Poverty** are significant. Hence we consider the following hypothesis,

$$H_{0,i} : \beta_i = 0 \text{ vs } H_{1,i} : \beta_i \neq 0, \quad \forall 1 \leq i \leq 3$$

From the above summary we can see that the coefficients of all the socio-economic variables are significant at the 5% level.

Let us also discuss the implications of the coefficient estimates. Keeping all other covariates (socio-economic variables) fixed, an increase of 0.054 violent crimes per 10^4 person is associated with an unit increase in percentage of population living in metro cities, an increase of 0.57 violent crimes per 10^4 person is associated with an unit increase in high school graduation rate, and an increase of 0.83 violent crimes per 10^4 person is associated with an unit increase in poverty rate.

We also look at confidence interval of the coefficients of the socio-economic factors.

```

# Table of Confidence Intervals of Coeff estimates
confint(crime.lm) %>%
  kableExtra::kable(format = "latex",
    row.names = NA, booktabs = T,

```

```

      digits = 4, align = rep("c",4)) %>%
kableExtra::save_kable("confint.pdf")

```

	2.5 %	97.5 %
(Intercept)	-89.3113	-29.7426
Metro	0.0022	0.1071
HighSchool	0.2632	0.8796
Poverty	0.4754	1.1913

Figure 6: Confidence interval of the fitted coefficients in the linear model (33)

From Table 6 none of the confidence intervals contain 0. This shows that with probability .95 the true parameters are in the intervals given in Table 6. This means with high probability all the socio-economic factors are positively related to percentage of violent crimes. This seems counter-intuitive, as a result let us look at the regression diagnostic presented in Fig 7.

```

# Plotting residuals of the regression against fitted values
fitvsres <- crime.lm %>%
  ggplot(aes(.fitted,.resid)) +
  geom_point() +

  # Fit a locally weighted regression
  stat_smooth(method="loess") +

  # Showing the y-axis
  geom_hline(yintercept=0, col="red", linetype="dashed") +

  xlab("Fitted values") +
  ylab("Residuals") +
  ggtitle("Residual vs Fitted Plot") +
  theme_bw() +
  theme(plot.title = element_text(hjust = 0.5))

# Saving the residual vs fitted values
ggsave(plot = fitvsres, filename = "regdiag.png",
       device = "png", width = 4, height = 4)

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

```

It is clear from Fig 7 that there is presence of outlier, which can heavily impact the conclusions drawn above. The outlier state is given in Table 8.

```

# Index of the outlier state
out.state <- which.max(crime.lm$residuals)

```

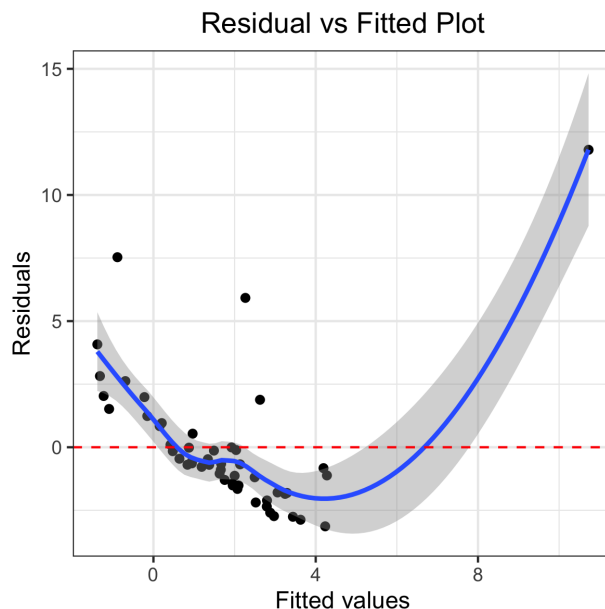


Figure 7: Plot of residuals against fitted values for the linear model (33)

```
# Details of the outlier state in a table
crime_data[out.state,] %>%
  kableExtra::kable(format = "latex",
                    row.names = NA, booktabs = T,
                    digits = 4, align = rep("c",4)) %>%
  kableExtra::save_kable("outlier.pdf")
```

STATE	Violent	Murder	Metro	HighSchool	Poverty	Pop	CrimeRate
DC	1608	44	100	86.4	18.5	714153	22.5162

Figure 8: The outlier state DC

Looking at the `crime_data` and comparing DC (from Table 8) with other states we can see that although it has high percentages of metropolitan population, high rates of high school graduation, there is significantly high percentage of violent crimes happening there, which does not agree with the other states. Let us redo the above analysis with DC removed from the dataset.

```
# Running the regression with DC removed
crime_lm.new <- crime_data %>%

  # Removing the outlier state
  slice(-out.state) %>%

  lm(CrimeRate ~ Metro + HighSchool + Poverty, .)
summary(crime_lm.new)
```

```
##
## Call:
## lm(formula = CrimeRate ~ Metro + HighSchool + Poverty, data = .)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.2947 -0.7477 -0.4244  0.1162  6.5108
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -3.303786   9.858384  -0.335   0.739
## Metro       -0.014132   0.015933  -0.887   0.380
## HighSchool   0.066073   0.098105   0.673   0.504
## Poverty     -0.006652   0.127960  -0.052   0.959
##
## Residual standard error: 1.548 on 46 degrees of freedom
## Multiple R-squared:  0.04577, Adjusted R-squared:  -0.01647
## F-statistic: 0.7354 on 3 and 46 DF,  p-value: 0.5363
```

From the **F-statistic** given in the summary we can see that there is not enough evidence to reject the hypothesis (34) and thus we can conclude that based on the available data, none of the socio-economic factors have significant effect on the **CrimeRate** in a state.

- (d) Based on the data collected the state of DC appears to be different from the others. It has significantly higher crimes per capita, when considered with poverty rate, metropolitan population and high school graduation rates. The crime rates in the state of DC needs to be analysed further and it would be unwise to implement matching (with other states) policies based upon the data collected. In the present data, it can be seen that DC is not in line with other states, and in order to make useful inferences from the data, it is best to remove DC from the dataset.

Based on the provided data, the socio-economic factors of poverty rates, metropolitan population and high school graduation rates cannot explain the variability of crime rates across states. In simpler terms, the data doesn't provide any evidence to suggest that above three socio-economic factors have any significant effect on the crime rates across the states.

But one must take this conclusion with a grain of salt. Although there was no evidence of significant effect of the factors, one cannot dismiss that they do not play any role. There may be other deeper factors at play here which were absent in the data. For example, the existing policies in metropolitan areas across states may be good enough to keep crimes at bay, which could overshadow the effects of the factors considered here.