

# Homework 6

jiaxi li

## Table of contents

Question 1 . . . . .	2
----------------------	---

### ! Important

Please read the instructions carefully before submitting your assignment.

1. This assignment requires you to only upload a PDF file on Canvas
2. Don't collapse any code cells before submitting.
3. Remember to make sure all your code output is rendered properly before uploading your submission.

Please add your name to the author information in the frontmatter before submitting your assignment

In this assignment, we will perform various tasks involving principal component analysis (PCA), principal component regression, and dimensionality reduction.


We will need the following packages:

```
packages <- c(
  "tibble",
  "dplyr",
  "readr",
  "tidyr",
  "purrr",
  "broom",
  "magrittr",
  "corrplot",
  "car"
)
```

```
# renv::install(packages)
sapply(packages, require, character.only=T)
```

---

## Question 1

 70 points

Principal component analysis and variable selection

### 1.1 (5 points)

The `data` folder contains a `spending.csv` dataset which is an illustrative sample of monthly spending data for a group of 5000 people across a variety of categories. The response variable, `income`, is their monthly income, and objective is to predict the `income` for an individual based on their spending patterns.

Read the data file as a tibble in R. Preprocess the data such that:

1. the variables are of the right data type, e.g., categorical variables are encoded as factors
2. all column names to lower case for consistency
3. Any observations with missing values are dropped

```
path <- "data/spending.csv"
```

```
df <- ... # Insert your code here
```

---

### 1.2 (5 points)

Visualize the correlation between the variables using the `corrplot()` function. What do you observe? What does this mean for the model?

```
df_x %>% ... # Insert your code here
```

---

1.3 (5 points)

Run a linear regression model to predict the `income` variable using the remaining predictors. Interpret the coefficients and summarize your results.

```
... # Insert your code here
```

---

1.3 (5 points)

Diagnose the model using the `vif()` function. What do you observe? What does this mean for the model?

```
... # Insert your code here
```

---

1.4 (5 points)

Perform PCA using the `princomp` function in R. Print the summary of the PCA object.

```
pca <- ... # Insert your code here  
... # Insert your code here
```

---

1.5 (5 points)

Make a screeplot of the proportion of variance explained by each principal component. How many principal components would you choose to keep? Why?

```
... # Insert your code here
```

1.6 (5 points)

By setting any factor loadings below 0.2 to 0, summarize the factor loadings for the principal components that you chose to keep.

```
clean_loadings <- ... # Insert your code here
```

Visualize the factor loadings.

```
... # Insert your code here
```

---

1.7 (15 points)

Based on the factor loadings, what do you think the principal components represent?

Provide an interpretation for each principal component you chose to keep.

---

1.8 (10 points)

Create a new data frame with the original response variable **income** and the principal components you chose to keep. Call this data frame **df\_pca**.

```
... # Insert your code here
```

Fit a regression model to predict the **income** variable using the principal components you chose to keep. Interpret the coefficients and summarize your results.

```
... # Insert your code here
```

Compare the results of the regression model in 1.3 and 1.9. What do you observe? What does this mean for the model?

```
... # Insert your code here
```

---

1.10 (10 points)

Based on your interpretation of the principal components from Question 1.7, provide an interpretation of the regression model in Question 1.9.

---

---

## Session Information

Print your R session information using the following command

```
sessionInfo()
```

R version 4.3.2 (2023-10-31 ucrt)

Platform: x86\_64-w64-mingw32/x64 (64-bit)

Running under: Windows 11 x64 (build 22631)

Matrix products: default

locale:

```
[1] LC_COLLATE=English_United States.utf8
[2] LC_CTYPE=English_United States.utf8
[3] LC_MONETARY=English_United States.utf8
[4] LC_NUMERIC=C
[5] LC_TIME=English_United States.utf8
```

time zone: America/New\_York

tzcode source: internal

attached base packages:

```
[1] stats      graphics  grDevices  utils      datasets  methods    base
```

other attached packages:

```
[1] car_3.1-2      carData_3.0-5  corrplot_0.92  magrittr_2.0.3 broom_1.0.5
[6] purrr_1.0.2    tidyr_1.3.1    readr_2.1.5    dplyr_1.1.4    tibble_3.2.1
```

loaded via a namespace (and not attached):

```
[1] vctrs_0.6.5      cli_3.6.2       knitr_1.45      rlang_1.1.3
[5] xfun_0.41        generics_0.1.3  jsonlite_1.8.8  glue_1.7.0
[9] backports_1.4.1  htmltools_0.5.7 hms_1.1.3       fansi_1.0.6
[13] rmarkdown_2.25   abind_1.4-5     evaluate_0.23   tzdb_0.4.0
[17] fastmap_1.1.1    yaml_2.3.8      lifecycle_1.0.4 compiler_4.3.2
[21] codetools_0.2-19 pkgconfig_2.0.3 rstudioapi_0.15.0 digest_0.6.34
[25] R6_2.5.1         tidyselect_1.2.0 utf8_1.2.4      pillar_1.9.0
[29] tools_4.3.2
```