

```

---
title: "LAB 9"
author: "STAT 131a"
date: "November 6, 2019"
output: pdf_document
---

```{r setup, include=FALSE}
knitr::opts_chunk$set(echo = TRUE)
```

```

Welcome to the lab 9! In this lab, we will learn how to perform PCA and implement multiple linear regressions in R.

### ## A simple PCA example

This is a simple example for PCA from [R Bloggers](https://www.r-bloggers.com/computing-and-visualizing-pca-in-r/). The iris dataset is perhaps the best known dataset for classification. The data set contains 3 classes of 50 instances each, where each class refers to a type of iris plant.

```

```{r}
# Load data
data(iris)
head(iris, 3)
```

```

There are five variables in the dataset `iris`, where `Sepal.Length`, `Sepal.Width`, `Petal.Length` and `Petal.Width` are continuous and `Species` is categorical.

```

```{r}
names(iris)
```

```

Apply PCA to the four continuous variables

```

```{r}
# log transform on the four continuous variable
log.ir <- log(iris[, 1:4])
# the Species variable
ir.species <- iris[, 5]

```

# apply PCA - scale. = TRUE is highly advisable, but default is FALSE.  
# center and scale is used to standardize the variables prior to the application of PCA.

```

ir.pca <- prcomp(log.ir,
                  center = TRUE,
                  scale. = TRUE)
```

```

Plot the PC1 and PC2.

```

```{r}
colorvec <- c("red", "green", "blue")
names(colorvec) = unique(ir.species)

```

```
plot(ir.pca$x[, 1], ir.pca$x[, 2], col = unname(colorvec[ir.species]), xlab =
"PC1", ylab = "PC2")
legend("topright", legend = names(colorvec), fill = unname(colorvec))
```

```

Lets add a few NA values to the data.

```
```{r}
iris[2,3]= NA
iris[5,2]= NA
iris[10,1] = NA
```

```

Apply PCA as before, but make change to code to address the missing values.

```
```{r}
log.ir <- attributes(na.omit(log(iris[, 1:4])))[["na.action"]]

ir.species <- iris[, 5]
ir.pca <- prcomp(log.ir,
                 center = TRUE,
                 scale. = TRUE)

...

# Multiple Regression
# Diamond Price

```

This is a very large data set showing various factors of over 50,000 diamonds including price, cut, color, clarity, etc. We are interested in the prediction of the `price` and how different factors influence the diamond price.

```
- **price**: price in US dollars ($326-$18,823)
- **carat**: weight of the diamond (0.2-5.01)
- **cut**: quality of the cut (Fair, Good, Very Good, Premium, Ideal)
- **color**: diamond colour, from J (worst) to D (best)
- **clarity**: a measurement of how clear the diamond is (I1 (worst), SI1, SI2, VS1, VS2, VVS1, VVS2, IF (best))
- **length.in.mm**: length in mm (0-10.74)
- **width.of.mm**: width in mm (0-58.9)
- **depth.in.mm**: depth in mm (0-31.8)
- **depth**: total depth percentage = z / mean(x, y) = 2 * z / (x + y) (43-79)
- **table**: width of top of diamond relative to widest point (43-95)

```

```
```{r}
diamonds <- read.csv("diamonds.csv")
head(diamonds)
```

```

### **\*\*Exercise 1: Exploratory Data Analysis before Regression\*\***

To understand the relationship between prices and carat, cut, etc. We first do the scatter plot. Create scatter plots between the response variable (price) and all the continuous variables. The function `is.numeric` might be helpful to check whether a variable is numeric. For example,

```
```{r}

```

```

vec1 = 1:10
vec2 = as.character(1:10)
vec1
vec2
# The following line of code would return TRUE
is.numeric(vec1)
# The following line of code would return FALSE
is.numeric(vec2)
# The following line of code would return TRUE
is.character(vec2)
```

```

The function `which` can help you to locate the column indices of the numeric vectors. (In fact, function `which` is a super useful function in R.)

```

```{r}
# function `which` give the TRUE indices of a logical object
which(c(TRUE, FALSE, TRUE, FALSE, TRUE))
```

```

```

```{r}
# Insert your code here, use for loop to loop through each numerical variable
x=apply(diamonds, is.numeric)
y=which(x==TRUE)
for (i in y) {
  plot(diamonds[,i], diamonds$price, ylab="Price", main=paste("price v" ,
colnames(diamonds)[i]), xlab=colnames(diamonds)[i])
}
```

```

## Multiple Regression with Continuous Variable

**\*\*Exercise 2 Fit the model and Calculate the Statistics\*\***

(a) Fit a linear model to price with all the continuous variable as explanatory variables. Print the summary of your model.

```

```{r}
# Insert you code here, save your model as `fit`
fit <- (lm(price ~ carat + depth + table + length.in.mm + width.of.mm+ depth.in.mm,
data = diamonds))
summary(fit)
```

```

(b) Calculate the fitted values.

```

```{r}
# Insert you code here, save your results as `fitted.value`

fitted.value <- fitted.values(fit)
```

```

(c) Calculate the residual, the residual sum of squares (RSS) and the total sum of squares (TSS) using the `fitted.value` from the above chunk.

```

```{r}

```

```
# Insert you code here, save your results as `RSS`
residual <- residuals(fit)
RSS <- sum(residual^2)
RSS
TSS <- sum(((diamonds$price) - mean(diamonds$price))^2)
TSS
...
```

(d) Calculate the R-square ( $R^2$ ) using RSS and TSS. How will you interpret the  $R^2$ ?

```
```{r}
# Insert you code here, save your results as `Rsqr`
Rsqr <- 1-(RSS/TSS)
Rsqr
...
```

>  $R^2$  is the coefficient of determination and it tells us about the goodness of fit.  $R^2$  puts into quantification how our fitted values actually are to observed values. Our  $R^2$  is about 85% so the regression is closely fitted with the data.

**\*\*Exercise 3 Think deeper. Is the model resonable?\***

(a) Using your fitted model, we can write down the model formula: (the estimated coefficients can be found in the summary chart)

```
$$earnings = 20849.316 + 10686.309 \cdot carat - 203.154 \cdot depth - 102.446 \cdot
table - 1315.668 \cdot length.in.mm + 66.322 \cdot width.of.mm + 41.628 \cdot
depth.in.mm + e$$
```

How do we interpret this equation? By looking at the  $p$ -value, we know that the coefficients we estimated are significant except for `depth.in.mm`. Take `length.in.mm` for example, the coefficient - 1315.668 tells us that if we increase `length.in.mm` by 1, the price of the diamond will decrease 1315.668 dollars in average. How weird is that! Will you consider drop `length.in.mm`, `width.of.mm` and `depth.in.mm` in your model? Why? (HINT: Check the correlation between `length.in.mm`, `width.of.mm`, `depth.in.mm` and `carat`.)

```
```{r}
# insert your code here
cor(diamonds[, c("length.in.mm", "width.of.mm", "depth.in.mm" ,"carat")])
...
```

> We would remove the variables `length.in.mm` and `width.of.mm` since the variables are highly correlated, we do not need all of them in our models. The values seem to be measuring the same thing. Furthermore, if `length.in.mm` and `width.of.mm` are outrageous values we don't need them to make an impact on our data if they aren't important to begin with. They can both be represented by `depth.in.mm`.

(b) Plot the variables versus the residual and calculate their correlation? Can you find any problem in your model by looking at the scatter plots? If you're asked to add some terms to improve the model, what will you do? (HINT: Consider the scatter plot in Exercise 1: are the relationships linear?)

```
```{r}
# Insert your code here, use for loop to loop through each numerical variable
x=sapply(diamonds, is.numeric)
```

```

y=which(x==TRUE)
for (i in y) {
  plot(diamonds[,i], residual,main=paste("residual v" , colnames(diamonds)[i]),
xlab=colnames(diamonds)[i])
  print (cor(diamonds[,i], residual))
}

```

```

}

```

> There appears to be problems in the model since the models do not appear to have a randomness. The residuals represent what is left in y after all the linear effects of the explanatory variables are removed. The relationships aren't linear since the patterns are not random. Similar to lecture, the correlations are actually uncorrelated, residuals always have a mean of 0. We expect  $R^2$  to increase if we add more variables since there is a direct relationship to RSS which also increases with more variables. When adding new terms, I would attempt to improve the fit since the  $R^2$  always increases then any terms will improve the regression. Examples of good variable terms to add could be "brand" or "setting."

## Multiple Regression with Continuous and Categorical Variable

\*\*Exercise 4\*\*

(a) Fit a linear regression model with explanatory variable `carat`, `depth`, `table`, `clarity`, `color` and `cut`.

```

```{r}
levels(diamonds$clarity)
levels(diamonds$color)
levels(diamonds$cut)
```

```

```

```{r}
# Insert you code here, save your model as `fit.categorical`
fit.categorical <- (lm(price ~ carat+ depth + table + as.factor(clarity) +
as.factor(color) + as.factor(cut), data = diamonds))
summary(fit.categorical)
```

```

(b) Write the equation when

I. clarity is VS2, color is H and cut is Premium.

```

```{r}

v1=subset(diamonds,clarity=="VS2" & color=="H" & cut=="Premium")
v1=(lm(price ~ carat + depth + table, data = v1))
summary(v1)
v2=subset(diamonds,clarity=="I1" & color=="D" & cut=="Fair")
v2=(lm(price ~ carat + depth + table, data = v2))
summary(v2)
```

```

```

```{}

```

```
# repalce ??? by numerical values
price = -18171.89 + 8699.36 * carat + 168.42 * depth + 79.18 * table
...
```

II. clarity is I1, color is D and cut is Fair.

```
```{
# repalce ??? by numerical values
price = -60252.5 + 5553.0 * carat + 734.6 * depth + 158.9 * table
```
```

This study resource was  
shared via CourseHero.com