

Lab 12

STATC 131A

April 18, 2022

Welcome to the lab 12! In this lab, we will use linear regression to predict the red wine quality using physicochemical tests scores such as citric acid, pH, etc.

But first, a review of using the `predict()` function.

```
x1 = rnorm(100)
x2 = rnorm(100)
y = 2*x1 + x2 + rnorm(100)
lm_out = lm(y~x1 + x2)
summary(lm_out)
```

```
##
## Call:
## lm(formula = y ~ x1 + x2)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.17475 -0.64389 -0.00244  0.48895  2.81397
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.19139    0.09722  -1.969   0.0519 .
## x1           1.88905    0.08837  21.378 <2e-16 ***
## x2           0.88714    0.08785  10.098 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.9718 on 97 degrees of freedom
## Multiple R-squared:  0.8549, Adjusted R-squared:  0.8519
## F-statistic: 285.6 on 2 and 97 DF,  p-value: < 2.2e-16
```

Calculate prediction for y when x1 is 1 and x2 is 0.5.

```
lm_out$coefficients[1] + lm_out$coefficients[2]* 1 + lm_out$coefficients[3]* 0.5
```

```
## (Intercept)
##      2.14123
```

Another way to do this.

```
predict(lm_out, newdata = data.frame(x1= 1, x2= 0.5))
```

```
##          1
## 2.14123
```

Can do several at once.

```
predict(lm_out, newdata = data.frame(x1= c(1, 2), x2= c(0.5, -1) ))
```

```
##          1          2
## 2.141230 2.699563
```

We can find intervals for confidence of average or prediction interval for an individual outcome.

```
predict(lm_out, newdata = data.frame(x1= c(1, 2), x2= c(0.5, -1) ), interval = "confidence")
```

```
##          fit          lwr          upr
## 1 2.141230 1.870022 2.412438
## 2 2.699563 2.263972 3.135153
```

```
predict(lm_out, newdata = data.frame(x1= c(1, 2), x2= c(0.5, -1) ), interval = "prediction")
```

```
##          fit          lwr          upr
## 1 2.141230 0.1934263 4.089034
## 2 2.699563 0.7221594 4.676966
```

A few other notes on regression.

- 1) If you try to predict one variable and include a perfectly correlated variable in the prediction set, then that variable will be perfectly fit to the outcome to the exclusion of all others.

```
perf_cor = y/4
summary(lm( y ~ perf_cor + x1 + x2 ))
```

```
## Warning in summary.lm(lm(y ~ perf_cor + x1 + x2)): essentially perfect fit:
## summary may be unreliable
```

```
##
## Call:
## lm(formula = y ~ perf_cor + x1 + x2)
##
## Residuals:
##          Min          1Q          Median          3Q          Max
## -3.907e-16 -1.349e-16 -3.350e-17  4.690e-17  4.157e-15
##
## Coefficients:
##              Estimate Std. Error   t value Pr(>|t|)
## (Intercept)  4.441e-17  4.592e-17  9.670e-01   0.336
## perf_cor      4.000e+00  1.881e-16  2.127e+16  <2e-16 ***
```

```
## x1          9.181e-18  9.780e-17  9.400e-02    0.925
## x2         -7.634e-17  5.827e-17 -1.310e+00    0.193
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.501e-16 on 96 degrees of freedom
## Multiple R-squared:      1, Adjusted R-squared:      1
## F-statistic: 1.039e+33 on 3 and 96 DF,  p-value: < 2.2e-16
```

- 2) If there are more variables used for prediction than there are observations, `lm` will only keep the first $n-1$ variables.

```
x3= rnorm(100)
x4= rnorm(100)
x5= rnorm(100)

all_x = data.frame(x1,x2,x3,x4,x5, y) # 100 x 6 df
lm(y~ . ,data= all_x[1:4,]) # only use the first 4 observations (4<5)

##
## Call:
## lm(formula = y ~ ., data = all_x[1:4, ])
##
## Coefficients:
## (Intercept)          x1          x2          x3          x4          x5
##      0.7634      2.1426      1.1664     -0.1683         NA         NA

# Then lm only use the first 4-1=3 variables
```

Wine data

The wine dataset is related to red variants of the Portuguese “Vinho Verde” wine. There are 1599 samples available in the dataset. Due to privacy and logistic issues, only physicochemical (inputs) and sensory (the output) variables are available (e.g. there is no data about grape types, wine brand, wine selling price, etc.).

The explanatory variables are all continuous variables based on physicochemical tests:

- fixed acidity
- volatile acidity
- citric acid
- residual sugar
- chlorides
- free sulfur dioxide
- total sulfur dioxide
- density
- pH
- sulphates
- alcohol

The response variable is the **quality** score between 0 and 10 (based on sensory data).

Read data. We randomly split the data into two parts-the `wine` dataset with 1199 samples and the `wine.test` dataset with 400 samples. Splitting the dataset is a common technique when we want to evaluate the model

performance. There are training set, validation set, and test set. The validation set is used for model selection. That is, to estimate the performance of the different model in order to choose the best one. The test set is used for estimating the performance of our final model.

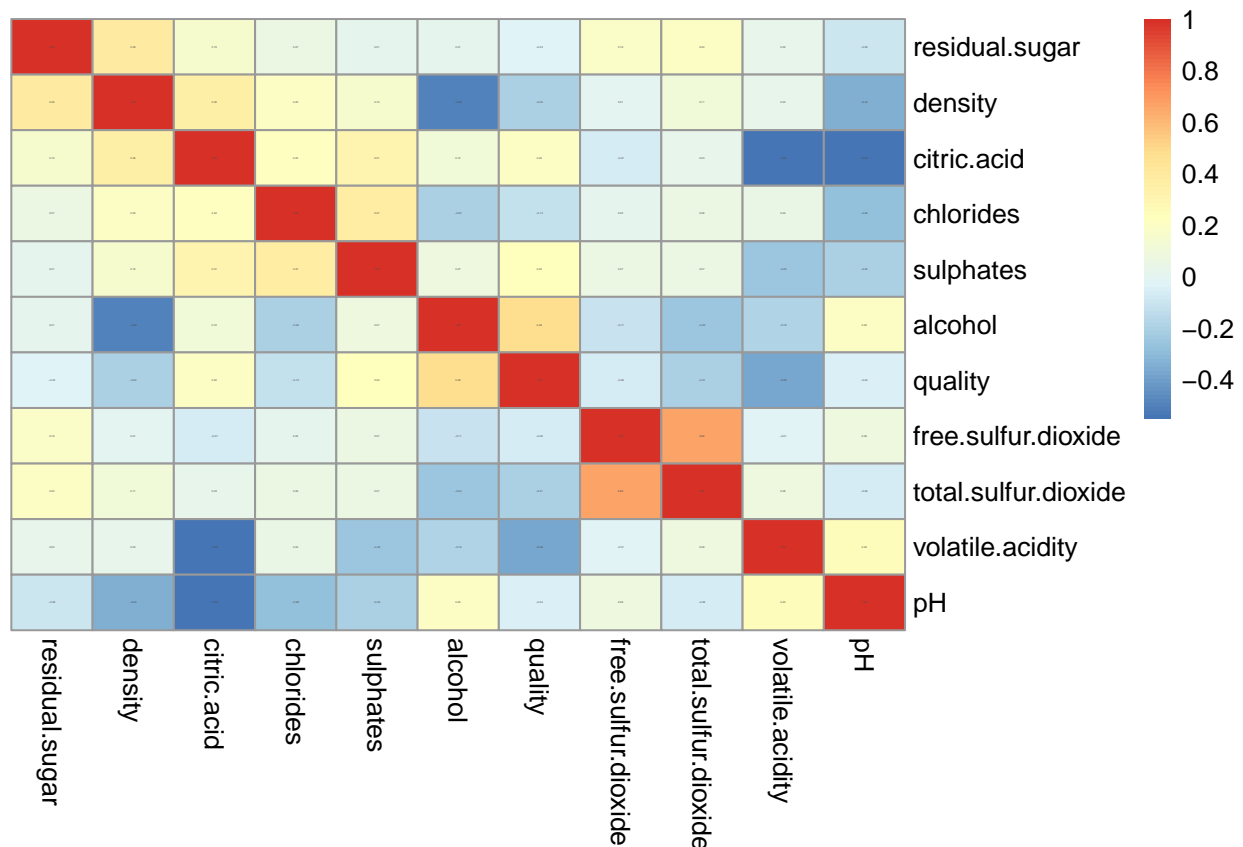
```
set.seed("2022")
wine.dataset <- read.csv("winequality-red.csv", sep = ";")
test.samples <- sample(1:nrow(wine.dataset), 400)
wine <- wine.dataset[-test.samples, ]
wine.test <- wine.dataset[test.samples, ]
```

To check the correlation between explanatory variables:

```
library(pheatmap)
```

```
## Warning: package 'pheatmap' was built under R version 4.1.3
```

```
corr.wine <- cor(wine[, -1])
pheatmap(corr.wine, treeheight_row = 0, treeheight_col = 0,
         display_numbers = T, fontsize_number = 0.5)
```



Great! The correlations are not as high as the diamond dataset we saw in the last lab, which means we do not need to worry too much about heteroscedasticity. We now fit the linear regression using all of the explanatory variables:

```
wine.fit <- lm(quality ~ . , data = na.omit(wine))
summary(wine.fit)
```

```
##
## Call:
## lm(formula = quality ~ . , data = na.omit(wine))
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.67557 -0.35919 -0.04682  0.46000  2.04930
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    12.5643244  24.4761128   0.513   0.6078
## fixed.acidity     0.0045448   0.0299388   0.152   0.8794
## volatile.acidity -1.0737942   0.1406020  -7.637 4.55e-14 ***
## citric.acid      -0.1394933   0.1661915  -0.839   0.4014
## residual.sugar    0.0007192   0.0175333   0.041   0.9673
## chlorides        -2.0052493   0.4682338  -4.283 2.00e-05 ***
## free.sulfur.dioxide 0.0061072   0.0025030   2.440   0.0148 *
## total.sulfur.dioxide -0.0038532  0.0008402  -4.586 4.99e-06 ***
## density          -8.0141715  24.9850874  -0.321   0.7485
## pH               -0.4859053   0.2198946  -2.210   0.0273 *
## sulphates         0.8267058   0.1259283   6.565 7.77e-11 ***
## alcohol           0.2822350   0.0302018   9.345 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.6496 on 1187 degrees of freedom
## Multiple R-squared:  0.356, Adjusted R-squared:  0.35
## F-statistic: 59.64 on 11 and 1187 DF, p-value: < 2.2e-16
```

Exercise 1 Confidence Interval

- (a) Calculate the confidence interval for all the coefficients from the regression done above. Which of these factors will positively influence the wine quality?

```
# Insert your code here to calculate the confidence intervals for the regression coefficients.
```

- (b) Calculate the confidence intervals for the samples in `wine.test` using the model you just fit. Which confidence interval will you use? Confidence intervals for the average response or the prediction interval?

```
# insert your code here and save your confidence intervals as `wine.confint`
# wine.confint <-
```

- (c) What is the percentage that your interval in (b) covers the true **quality** score in `wine.test`? What if you use the other confidence interval? Which one is consistent with your confidence level?

```
# insert your code here and save your percentage as `pct.covered`
# pct.covered <-
# pct.covered
# insert your code here and save your percentage calculated
# using the other confidence interval as `pct.covered.other`
# wine.confint.other <-
# pct.covered.other <-
# pct.covered.other
```

Exercise 2 Bootstrap CI

Scale the columns of the dataset using `scale()` and then make 95% bootstrap confidence intervals for the coefficients for the predictors. Plot these confidence intervals using the `plotCI()` function in `gplots`. Code from professor for making bootstrap CI is included. You can use this or write your own code. Use the wine subset as used above.

```
bootstrapLM <- function(y,x, repetitions, confidence.level=0.95){
  # calculate the observed statistics
  stat.obs <- coef(lm(y~., data=x))
  # calculate the bootstrapped statistics
  bootFun<-function(){
    sampled <- sample(1:length(y), size=length(y),replace = TRUE)
    coef(lm(y[sampled]~.,data=x[sampled,])) #small correction here to make it for a matrix x
  }
  stat.boot<-replicate(repetitions,bootFun())
  # nm <-deparse(substitute(x))
  # row.names(stat.boot)[2]<-nm
  level<-1-confidence.level
  confidence.interval <- apply(stat.boot,1,quantile,probs=c(level/2,1-level/2))
  return(list(confidence.interval = cbind("lower"=confidence.interval[1,],"estimate"=stat.obs,"upper"
```

```
# insert your code here
```

Regression dianosis

Red wine dataset

Reload the data.

```
wine<- read.csv("winequality-red.csv", sep = ";")
wine$quality <- wine$quality + rnorm(length(wine$quality))
```

Fit the model.

```
wine.fit <- lm(quality~volatile.acidity+chlorides+free.sulfur.dioxide+total.sulfur.dioxide+pH+sulphates
summary(wine.fit)
```

```
##
```

```
## Call:
## lm(formula = quality ~ volatile.acidity + chlorides + free.sulfur.dioxide +
##     total.sulfur.dioxide + pH + sulphates + alcohol, data = wine)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.7842 -0.7530 -0.0484  0.7729  3.6602
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    4.545629   0.728636   6.239 5.65e-10 ***
## volatile.acidity -0.927466   0.182365  -5.086 4.10e-07 ***
## chlorides      -1.995046   0.718916  -2.775  0.00558 **
## free.sulfur.dioxide  0.009023   0.003844   2.347  0.01903 *
## total.sulfur.dioxide -0.004939   0.001242  -3.977 7.29e-05 ***
## pH             -0.615917   0.212592  -2.897  0.00382 **
## sulphates       1.048310   0.198759   5.274 1.52e-07 ***
## alcohol         0.304831   0.030374  10.036 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.171 on 1591 degrees of freedom
## Multiple R-squared:  0.1621, Adjusted R-squared:  0.1584
## F-statistic: 43.96 on 7 and 1591 DF,  p-value: < 2.2e-16
```

Exercise 3

- (a) Do regression diagnostics using the `plot` function.

```
# insert your code here to do regression diagnostics.
```

- (b) Answer the following TRUE/FALSE questions based on the diagnostics plot. Uncomment your answer.

```
### I. The plot indicates heteroscedasticity.
# TRUE
# FALSE
### II. There are non-linearity between the explanatory variable and response variable.
# TRUE
# FALSE
### III. The normal assumption holds for this model.
# TRUE
# FALSE
```

- (c) Identify at least two outliers from the data.

I think the sample ??? and ??? are outliers.

Diamond dataset

Read the data.

```
diamonds <- read.csv("diamonds.csv")
diamonds <- diamonds[sample(1:nrow(diamonds), 1000), ]
head(diamonds)
```

```
##      carat      cut color clarity depth table price length.in.mm width.of.mm
## 32211  0.31   Premium    D     VS1  59.7   58   788         4.40         4.45
## 13189  1.21 Very Good    I     SI1  62.9   55  5452         6.69         6.76
## 26692  0.28 Very Good    H     VS1  61.9   56   429         4.18         4.20
## 39543  0.40   Premium    F     VS2  62.6   58  1080         4.72         4.68
## 46639  0.30   Premium    I     VVS2 61.7   58   526         4.28         4.37
## 32703  0.31    Ideal     F     VS2  60.8   57   802         4.39         4.36
##      depth.in.mm
## 32211         2.64
## 13189         4.23
## 26692         2.59
## 39543         2.94
## 46639         2.67
## 32703         2.66
```

Fit a linear regression.

```
diamond.fit <- lm(price ~ carat + cut + color + clarity + depth + table, data = diamonds)
summary(diamond.fit)
```

```
##
## Call:
## lm(formula = price ~ carat + cut + color + clarity + depth +
##      table, data = diamonds)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -7516.7  -591.1  -132.6   398.7  5612.1
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -3895.22    2582.20  -1.508   0.1318
## carat         8699.92     84.81 102.582 < 2e-16 ***
## cutGood        302.48    221.58   1.365   0.1725
## cutIdeal       469.12    224.07   2.094   0.0366 *
## cutPremium     414.88    211.96   1.957   0.0506 .
## cutVery Good   367.01    214.13   1.714   0.0868 .
## colorE        -205.90    135.93  -1.515   0.1302
## colorF        -309.58    134.25  -2.306   0.0213 *
## colorG        -536.38    133.14  -4.029 6.04e-05 ***
## colorH        -989.17    138.16  -7.160 1.59e-12 ***
## colorI       -1273.00    152.52  -8.347 2.37e-16 ***
## colorJ       -2119.39    187.23 -11.320 < 2e-16 ***
## clarityIF      5619.35    380.63  14.763 < 2e-16 ***
## claritySI1     3988.73    306.48  13.015 < 2e-16 ***
## claritySI2     3159.21    309.48  10.208 < 2e-16 ***
## clarityVS1     5024.21    313.58  16.022 < 2e-16 ***
## clarityVS2     4786.53    308.92  15.494 < 2e-16 ***
```



```
## clarityVVS1    5385.38      339.97  15.841 < 2e-16 ***
## clarityVVS2    5152.80      325.28  15.841 < 2e-16 ***
## depth         -14.04       27.93  -0.503  0.6154
## table         -43.17       20.10  -2.148  0.0320 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1085 on 979 degrees of freedom
## Multiple R-squared:  0.9235, Adjusted R-squared:  0.9219
## F-statistic: 590.8 on 20 and 979 DF,  p-value: < 2.2e-16
```

Exercise 4

- (a) Do regression diagnostics using the `plot` function.

```
# insert your code here to do regression diagnostics.
```

- (b) Answer the following TRUE/FALSE questions based on the diagnostics plot. Uncomment your answer.

```
### I. The plot indicates heteroscedasticity.
# TRUE
# FALSE
### II. There are non-linearity between the explanatory variable and response variable.
# TRUE
# FALSE
### III. The normal assumption holds for this model.
# TRUE
# FALSE
```

Multiple regression with continuous and categorical variables

Exercise 5

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

```
levels(diamonds$clarity)
```

```
## NULL
```

```
levels(diamonds$color)
```

```
## NULL
```

```
levels(diamonds$cut)
```

```
## NULL
```

```
# Insert you code here, save your model as `fit.categorical`  
# fit.categorical <-
```

(b) Write the equation when

i. Clarity is VS2, color is H, and cut is Premium. Replace ??? with numerical values.

$$\text{price} = ??? + ??? \cdot \text{carat} + ??? \cdot \text{depth} + ??? \cdot \text{table}$$

ii. clarity is I1, color is D and cut is Fair. Replace ??? by numerical values.

$$\text{price} = ??? + ??? \cdot \text{carat} + ??? \cdot \text{depth} + ??? \cdot \text{table}$$