

Applied Statistical Methods: Regression Analysis

Shizhe Chen

2020-03-28

Contents

Preface	5
1 Linear regression with R	7
1.1 Advertising data	7
1.2 Flu shot	8
1.3 Project STAR	10
1.4 Note	11
2 Estimation	13
3 Sampling distribution	17
3.1 Understanding sampling distribution via simulation	17
3.2 Shapes of sampling distributions	21
3.3 Small sample size	24
4 Statistical inference: confidence Intervals	27
4.1 Confidence interval	27
4.2 Prediction interval	30
4.3 Simultaneous confidence intervals/bands/regions	30
5 Statistical inference: Hypothesis testing	33
5.1 Hypothesis testing	33
5.2 Multiple testing	36
6 Model diagnostics	39
6.1 Residual plot	39
6.2 Remedies for non-linearity	40
6.3 Independence	40
6.4 Normality	41
6.5 Homoscedasticity	42
6.6 Influential Observations and Outliers	42
7 Multiple covariates	45
7.1 Examples	45
7.2 Classification of variables	46

7.3	Least squares estimation	48
7.4	Underfitting and overfitting	50
7.5	Sampling distribution and inference	52
8	Model selection	55
8.1	Criteria	55
8.2	Selection procedure	55
A	R basics	57
A.1	Summary statistics	57
A.2	Data structures	57
A.3	List	57
A.4	Functions in R	57
A.5	Miscellaneous	57
B	Linear algebra	59
B.1	Vector	59
B.2	Matrix	59
B.3	Other operations on vectors and matrices	59
C	Simulation and visualization	61
C.1	Simulation and visualization: univariate	61
C.2	Simulation and visualization: multivariate	61

Preface

This Gitbook only contains code and some comments that will be filled up slowly throughout the quarter. Lecture notes are available in separate files that you can find on Canvas.

Chapter 1

Linear regression with R

Reading materials: Slides 4 - 12 in STA108_LinearRegression_S20.pdf

Fitting a linear model is simple in R. The bare minimum requires you to know only two functions `lm()` and `summary()`. We will apply linear regression on three data set `advertising`, `flu shot`, and Project STAR.

1.1 Advertising data

This data is taken from An Introduction to Statistical Learning, by James et al. A brief description from Section 2.1 in ISL is provided below. You may read more about the data set in the free e-book.

The Advertising data set consists of the sales of that product in 200 different markets, along with advertising budgets for the product in each of those markets for three different media: TV, radio, and newspaper... It is not possible for our client to directly increase sales of the product. On the other hand, they can control the advertising expenditure in each of the three media. Therefore, if we determine that there is an association between advertising and sales, then we can instruct our client to adjust advertising budgets, thereby indirectly increasing sales. In other words, our goal is to develop an accurate model that can be used to predict sales on the basis of the three media budgets.

It is important to note that we will avoid making any causal statements (i.e., increasing TV advertising budget will increase the sales by ...) in our analysis. Causal inference on observational data is another major field in statistics. We will limit our discussion on association.

```
dat.advertising=read.csv('./data/advertising.csv');  
# Use the code in Appendix C to visualize this data set  
library(ggplot2)  
library(GGally)  
ggpairs(dat.advertising)
```

```
# What is X in this data frame?
dat.advertising=dat.advertising[,-1]
```

We will consider a simple linear regression model in this chapter. The only covariate/predictor/independent variable that we use here is the TV advertising budget. Therefore, the model is

$$y_i = x_i\beta_1 + \beta_0 + \epsilon_i, i = 1, \dots, 200,$$

where y_i is the sales (1 unit = 1000 dollars) for the i th entry, x_i is the TV advertising budget for the i th entry, β_0 is the intercept term, and β_1 is the regression slope. In addition, we assume that the errors $\{\epsilon_i\}_{i=1}^{200}$ satisfy that $\epsilon_1, \dots, \epsilon_{200}$ are independently and identically distributed (i.i.d.), $\mathbb{E}[\epsilon_i] = 0$ for $i = 1, 2, \dots, 200$ and $\text{var}(\epsilon_i) = \sigma^2$ for $i = 1, 2, \dots, 200$. Recall that we consider fixed design (i.e., x_i is not random) in this course for simplicity.

```
# Fit a simple linear regression
fit.advertising = lm(sales~TV+1,data=dat.advertising);
# Summarize the fitted results
summary(fit.advertising)
```

We see that $\hat{\beta}_1$ equals 0.05 and $\hat{\beta}_0$ equals 7.03. We can draw the line in the scatter plot.

How would you interpret the fitted slope and intercept?

```
# Draw the fitted line using ggplot2
ggplot(data = dat.advertising) +
  geom_point(mapping = aes(x = TV, y = sales)) +
  geom_line(data = fortify(fit.advertising), aes(x = TV, y = .fitted), color='red')
```

```
# In fact, you can use the following code without fitting the lm()
```

```
ggplot(data = dat.advertising aes(x = TV, y = sales)) +
  geom_point() +
  geom_smooth(method = "lm")
```

```
# The band is a confidence band
```

1.2 Flu shot

The data come from the following paper: McDonald, C., Hiu, S., AND Tierney, W. (1992). Effects of computer reminders for influenza vaccination on morbidity during influenza epidemics. MD Computing 9, 304-312. ([link](#)).

Briefly, in this study, physicians were randomly selected to receive a letter encouraging them to inoculate patients at risk for flu. The treatment of interest is the actual flu shot, and the outcome is an indicator for flu-related hospital visits from the patients. In this data set, you can also find four important background covariates including gender, age, a chronic obstructive pulmonary disease (COPD) indicator, and a heart disease indicator. For further

details, read the original publication or Hirano, Imbens, Rubin and Zhou (2000) (link), or Chapter 25 in Imbens and Rubin (2015) (link).

```
dat.flu = read.table("./data/flu240.txt", header = TRUE)

# Take a look at size of the data before proceed:
dim(dat.flu)
colnames(dat.flu)
# Basic visulization:
ggpairs(dat.flu) # This might take a while..

# Do you think the plots are informative?

# One strategy: jittering

# Without jittering:
ggplot(dat.flu, aes(x=treatment.assigned, y=outcome)) + geom_point()

# With jittering:
ggplot(dat.flu, aes(x=treatment.assigned, y=outcome)) + geom_jitter(width = 0.25)
```

Again, we will consider a simple linear regression model for now. The only covariate/predictor/independent variable that we use here is the treatment assignment. Therefore, the model is

$$y_i = x_i\beta_1 + \beta_0 + \epsilon_i, i = 1, \dots, 2891,$$

where y_i is the outcome (0 or 1) for the i th patient, x_i is the treatment assignment (0 or 1) for the i th patient's physician, β_0 is the intercept term, and β_1 is the regression slope. In addition, we assume that the errors $\{\epsilon_i\}_{i=1}^{2891}$ satisfy that $\epsilon_1, \dots, \epsilon_{2891}$ are independently and identically distributed (i.i.d.), $\mathbb{E}[\epsilon_i] = 0$ for $i = 1, 2, \dots, 2891$ and $\text{var}(\epsilon_i) = \sigma^2$ for $i = 1, 2, \dots, 2891$. The following code is almost identical to those in the previous section.

How would you interpret the fitted slope and intercept?

Are there any violations to the assumptions on $\{\epsilon_i\}_{i=1}^{2891}$?

```
# Fit a simple linear regression
fit.flu= lm(outcome~treatment.assigned+1,data=dat.flu);
# Summarize the fitted results
summary(fit.flu)

# Draw the fitted line using ggplot2
ggplot(dat.flu, aes(x=treatment.assigned, y=outcome)) +
  geom_jitter(width = 0.25)+
  geom_smooth(method = "lm")
```

1.3 Project STAR

Tennessee Student/Teacher Achievement Ratio study (Project STAR) was conducted in the late 1980s to evaluate the effect of class size on test scores. This dataset has been used as a classic examples in many textbooks and research papers. Briefly, the study randomly assigned students to small classes, regular classes, and regular classes with a teacher's aide. In order to randomize properly, schools were enrolled only if they had enough studybody to have at least one class of each type. Once the schools were enrolled, students were randomly assigned to the three types of classes, and one teacher was randomly assigned to one class. You can read more about the study on the Harvard dataverse ([link](#)) or Chapter 9 in Imbens and Rubin (2015) ([link](#)).

```
library(AER)
data("STAR")
dat.STAR=STAR; # Just to be consistent
# Take a look at size of the data before proceed:
dim(dat.STAR)
colnames(dat.STAR)
# With 47 variables, you might want to read the help file to see what these variables

# Basic visulization:
# ggpairs(dat.flu) # This might take a really long time ...

# We will consider the math scores and the class assignments in the second grade

# With jittering:
ggplot(dat.STAR, aes(x=star2, y=math2)) +
  geom_jitter(width = 0.25)
```

Again, we will consider a simple linear regression model for now. We will study the association between math scores and the class assignments in the second grade. Therefore, the model is

$$y_i = \beta_0 + x_{i,1}\beta_1 + x_{i,2}\beta_2 + \epsilon_i, i = 1, \dots, 2891,$$

where y_i is the math score for the i th student, $x_{i,1}$ is an indicator (0 or 1) whether the i th student is in the small class, and $x_{i,2}$ is an indicator (0 or 1) whether the i th student is in the regular class with aide, β_0 is the intercept term, and β_1 and β_2 are the regression coefficients for $x_{i,1}$ and $x_{i,2}$, respectively. In addition, we assume that the errors $\{\epsilon_i\}_{i=1}^{2891}$ satisfy that $\epsilon_1, \dots, \epsilon_{2891}$ are independently and identically distributed (i.i.d.), $\mathbb{E}[\epsilon_i] = 0$ for $i = 1, 2, \dots, 2891$ and $\text{var}(\epsilon_i) = \sigma^2$ for $i = 1, 2, \dots, 2891$. The following code is almost identical to those in the previous section.

```
# Fit a simple linear regression
fit.STAR= lm(math2~as.factor(star2)+1,data=dat.STAR);
# Summarize the fitted results
summary(fit.STAR)
```

```
# We can no longer draw a fitted line here...
```

1.4 Note

In the the midterm projects, you will be asked to reproduce the analysis on the **advertising** data using your own functions. The hope is that, by reinventing the wheels, you will have a thorough understanding of regression.

Chapter 2

Estimation

Reading materials: Slides 13 - 23 in STA108_LinearRegression_S20.pdf

```
# Consider the advertising dataset
dat.advertising=read.csv('./data/advertising.csv');
# obtain the residuals:
x=dat.advertising$TV;
y=dat.advertising$sales;
beta1=1;beta0=0.2; # Just two arbitrary numbers
res=y-x*beta1-beta0;

# What if we want to feed in different numbers?

# Write a function that can calculate x*beta1 + beta0

linear.model<-function(beta,covariate){
  # beta: a 2-vector, where the first entry is the intercept
  yout=covariate*beta[2]+beta[1]
  return(yout);
  # Note: this function only works with simple linear regression model
  # How would you generalize it?
}

res.new=dat.advertising$sales-linear.model(beta=c(0.2,1),covariate=dat.advertising$TV);

identical(res.new,res)

# With the above function, we can reproduce the plot in Slide 14
beta=c(4,0.03);
plot(x=dat.advertising$TV, y=dat.advertising$sales, pch=16,cex=1.3,col='blue',xlab='x',y
abline(a=beta[1],b=beta[2],col='red',lwd=3)
```

```

yout=linear.model(beta=beta,covariate=dat.advertising$TV);
for(i in 1:dim(dat.advertising)[1]){
  segments(x0=dat.advertising$TV[i], y0=yout[i], y1 = dat.advertising$sales[i],col='blue')
}

# We can see that the line is less than ideal
# We will move on to find the minimizer of the squared error loss

sum.of.squares<-function(beta,covariate,outcome){
  yout=linear.model(beta=beta,covariate=covariate);
  res=outcome-yout;
  sos= sum(res^2);
  return(sos)
}

# There are many ways to minimize the sum of squares

# 1. Use generic optimizer in R
fit.optim=optim(c(0,0),sum.of.squares,covariate=dat.advertising$TV,outcome=dat.advertising$sales)
beta.hat.optim=fit.optim$par;
# 2. Use analytic form of the optimizer

fit.linear.model<-function(covariate,outcome){
  # I will write down the function for (multiple) linear regression here
  X=cbind(1,covariate);
  beta.fit=solve( t(X)%*%X )%*%t(X)%*%outcome;
  return(beta.fit)
}

beta.hat=fit.linear.model(covariate=dat.advertising$TV,outcome=dat.advertising$sales)

# 3. Write your own version of Newton-Raphson method, or gradient descent
# Not required or taught in this class

# We might want to check if our solution matches that from lm()

fit.advertising=lm(sales~TV+1,data=dat.advertising);

fit.advertising$coef
beta.hat
beta.hat.optim

# Calculate the decomposition of sum of squares
residual.sum.of.squares=sum.of.squares(beta=beta.hat,covariate=dat.advertising$TV,outcome=dat.advertising$sales)
explained.sum.of.squares=sum((linear.model(beta=beta.hat,covariate=dat.advertising$TV)-dat.advertising$sales)^2)

```

```

total.sum.of.squares=sum((dat.advertising$sales-mean(dat.advertising$sales))^2)

# Check if the equality holds..
total.sum.of.squares
explained.sum.of.squares+residual.sum.of.squares

# But if we use identical()...
identical(total.sum.of.squares,explained.sum.of.squares+residual.sum.of.squares)

# Calculate the coefficient of determination
R.sq= explained.sum.of.squares/total.sum.of.squares;

# Calculate the Pearson sample correlation
cor.hat=cor(dat.advertising$TV,dat.advertising$sales)

# Verify the identity:
R.sq
cor.hat^2

```

Q: what if these equalities only hold on this one data set? What would you do to rea

Chapter 3

Sampling distribution

3.1 Understanding sampling distribution via simulation

Reading materials: Slides 24 - 33 in STA108_LinearRegression_S20.pdf.

In this section, we assume a true data generating model in order to draw samples from the true population distribution. In particular, we assume that

$$y_i = x_i\beta_1 + \beta_0 + \epsilon_i, \quad i = 1, \dots, 50$$

where $\epsilon_i \sim \mathcal{U}(-2.5, 2, 5)$. We set $\beta_1 = 0.15$ and $\beta_0 = 20$.

```
### We will simulate synthetic data to understand the concept of sampling distribution

### To match our setup, we will generate the covariates and keep them fixed in the remain
set.seed(1)
n=50;
x=as.matrix(rnorm(n,mean=10,sd=2),ncol=n);
beta.true=c(20,0.15)
Ey= x*beta.true[2]+beta.true[1];

error.terms= (runif(n)-0.5)*5;
y=Ey+error.terms;
beta.hat=fit.linear.model(covariate=x,outcome=y);

plot(x=x,y=y,pch=16,cex=2)
abline(b=beta.hat[2],a=beta.hat[1],col='red',lwd=2)

### We can run the simulation by putting the above lines into a large for-loop
### We will wrap them up into one function
```

```

simulate.one.instance<-function(x,beta.true){
  n=length(x);
  Ey= x*beta.true[2]+beta.true[1];
  error.terms= (runif(n)-0.5)*5;
  y=Ey+error.terms;
  beta.hat=fit.linear.model(covariate=x,outcome=y);
  return(beta.hat)
}

### Set the number of replicates to be 10000
N.sim=1e4;
set.seed(1)
beta.sim=replicate(N.sim,simulate.one.instance(x,beta.true));

# Draw the histogram for the estimated intercept
hist(beta.sim[1,1,],xlab='Fitted intercept',main='')
abline(v=beta.true[1],lwd=3,col='red')

```

Recall from the lecture notes that the variance of $\hat{\beta}_1$ is

$$\text{var}(\hat{\beta}_1) = \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sigma^2.$$

```

## Using the simulation results, we can verify the statement on Slide 26

### Expectation
mean(beta.sim[1,1,])-beta.true[1] # intercept
mean(beta.sim[2,1,])-beta.true[2] # slope

### Variance
sigma.sq = 5^2/12; # variance of the error: (runif(n)-0.5)*5;

# slope:
sigma.sq/sum( (x-mean(x))^2 ) # theoretical value
var(beta.sim[2,1,]) # variance from the simulation

# intercept:
sigma.sq*sum(x^2)/sum( (x-mean(x))^2 )/n # theoretical value
var(beta.sim[1,1,]) # variance from the simulation

```

Q: How would you verify the Gauss-Markove theorem in this simulation?

```

# A very simple simulation to compare the LSE with another linear unbiased estimator:
simulate.one.instance<-function(x,beta.true){
  n=length(x);
  Ey= x*beta.true[2]+beta.true[1];
  error.terms= (runif(n)-0.5)*5;
  y=Ey+error.terms;
  beta.hat=fit.linear.model(covariate=x,outcome=y);
  slope.hat=beta.hat[2];
  slope.dot= (y[n]-y[1])/(x[n]-x[1])
  return(c(slope.hat,slope.dot))
}

### Set the number of replicates to be 10000
N.sim=1e4;
set.seed(1)
slope.sim=replicate(N.sim,simulate.one.instance(x,beta.true));

apply(slope.sim,1,mean) # both are unbiased
apply(slope.sim,1,var) # variance of LSE is smaller

## Estimate the variance of errors using the residuals
set.seed(1)
error.terms= (runif(n)-0.5)*5;
y=Ey+error.terms;
beta.hat=fit.linear.model(covariate=x,outcome=y);

residual.sum.of.squares=sum.of.squares(beta=beta.hat,covariate=x,outcome=y)

sigma.sq.hat=residual.sum.of.squares/(n-2) # estimates from one instance

## Run a simulation to verify the claim

## NOTE: we will use the same function name
simulate.one.instance<-function(x,beta.true){
  n=length(x);
  Ey= x*beta.true[2]+beta.true[1];
  error.terms= (runif(n)-0.5)*5;
  y=Ey+error.terms;
  beta.hat=fit.linear.model(covariate=x,outcome=y);
  residual.sum.of.squares=sum.of.squares(beta=beta.hat,covariate=x,outcome=y)

sigma.sq.hat=residual.sum.of.squares/(n-2) # estimates from one instance

  return(sigma.sq.hat)
}

```

```

}

### Set the number of replicates to be 10000
N.sim=1e4;
set.seed(1)
sigma.sq.hat.sim=replicate(N.sim,simulate.one.instance(x,beta.true));

mean(sigma.sq.hat.sim)
sigma.sq

## To wrap this up
## For any new fits, we can estimate the variance and standard errors of the estimators
set.seed(1)
error.terms= (runif(n)-0.5)*5;
y=Ey+error.terms;
beta.hat=fit.linear.model(covariate=x,outcome=y);

estimate.sigma.sq<-function(beta,covariate,outcome){
  residual.sum.of.squares=sum.of.squares(beta=beta,covariate=covariate,outcome=outcome)
  n=length(outcome)
  sigma.sq.hat=residual.sum.of.squares/(n-2)
  return(sigma.sq.hat)
}

estimate.coef.var<-function(beta,covariate,outcome){
  sigma.sq.hat=estimate.sigma.sq(beta,covariate,outcome)
  var.hat.beta=beta;
  var.hat.beta[2]=sigma.sq.hat/sum( (covariate-mean(covariate))^2 )
  n=length(outcome)
  var.hat.beta[1]=sigma.sq.hat*sum(covariate^2)/sum((covariate-mean(covariate))^2 )/n
  return( var.hat.beta)
}

estimate.coef.sd<-function(beta,covariate,outcome){
  var.hat.beta=estimate.coef.var(beta,covariate,outcome)
  sd.hat.beta=sqrt(var.hat.beta);
  return(sd.hat.beta)
}
estimate.coef.sd(beta=beta.hat,covariate=x,outcome=y)

## Compare your numbers with the output from lm()
summary(lm(y~x+1))

```

3.2 Shapes of sampling distributions

Reading materials: Slide 34 in STA108_LinearRegression_S20.pdf.

There can be multiple distributions with the same mean and variance. Consider three

```
xgrid = seq(from = -5, to = 5, by=0.01);
norm_density=dnorm(xgrid,mean=0,sd=sqrt(2));
unif_density=dunif(xgrid,min=-sqrt(6),max=sqrt(6));

xbin=c(-sqrt(2),sqrt(2));
ybin=c(0.5,0.5);

plot(norm_density~xgrid,type="l",col='red',lwd=3,ylim=c(0,0.52),xlab="z",ylab="density")
lines(unif_density~xgrid,col='black',lwd=3)

segments(xbin[1],0,xbin[1],ybin[1],col='green',lwd=10)
segments(xbin[2],0,xbin[2],ybin[2],col='green',lwd=10)
legend(x=2.5,y=0.5,legend=c("Uniform","Normal","Binary"),lwd=2,col=c('black','red','green'))
```

3.2.1 Asymptotic distribution

Reading materials: Slides 35 - 49 in STA108_LinearRegression_S20.pdf.

The L-F central limit theorem claims that the asymptotic distribution for $\frac{\hat{\beta}_1 - \beta_1}{\hat{\sigma}/S_{xx}^{1/2}}$ is $\mathcal{N}(0, 1)$.

```
## Edit the simulation code
### This time with large n (n=500)
### Three different versions of errors
### Unequal variances
simulate.one.instance<-function(x,beta.true,error.type){
  n=length(x);
  Ey= x*beta.true[2]+beta.true[1];
  vars=0.1+runif(n)*2; # unequal variance
  if(error.type=='uniform'){
    error.terms= (runif(n)-0.5)*sqrt(vars);

  }else if(error.type=='Bernoulli'){
    error.terms=(rbinom(n,size=1,prob=0.5)-0.5)*sqrt(vars);
  }else{
    error.terms= rnorm(n)*sqrt(vars);
  }
  y=Ey+error.terms;
  beta.hat=fit.linear.model(covariate=x,outcome=y);
```

```

    residual.sum.of.squares=sum.of.squares(beta=beta.hat,covariate=x,outcome=y)

sigma.sq.hat=residual.sum.of.squares/(n-2) # estimates from one instance

    beta.hat.std =(beta.hat[2]-beta.true[2])/sqrt(sigma.sq.hat)*sqrt( sum( (x-mean(x))^2 )

    return(beta.hat.std)
}

### Set the number of replicates to be 10000
N.sim=1e4;
set.seed(1)

n=500;
x=as.matrix(rnorm(n,mean=10,sd=2),ncol=n);
beta.hat.std.uniform=replicate(N.sim,simulate.one.instance(x,beta.true,error.type='uniform'))
beta.hat.std.Bernoulli=replicate(N.sim,simulate.one.instance(x,beta.true,error.type='Bernoulli'))
beta.hat.std.normal=replicate(N.sim,simulate.one.instance(x,beta.true,error.type='normal'))

# Draw the density plot

density.uniform=density( beta.hat.std.uniform);
density.Bernoulli=density( beta.hat.std.Bernoulli);
density.normal=density( beta.hat.std.normal);

# Theoretical density:
xgrid = seq(from=-4,to=4,length.out=100);
normal.pdf= dnorm(xgrid)

plot(density.uniform,xlab="Standardized intercept",col='red',lwd=3,ylim=c(0,0.5),main='')
lines(density.Bernoulli,col='green',lwd=3,lty=2)
lines(density.normal,col='blue',lwd=3,lty=2)
lines(normal.pdf~xgrid,lwd=2,lty=3)
legend(x=2,y=0.46,legend=c("Uniform errors","Bernoulli errors","Normal errors", "N(0,1)"))

## The L-F CLT allows us to write down the asymptotic distribution in closed-form
## The bootstrap method is an alternative when analytic solution is not available
## Note: Bootstrap is not that useful for linear regression, but has proven to be a power

set.seed(1)
# Generate one set of data:
n=500;
x=as.matrix(rnorm(n,mean=10,sd=2),ncol=n);
Ey= x*beta.true[2]+beta.true[1];
vars=0.1+runif(n)*2; # unequal variance

```

```

error.terms= (runif(n)-0.5)*sqrt(vars);
y=Ey+error.terms;
beta.hat=fit.linear.model(covariate=x,outcome=y);

# Now write a function for bootstrap
boot.one.instance<-function(covariate,outcome){
  n=length(outcome);
  sample_indices = sample(1:n,n,replace=TRUE) # sampling with replacement
  covariate.boot= covariate[sample_indices]; outcome.boot= outcome[sample_indices];

  beta.hat=fit.linear.model(covariate=covariate.boot,outcome=outcome.boot);
  return(beta.hat[2])
}

set.seed(1)

B=1e5;
beta.hat.boot=replicate(B,boot.one.instance(covariate=x,outcome=y));

## Compare the bootstrap distribution of the standardized slope against that from the CLT

density.boot=density((beta.hat.boot-mean(beta.hat.boot))/sd(beta.hat.boot));
# Theoretical density (from CLT)
xgrid = seq(from=-4,to=4,length.out=100);
normal.pdf= dnorm(xgrid)

plot(density.boot,xlab="Standardized intercept",col='red',lwd=3,ylim=c(0,0.5),main='',lty=1)
lines(normal.pdf~xgrid,lwd=2,lty=3)
legend(x=2,y=0.46,legend=c("Bootstrap", "N(0,1)"),lwd=3,col=c('red','green','blue','black'))

## Compare the distributions using the raw bootstrap distribution, and the distribution

density.boot=density(beta.hat.boot);
# Theoretical density (from CLT)
xgrid = seq(from=min(beta.hat.boot),to=max(beta.hat.boot),length.out=100);
normal.pdf= dnorm(xgrid,mean=mean(beta.hat.boot),sd=sd(beta.hat.boot))

plot(density.boot,xlab="Standardized intercept",col='red',lwd=3,main='',lty=2)
lines(normal.pdf~xgrid,lwd=2,lty=3)
legend(x=0.15,y=50,legend=c("Empirical", "Moment-based"),lwd=3,col=c('red','green','blue'))

```

3.3 Small sample size

Reading materials: Slides 50 - 57 in STA108_LinearRegression_S20.pdf.

When the sample size is small, the asymptotic distribution is not close to the true sampling distribution. Neither the central limit theorem or the bootstrap distribution are good approximation of the asymptotic distribution.

```
# Simulation with small sample size
# We will demonstrate the case with 5 samples
# You will find the code in this trunk almost identical with previous code

# Generate one set of data:
set.seed(1)
n=10;
x=as.matrix(rnorm(n,mean=10,sd=2),ncol=n);
beta.true=c(20,0.15)
Ey= x*beta.true[2]+beta.true[1];

error.terms= rnorm(n)*5; # Normal errors!
y=Ey+error.terms;
beta.hat=fit.linear.model(covariate=x,outcome=y);
plot(x=x,y=y,pch=16,cex=2)
abline(b=beta.hat[2],a=beta.hat[1],col='red',lwd=2)

## Simulate the true sampling distribution:
simulate.one.instance<-function(x,beta.true){
  n=length(x);
  Ey= x*beta.true[2]+beta.true[1];
  error.terms= (runif(n)-0.5)*5;
  y=Ey+error.terms;
  beta.hat=fit.linear.model(covariate=x,outcome=y);
  return(beta.hat)
}

### Set the number of replicates to be 10000
N.sim=1e4;
true.beta.hat=replicate(N.sim,simulate.one.instance(x,beta.true));

### Obtain the mean and standard deviation using CLT:
residual.sum.of.squares=sum.of.squares(beta=beta.hat,covariate=x,outcome=y)
sigma.sq.hat=residual.sum.of.squares/(n-2) # estimates from one instance

### Obtain the bootstrap distribution
B=1e5;
```



```

boot.beta.hat=replicate(B,boot.one.instance(covariate=x,outcome=y));

density.true=density( (true.beta.hat[2,1,]-mean(true.beta.hat[2,1,]))/sd(true.beta.hat[2,1,]));

density.boot=density((boot.beta.hat-mean(boot.beta.hat))/sd(boot.beta.hat));
# Theoretical density (from CLT)
xgrid = seq(from=-4,to=4,length.out=100);
normal.pdf= dnorm(xgrid)

plot(density.true,xlab="Estimated intercept",lwd=3,ylim=c(0,0.5),xlim=c(-4,4),main='',lty=1)
lines(density.boot,col='green',lwd=3,lty=2)
lines(normal.pdf~xgrid,col='red',,lwd=2,lty=2)
legend(x=2,y=0.5,legend=c("True","Bootstrap","CLT"),lwd=3,col=c('black','green','red'),lty=c(1,2,2))

## With very few samples, we have to exploit the parametric assumption
## In this case, we take advantage of the normality assumption on the errors

t.pdf= dt(xgrid,df=n-2)

plot(density.true,xlab="Estimated intercept",lwd=3,ylim=c(0,0.5),xlim=c(-4,4),main='',lty=1)
lines(density.boot,col='green',lwd=3,lty=2)
lines(normal.pdf~xgrid,col='red',,lwd=2,lty=2)

lines(t.pdf~xgrid,col='blue',,lwd=2,lty=2)
legend(x=2,y=0.5,legend=c("True","Bootstrap","CLT","Student t"),lwd=3,col=c('black','green','blue','red'),lty=c(1,2,2,2))

## The Student t distribution grows more similar to the standard normal as its degrees of freedom increase

t_dist=dt(xgrid,df=1);
plot(t_dist~xgrid,xlab="beta 1 hat",col=rgb(1,0,0,0.1),lwd=3,type="l",ylim=c(0,0.5),xlim=c(-4,4))

for(i in 2:10){
  t_dist=dt(xgrid,df=i);

  lines(t_dist~xgrid,col=rgb(1,0,0,i/10),lwd=3)
}

legend(x=2,y=0.46,legend=c("Student t", "N(0,1)"),lwd=3,col=c('red','green'))

lines(normal.pdf~xgrid,col=rgb(0,1,0,0.4),lwd=5)

```


Chapter 4

Statistical inference: confidence Intervals

Reading materials: Slides 58 - 73 in STA108_LinearRegression_S20.pdf

4.1 Confidence interval

```
## We will look at synthetic data here, because we have control over the truth
set.seed(1)
n=50;
x=as.matrix(rnorm(n,mean=10,sd=2),ncol=n);
beta.true=c(20,0.15)
Ey= x*beta.true[2]+beta.true[1];
error.terms= rnorm(n)*5;
y=Ey+error.terms;
```

For a given confidence level alpha, construct a 100(1-alpha)% confidence interval

```
alpha=0.023;
# There is a function confint() in R
fit.lm=lm(y~x+1);
confint(fit.lm,level=1-alpha)
```

We will implement our own version

```
## We actually have most of the part ready
beta.hat=fit.linear.model(covariate=x,outcome=y);
beta.sd=estimate.coef.sd(beta=beta.hat,covariate=x,outcome=y);
```

```
## The only missing piece is the quantile
conf.int.quantile<-function(alpha,type,...){
```

```

if(type=="t"){
  out=qt(c(1-alpha/2,alpha/2), ... )
}else if (type=="normal"){
  out=qnorm(c(1-alpha/2,alpha/2), ... )
}
return(out)
}
quants<-conf.int.quantile(alpha,type='t',df=n-2)

beta.hat%*%c(1,1)-beta.sd%*%quants

## Compare with the output from confint
confint(fit.lm,level=1-alpha)

## How about using bootstrap to construct the CIs?
boot.fit<-function(covariate,outcome){
  n=length(outcome);
  sample_indices = sample(1:n,n,replace=TRUE) # sampling with replacement
  covariate.boot= covariate[sample_indices]; outcome.boot= outcome[sample_indices];

  beta.hat=fit.linear.model(covariate=covariate.boot,outcome=outcome.boot);
  return(t(beta.hat ))
}
B=1e5;
beta.hat.boot=replicate(B,boot.fit(covariate=x,outcome=y));

dim(beta.hat.boot)
apply(beta.hat.boot[1,,],1,quantile,probs=c(alpha/2,1-alpha/2))

# We can wrap this up into a function
conf.int<-function(alpha,type,covariate,outcome,B=1e5){

  beta.hat=fit.linear.model(covariate,outcome);
  beta.sd=estimate.coef.sd(beta=beta.hat,covariate,outcome);
  if(type=='bootstrap'){
    beta.hat.boot=replicate(B,boot.fit(covariate,outcome));
    out=t(apply(beta.hat.boot[1,,],1,quantile,probs=c(alpha/2,1-alpha/2)));
  }else if(type=='t'){
    quants<-conf.int.quantile(alpha,type='t',df=n-2)
    out=beta.hat%*%c(1,1)-beta.sd%*%quants;
  }else{
    quants<-conf.int.quantile(alpha,type='normal')
    out=beta.hat%*%c(1,1)-beta.sd%*%quants;
  }
}

```

```

    colnames(out)=c( paste(round(alpha*50,digits=3),'%'), paste(100-round(alpha*50,digits=3),'%'))
    return(out)
}
conf.int(alpha=alpha,type='t',covariate=x,outcome=y)
confint(fit.lm,level=1-alpha)

## Verify the coverage of the confidence intervals
simulate.one.instance<-function(x,beta.true,alpha){
  n=length(x);
  Ey= x*beta.true[2]+beta.true[1];
  error.terms= rnorm(n)*5;
  y=Ey+error.terms;
  CIs=conf.int(alpha=alpha,type='t',covariate=x,outcome=y);
  return(CIs)
}

N.sim=1e4;
set.seed(1)
n=50;
x=as.matrix(rnorm(n,mean=10,sd=2),ncol=n);
beta.true=c(20,0.15)
sim.CIs=replicate(N.sim,simulate.one.instance(x,beta.true,alpha));

## Visualize and verify the coverage for the slope
coverage=sum(sim.CIs[2,1,]<beta.true[2] & sim.CIs[2,2,]>beta.true[2])/N.sim;
1-alpha

plot(0,0,col="white",xlim=c(min(sim.CIs[2,1,]),max(sim.CIs[2,2,])),ylim=c(0,20),xlab="Slope",ylab="Coverage")
abline(v=beta.true[2],lwd=3,col='red')
for(i in 1:19){
  segments(sim.CIs[2,1,i],i,sim.CIs[2,2,i],i,lwd=3)
}

```

4.1.1 Determine cutoffs

We can define multiple confidence intervals of the same confidence levels for on the same sampling distribution.

```

alpha=0.05;

cutoffs=c(0.01,0.025,0.049);

colorlist=c('red','green','blue');
leg_text=c("1","2","3");

```

```

normal.pdf= dnorm(xgrid)
plot(normal.pdf~xgrid,xlab="beta 1 hat",ylab="Density",col=rgb(0,0,0,1),lwd=3,type="l",y

for(i in 1:length(cutoffs)){
  lower_bound = qnorm(cutoffs[i]);
  upper_bound = qnorm(1-alpha+cutoffs[i]);
  segments(lower_bound,0,lower_bound,0.5,col=colorlist[i],lwd=3)
  segments(upper_bound,0,upper_bound,0.5,col=colorlist[i],lwd=3)
  segments(lower_bound,i*0.1,upper_bound,i*0.1,col=colorlist[i],lwd=3)
  leg_text[i] =paste(cutoffs[i],"to", 1-alpha+cutoffs[i], "(length:", signif(upper_bound-1
}
legend(x=-5,y=0.5,legend=leg_text,lwd=3,col=colorlist)

```

4.2 Prediction interval

```

set.seed(1)
n=50;
x=as.matrix(rnorm(n,mean=10,sd=2),ncol=n);
beta.true=c(20,0.15)
Ey= x*beta.true[2]+beta.true[1];
error.terms= rnorm(n)*5;
y=Ey+error.terms;
fit.lm=lm(y~as.vector(x)+1);

alpha=0.05;
dat.new = data.frame(x=c(10,20));

# Prediction interval (assume normality)
predict(fit.lm,newdata=dat.new,interval="prediction",level=1-alpha)

# and confidence interval
predict(fit.lm,newdata=dat.new,interval="confidence",level=1-alpha)

# Implementation is similar to our conf.int function, and is thus skipped.
# You can verify the coverage by modifying the code above.

```

4.3 Simultaneous confidence intervals/bands/regions

```

## We will return to the advertising data
fit.lm=lm(sales~TV+1,data=dat.advertising)
fit.lm.sum=summary(fit.lm)
alpha=0.05;

```

```

xgrid_pred=seq(from=5,to=50,by=2.5);
dat.new = data.frame(TV=xgrid_pred);
pointwise=predict(fit.lm,newdata=dat.new, interval="confidence", level=1-alpha)

# Bonferroni corrected:
Bonf=predict(fit.lm,newdata=dat.new, interval="confidence", level=1-alpha/length(xgrid_pred))

# Working-Hotelling
fstat= qf(1-alpha,2,length(dat.advertising$TV)-2);
xdense=seq(from=0,to=50,by=0.1);
Sxx= sum( (dat.advertising$TV-mean(dat.advertising$TV))^2);
pivot=sqrt(2*fstat)*fit.lm.sum$sigma*sqrt(1/length(dat.advertising$TV)+(xdense-mean(dat.advertising$TV))^2);
ylb=fit.lm$coefficients[2]*xdense+fit.lm$coefficients[1]-pivot;
yub=fit.lm$coefficients[2]*xdense+fit.lm$coefficients[1]+pivot;

plot(dat.advertising$sales~dat.advertising$TV,pch=16,xlim=c(0,50),ylim=c(4,12),xlab='TV',ylab='sales')
abline(b=fit.lm$coefficients[2],a=fit.lm$coefficients[1],lwd=3)
for(i in 1:length(xgrid_pred)){
  segments(xgrid_pred[i]-1,Bonf[i,2],xgrid_pred[i]-1,Bonf[i,3],lwd=3,col='green')
  segments(xgrid_pred[i]-1,pointwise[i,2],xgrid_pred[i]-1,pointwise[i,3],lwd=3,col='red')
}
lines(ylb~xdense,col='blue',lwd=3)
lines(yub~xdense,col='blue',lwd=3)
legend(x=32,y=6,legend=c("Pointwise", "Bonferroni", "W-H band"), lwd=3, col=c("red","green","blue"))

```


Chapter 5

Statistical inference: Hypothesis testing

Reading materials: Slides 74 - 92 in STA108_LinearRegression_S20.pdf.

We will continue working on synthetic data in this chapter.

```
set.seed(1)
n=50;
x=as.matrix(rnorm(n,mean=10,sd=2),ncol=n);
beta.true=c(20,0.15)
Ey= x*beta.true[2]+beta.true[1];
error.terms= rnorm(n)*2; # Normal errors, in order to verify the performance of the t-
# You can change it to other distributions
y=Ey+error.terms;
```

5.1 Hypothesis testing

```
# The lm() and summary() functions will automatically calculate the test statistics and
fit.lm<-lm(y~x+1);fit.lm.summary<-summary(fit.lm)
fit.lm.summary
```

```
# We will create our own function for calculating the "t value"
source('sources.r') # load the functions we have written

beta.hat=fit.linear.model(covariate=x,outcome=y)
beta.hat.sd=estimate.coef.sd(beta=beta.hat,covariate=x,outcome=y)

# The test statistics for H0: beta1=0, and H0: beta0=0 are
beta.hat.t = (beta.hat-0)/beta.hat.sd;
```

```

calculate.t<-function(covariate,outcome){
  beta.hat=fit.linear.model(covariate=covariate,outcome=outcome)
  beta.hat.sd=estimate.coef.sd(beta=beta.hat,covariate=covariate,outcome=outcome)
  beta.hat.t = (beta.hat-0)/beta.hat.sd;
  return(beta.hat.t)
}

# The p-values using the Student t distribution are
2*apply(cbind(pt(-abs(beta.hat.t),df=n-2),(1-pt(abs(beta.hat.t),df=n-2))),1,min)

# We can also calculate the p-values using CLT
2*apply(cbind(pnorm(-abs(beta.hat.t)),(1-pnorm(abs(beta.hat.t)))),1,min)

# Or to use the bootstrap
# This part is NOT required
# It is recommended to read more about bootstrap if you want to use it for hypothesis
# since the empirical bootstrap procedure is not the most efficient one for this purpo

beta.hat.boot=replicate(1e5,boot.fit(covariate=x,outcome=y));
pval=numeric(length(beta.hat));
for(i in 1:length(beta.hat)){
  boot.est=beta.hat.boot[1,i,]
  pval[i]=2*min( mean(0<boot.est), mean(0>boot.est) )
}

## Wrap up the above code into one function
calculate.pvalue<-function(covariate,outcome,type){
  beta.hat.t=calculate.t(covariate,outcome);
  if (type=='t'){
    n=length(outcome);
    pval=2*apply(cbind(pt(beta.hat.t,df=n-2),(1-pt(beta.hat.t,df=n-2))),1,min);
  }else if (type=='z'){
    pval=2*apply(cbind(pnorm(-abs(beta.hat.t)),(1-pnorm(abs(beta.hat.t)))),1,min);
  }else if (type == 'bootstrap'){

    beta.hat=fit.linear.model(covariate=covariate,outcome=outcome)
    beta.hat.boot=replicate(1e5,boot.fit(covariate=covariate,outcome=outcome));
    pval=numeric(length(beta.hat));
    for(i in 1:length(beta.hat)){
      boot.est=beta.hat.boot[1,i,]
      pval[i]=2*min( mean(0<boot.est), mean(0>boot.est) )
    }
  }
}

```

```

    }
    return(pval)
}

```

Write a simulation to verify the equivalence between type I error rate and significance level

```

alpha=0.03; # significance level
simulate.one.instance<-function(x,beta.null,alpha,type){
  n=length(x);
  Ey= x*beta.null[2]+beta.null[1];
  error.terms= rnorm(n)*2;
  y=Ey+error.terms;
  pval=calculate.pvalue(covariate=x,outcome=y,type=type);
  rej.flag= pval[2]<alpha
  return(rej.flag)
}

N.sim=1e4;
set.seed(1)
n=50;
x=as.matrix(rnorm(n,mean=10,sd=2),ncol=n);
beta.null=c(beta.true[1],0)
sim.typeI=replicate(N.sim,simulate.one.instance(x,beta.null,alpha,type='t'));
mean(sim.typeI)

```

Assess the power

```

N.sim=1e4;
set.seed(1)
n=50;
x=as.matrix(rnorm(n,mean=10,sd=2),ncol=n);
powers<-numeric(3)
for(a in 1:3){
  beta.alt=c(beta.true[1],a/10)
  sim.rej=replicate(N.sim,simulate.one.instance(x,beta.alt,alpha,type='t'));
  powers[a]<-mean(sim.rej)
}
powers

```

Permutation test for testing $H_0: \beta_1 = 0$ (slope is zero)

```

permutation.test<-function(covariate,outcome){
  n=length(outcome);
  sample_indices = sample(1:n,n,replace=FALSE) # sampling without replacement
  covariate.perm= covariate[sample_indices]; outcome.perm= outcome[sample_indices];

```

```

    beta.hat.t=calculate.t(covariate.perm,outcome.perm)

    return(beta.hat.t[2])
}

set.seed(1)
n=50;
x=as.matrix(rnorm(n,mean=10,sd=2),ncol=n);
beta.true=c(20,0.15)
Ey= x*beta.true[2]+beta.true[1];
error.terms= rnorm(n)*2; # Normal errors, in order to verify the performance of the t-
# You can change it to other distributions
y=Ey+error.terms;

beta.hat.t.perm=replicate(1e5,permutation.test(covariate=x,outcome=y));
beta.hat.t=calculate.t(x,y)
2*min( mean(beta.hat.t[2]<beta.hat.t.perm),mean(beta.hat.t[2]>beta.hat.t.perm))

## F-test can be done using the anova() function in R
## It will be left as an exercise for you to write your version of F test
set.seed(1)
n=50;
x=as.matrix(rnorm(n,mean=10,sd=2),ncol=n);
beta.true=c(20,0.15)
Ey= x*beta.true[2]+beta.true[1];
error.terms= rnorm(n)*2; # Normal errors, in order to verify the performance of the t-
# You can change it to other distributions
y=Ey+error.terms;

full_model=lm(y~x+1);
reduced_model=lm(y~1);
anova(reduced_model,full_model)
summary(lm(y~x+1))

```

5.2 Multiple testing

```

# To discuss multiple testing, we generate 1000 outcomes, among which 100 are truly as

set.seed(1)
n=50;m=1000;
x=as.matrix(rnorm(n,mean=10,sd=2),ncol=n);
Y=matrix(0,nrow=n,ncol=m)

```

```

beta.true=c(20,0.6)
for(i in 1:m){
  Ey= x*beta.true[2]*(i<101)+beta.true[1];
  error.terms= rnorm(n)*2;
  Y[,i]=Ey+error.terms;
}

# Conduct the same test
pvalues=numeric(m)
for (i in 1:m){
  pval=calculate.pvalue(covariate=x,outcome=Y[,i],type='z');
  pvalues[i]<-pval[2];
}

```

```

## Bonferroni correction
alpha=0.05;
alpha.Bonf=alpha/m; # devided by the number of hypotheses
rej.flag=pvalues<alpha.Bonf;

```

```

## True positive
sum(rej.flag & ((1:m)<101))
## False positive
sum(rej.flag & ((1:m)>100))
## True negative
sum(!rej.flag & ((1:m)>100))

```

```

## False negative
sum(!rej.flag & ((1:m)<101))

```

```

## Benjamini-Hochberg procedure for FDR control

```

```

pvalues.sorted<-sort(pvalues,index.return=T)
alpha.BH<- (1:m)*alpha/m;
k.hat<-max(which(pvalues.sorted$x<alpha.BH))

```

```

rej.flag.BH=numeric(m);
rej.flag.BH[pvalues.sorted$ix[1:k.hat]]=1;

```

```

## True positive
sum(rej.flag.BH & ((1:m)<101))
## False positive
sum(rej.flag.BH & ((1:m)>100))
## True negative

```

```
sum(!rej.flag.BH & ((1:m)>100))

## False negative
sum(!rej.flag.BH & ((1:m)<101))

## False positive rate

sum(rej.flag.BH & ((1:m)>100))/sum(rej.flag.BH)

# Run a simulation to verify this
```

Chapter 6

Model diagnostics

Reading materials: Slides 93 - 101 in STA108_LinearRegression_S20.pdf.

Note: we will use existing functions in R for model diagnostics in this chapter. However, in your midterm report, you are still required to implement your own tools for model diagnostics.

6.1 Residual plot

```
## We can extract residual from an `lm()` object.  
  
fit.lm = lm(sales~TV+1,data=dat.advertising); # Fit the linear regression  
resid= fit.lm$residuals;  
  
plot(resid~dat.advertising$TV,pch=16,col='red',main='Residual plot')  
abline(h=0,lwd=3)
```

Why don't we draw the residual plot as residuals v.s. the response? This is because the two quantities are surely positively correlated, and thus is hard to extract any useful information from the plot.

```
plot(resid~dat.advertising$sales,pch=16,col='red',main='Wrong plot')  
abline(h=0,lwd=3)
```

How about drawing residuals against the fitted values? This is a good choice to detect certain anomalies.

```
par(mfrow=c(1,2))  
plot(resid~dat.advertising$TV,pch=16,col='red',main='Residual plot')  
abline(h=0,lwd=3)  
plot(resid~fit.lm$fitted.values,pch=16,col='red',main='Residual plot',xlab='Fitted value')  
abline(h=0,lwd=3)
```

```
par(mfrow=c(1,1))
```

6.2 Remedies for non-linearity

From the residual plot, we see that there exist some non-linearity between the residuals and the covariate, especially when the TV is number is small. This could suggest that the relationship between TV and sales is nonlinear. Typical form of nonlinearity takes form as $\exp(x)$, $x^{1/2}$, $\log(x)$, x^2 , etc. We may be able to guess the nonlinearity from the residual plot, or use model selection to pick the best nonlinear function, if there is not scientific knowledge on the relationship. In this example, we will fit two regression $y \sim \log(x)\beta_1 + \beta_0$ and $y \sim x^{1/2}\beta_1 + \beta_0$.

```
fit.log = lm(sales~log(TV)+1,data=dat.advertising); # Fit the linear regression
dat.advertising$sqrtTV=dat.advertising$TV^{1/2};
fit.sqrt = lm(sales~sqrtTV+1,data=dat.advertising); # Fit the linear regression
```

```
#par(mfrow=c(3,1))
```

```
plot(fit.log$residuals~log(dat.advertising$TV),pch=16,col='red',main='Residual plot (log)')
abline(h=0,lwd=3)
```

```
plot(fit.lm$residuals~dat.advertising$TV,pch=16,col='red',main='Residual plot (TV)',xlab='TV')
abline(h=0,lwd=3)
```

```
plot(fit.sqrt$residuals~dat.advertising$sqrtTV,pch=16,col='red',main='Residual plot (sqrtTV)')
abline(h=0,lwd=3)
```

```
#par(mfrow=c(1,1))
```

6.3 Independence

In the advertising data set, the residual plot does not tell us whether certain data points are correlated or not.

```
plot(fit.lm$residuals~dat.advertising$TV,pch=16,col='red',main='Residual plot')
abline(h=0,lwd=3)
```

We will consider a model where we know the errors are correlated, and show that they may still be hard to recognized in practice. We consider a model with autoregressive error, AR(1), which is widely used in time-series data analysis (e.g., financial data). We consider a model where, for $i = 1, 2, \dots, n$,

$$y_i = x_i\beta_1 + \beta_0 + \epsilon_i,$$

where $\epsilon_i = 3\epsilon_{i-1}/4 + z_i/4$ and $z_i \sim \mathcal{N}(0, 1)$. Here i represents some unit of time (e.g., months, years, days).


```

n=100;
set.seed(1);
z=rnorm(n);
epsilon= numeric(n);
for(i in 1:n){
  if (i==1){
    epsilon[i]=z[i];

  }else{
    epsilon[i]=3*epsilon[i-1]/4+z[i]/4;
  }
}
x=rnorm(n)*10+5;
y=x*2+1+z;

ARfits = lm(y~x+1);

par(mfrow=c(1,2))
plot(ARfits$residuals~x,pch=16,col='red',main='Residual plot 1')
abline(h=0,lwd=3)

plot(ARfits$residuals~c(1:n),pch=16,col='red',main='Residual plot 2',xlab="Index i")
abline(h=0,lwd=3)

par(mfrow=c(1,1))

```

6.4 Normality

It is also hard to see normality based on the residual plot. We can use a Quantile-Quantile plot to check if the errors are normally distributed. In the Q-Q plot, we draw the quantiles of residuals against the quantiles of the theoretical quantiles from a normal distribution. The $100(i/n)\%$ th quantile of the residuals is defined as the i th smallest residual.

Non-normal distributions in Q-Q plots.

```

n=500;
distributions=matrix(0,n,4)
distributions[,1] =-exp(rnorm(n));
distributions[,2] =exp(rnorm(n));
distributions[,3]=rt(n,df=3);
distributions[,4]=runif(n);

titles = c('Left skewed', 'Right skewed', 'Heavy-tailed', 'Light-tailed' )
par(mfrow=c(2,2))

```

```

for ( i in 1:4){
  this_dist=sort(distributions[,i]);
  normal_mean =mean(this_dist);normal_sd = sd(this_dist);
  this_dist=(this_dist - normal_mean)/normal_sd;
  normal_quantiles = qnorm( (1:length(this_dist))/length(this_dist));

  plot(this_dist~normal_quantiles,pch=16,col='red',main=titles[i],xlab='Normal quantiles',
  abline(a=0,b=1,lwd=3)
}

```

6.5 Homoscedasticity

We can see clearly that the residuals have wider spans when TV is larger, which suggested increasing variance. We can stabilize the variance by transforming the response variable $y^{1/2}$, $\log(y)$, etc. You can also use the Box-Cox transformation to find the most appropriate function.

```

fit.ylog=lm(log(sales)~TV+1,data=dat.advertising)
fit.ysqrt=lm(sqrt(sales)~TV+1,data=dat.advertising)
#par(mfrow=c(3,1))
plot(fit.lm$residuals~dat.advertising$TV,pch=16,col='red',main='Residual plot (original)',
abline(h=0,lwd=3)

plot(fit.ylog$residuals~dat.advertising$TV,pch=16,col='red',main='Residual plot (log)',y
abline(h=0,lwd=3)

plot(fit.ysqrt$residuals~dat.advertising$TV,pch=16,col='red',main='Residual plot (sqrt)',
abline(h=0,lwd=3)

```

6.6 Influential Observations and Outliers

The influential observations are samples that has a large leverage. You should search for a formal definition of leverage if you are interested. In plain words, the influential observations are data points that live far away from others in terms of their values of the covariates.

Outliers are the observations whose responses are far away from observations with similar covariates. We can see these from the usual scatter plots. We will not cover formal testing or measures for influential observations and outliers in this class.

```

influential_cov=dat.advertising$TV;influential_cov[1]=500;
outliers_resp=fit.lm$residuals;outliers_resp[1]=30;

par(mfrow=c(1,2))
plot(fit.lm$residuals~influential_cov,pch=16,col='red',main='Influential obs.',xlab='TV'
abline(h=0,lwd=3)

```

```
points(y=fit.lm$residuals[1],x=influential_cov[1],col='green',pch=16)

plot(outliers_resp~dat.advertising$TV,pch=16,col='red',main='Outlier',xlab='TV')
abline(h=0,lwd=3)
points(y=outliers_resp[1],x=dat.advertising$TV[1],col='green',pch=16)
```


Chapter 7

Multiple covariates

7.1 Examples

Reading materials: Slides 103 - 111 in STA108_LinearRegression_S20.pdf.

We will revisit the examples that have been studied in Chapter ??.

7.1.1 Advertising data

We now consider all the covariates in the data set. The resulting model is

$$y_i = \beta_0 + \sum_{j=1}^3 x_{i,j} \beta_j + \epsilon_i, i = 1, \dots, 200,$$

where y_i is the **sales** (1 unit = 1000 dollars) for the i th entry, $x_{i,1}$ is the TV advertising budget for the i th entry, $x_{i,2}$ is the radio advertising budget for the i th entry, $x_{i,3}$ is the newspaper advertising budget for the i th entry. In addition, we assume that the errors $\{\epsilon_i\}_{i=1}^{200}$ satisfy that $\epsilon_1, \dots, \epsilon_{200}$ are independently and identically distributed (i.i.d.), $\mathbb{E}[\epsilon_i] = 0$ for $i = 1, 2, \dots, 200$ and $\text{var}(\epsilon_i) = \sigma^2$ for $i = 1, 2, \dots, 200$. Recall that we consider fixed design (i.e., x_i is not random) in this course for simplicity.

```
dat.advertising=read.csv('./data/advertising.csv');  
# Fit a multiple linear regression  
fit.advertising = lm(sales~TV+radio+newspaper+1,data=dat.advertising);  
# Summarize the fitted results  
summary(fit.advertising)  
  
fit.advertising.slr = lm(sales~TV+1,data=dat.advertising);  
# Summarize the fitted results  
summary(fit.advertising.slr)
```

How would you interpret the fitted results? How do they compare with the results from the simple linear regression in Chapter ???

7.1.2 Flu shot

We will include `age` and `female`(gender) in the regression in addition to the treatment received. What can you conclude from the fitted results?

```
dat.flu = read.table("./data/flu240.txt", header = TRUE)
# Fit a multiple linear regression
fit.flu= lm(outcome~treatment.received+age+female+1,data=dat.flu);
# Summarize the fitted results
summary(fit.flu)
```

In addition, we can fit a regression with treatment received as the outcome, where treatment assigned, age and gender are covariates. What can you conclude now?

```
fit.trt= lm(treatment.received~treatment.assigned+age+female+1,data=dat.flu);
# Summarize the fitted results
summary(fit.trt)
```

7.1.3 Project STAR

Project STAR is a stratified randomized experiment, where randomization happened within each school. Hence, we will add `schoolid2` as an additional factor in the regression.

```
library(AER)
data("STAR")
dat.STAR=STAR; # Just to be consistent
# Fit a simple linear regression
fit.STAR= lm(math2~as.factor(star2)+as.factor(schoolid2)+1,data=dat.STAR);
# Summarize the fitted results
summary(fit.STAR)$coef[1:3,]

# We can no longer draw a fitted line here...
```

7.2 Classification of variables

Reading materials: Slides 112 - 118 in STA108_LinearRegression_S20.pdf.

We will use simulation to examine the roles of the variables.

```
source('sources.r')

simulate.one.instance.data<-function(){
n=100;
confounder = rnorm(n)*2;
precision = runif(n);
instrument = rbinom(n,size=1,prob=0.4);
trt=rnorm(n) + confounder*0.4-instrument;
```

```

y=trt*1-0.5*confounder+precision*0.7+0.5*rnorm(n);
out=data.frame(y=y,trt=trt,confounder=confounder, precision=precision, instrument=instru
return(out)
}

```

```
## Effect if ignoring confounder:
```

```

simulate.one.instance<-function(){
  dat=simulate.one.instance.data();
  beta.hat=fit.linear.model(covariate=dat[, 'trt'], outcome=dat$y)

  beta.hat.confounder=fit.linear.model(covariate=as.matrix(dat[, c('trt', 'confounder')]))
  return(c(beta.hat[2], beta.hat.confounder[2]))
}

```

```
N.sim=10000;
```

```
set.seed(1)
```

```

sim.confounder=replicate(N.sim, simulate.one.instance());
apply(sim.confounder, 1, mean) # The true value is 1
apply( (sim.confounder-1), 1, function(x){ sum(x^2) }) # The mean squared error

```

```
## Effect if ignoring precision variable:
```

```

simulate.one.instance<-function(){
  dat=simulate.one.instance.data();
  beta.hat=fit.linear.model(covariate=as.matrix(dat[, c('trt', 'confounder')]), outcome=dat$y)

  beta.hat.precision=fit.linear.model(covariate=as.matrix(dat[, c('trt', 'confounder', 'precision')]))
  return(c(beta.hat[2], beta.hat.precision[2]))
}

```

```
N.sim=10000;
```

```
set.seed(1)
```

```

sim.precision=replicate(N.sim, simulate.one.instance());
apply(sim.precision, 1, mean) # The true value is 1
apply( (sim.precision-1), 1, function(x){ sum(x^2) }) # The mean squared error

```

```
## The power of instruments
```

```
### If the confounder is unobserved, we can use the two stage least squares method
```

```

simulate.one.instance<-function(){
  dat=simulate.one.instance.data();
  beta.hat.naive=fit.linear.model(covariate=as.matrix(dat[, 'trt']), outcome=dat$y)

```

```

yiv=fit.linear.model(covariate=as.matrix(dat[,c('instrument')]),outcome=dat$y)

trtiv=fit.linear.model(covariate=as.matrix(dat[,c('instrument')]),outcome=dat$trt)

  return(c(beta.hat.naive[2],yiv[2]/trtiv[2]))
}
N.sim=10000;

set.seed(1)

sim.iv=replicate(N.sim, simulate.one.instance());
apply(sim.iv,1,mean) # The true value is 1
apply( (sim.iv-1),1,function(x){ sum(x^2) }) # The mean squared error

```

7.3 Least squares estimation

Reading materials: Slides 119 - 141 in STA108_LinearRegression_S20.pdf.

```

# Consider another data set with a bit more covariates
library(AER)
data("Fatalities")
dat.fatalities = Fatalities
# This dataset contains 34 variables
# It is actually a longitudinal/panel data
# We will ignore the spatial temporal structure in this class

y=dat.fatalities$fatal; #Number of vehicle fatalities.
X=as.matrix(dat.fatalities[,c('spirits','unemp','income','beertax','baptist','mormon','c

# Rewrite the functions in Chapter 2
linear.model<-function(beta,covariate){
  # beta: a 2-vector, where the first entry is the intercept
  beta=as.matrix(beta,nrow=length(beta))
  yout=covariate%*%beta[-1]+beta[1];
  return(yout);
  # Note: this function only works with simple linear regression model
  # How would you generalize it?
}

sum.of.squares<-function(beta,covariate,outcome){
  yout=linear.model(beta=beta,covariate=covariate);
  res=outcome-yout;

```



```

    sos= sum(res^2);
    return(sos)
}
fit.linear.model<-function(covariate,outcome){
  X=cbind(1,covariate);

  beta.fit=solve( t(X)%*%X )%*%t(X)%*%outcome;
  return(beta.fit)
}

beta.hat=fit.linear.model(covariate=X,outcome=y)

# Compare with lm()
fit.lm=lm(y~X+1);
beta.hat

# Hat matrix
calculate.hat<-function(covariate){
  X=cbind(1,covariate);
  P=X%*%solve(t(X)%*%X)%*%t(X);
  return(P)
}
P=calculate.hat(covariate=X);

# Verify the properties of the hat matrix:
## Residuals
resid= y-linear.model(beta=beta.hat,covariate=X);
resid.hat = (diag(length(y))-P)%*%y;
sum((resid-resid.hat)^2)

## Idempotent
max( abs(P%*%P-P))

## Orthogonality
max(abs((diag(length(y))-P)%*%X))

# Sum of squares with the hat matrix:
n=length(y);J=matrix(1,n,n);I=diag(n)
total.sum.of.squares= t(y)%*%(I-J/n)%*%y
explained.sum.of.squares=t(y)%*%(P-J/n)%*%y
residual.sum.of.squares=t(y)%*%(I-P)%*%y

# Check if it is true
sum.of.squares(beta.hat,covariate = X,outcome = y)

```

```

# R2
R.sq = explained.sum.of.squares/total.sum.of.squares;

# Through more variables into X
# These variables are meaningless
set.seed(1)
Z=matrix(rnorm(5*n),nrow=n,ncol=5);
X.Z=cbind(X,Z);

P.Z=calculate.hat(covariate=X.Z);

R.sq.Z= t(y)%*(P.Z-J/n)%*y/t(y)%*(I-J/n)%*y;

R.sq.Z
R.sq

# Hence the adjusted R2:
R.sq.adj=1-(residual.sum.of.squares/(n-dim(X)[2]-1))/(total.sum.of.squares/(n-1));

R.sq.Z.adj=1-(t(y)%*(I-P.Z)%*y/(n-dim(X.Z)[2]-1))/(t(y)%*(I-J/n)%*y/(n-1));

```

7.4 Underfitting and overfitting

Reading materials: Slides 142 - 152 in STA108_LinearRegression_S20.pdf.

```

## Underfitting
simulate.one.instance<-function(beta.true,X,Z){
  eta=c(3,4);
  n= dim(X)[1];
  noise = rnorm(n)/2;
  y.XZ=X%*%beta.true+Z%*%eta + noise;

  beta.full = fit.linear.model(covariate=cbind(X,Z),outcome=y.XZ)
  beta.under= fit.linear.model(covariate=X,outcome=y.XZ)

  return(cbind(beta.full[2:3]-beta.true,beta.under[2:3]-beta.true))
}

beta.true=c(1,2)
n=100;
X=matrix(rnorm(2*n),ncol=2)
Z=matrix(rnorm(2*n),ncol=2)

set.seed(1);

```

```

N.sim=10000;
under.sim=replicate(N.sim,simulate.one.instance(beta.true,X,Z))
# Fitting the true model gives

apply(under.sim[1,,],1,mean) # Non-zero bias

# Does it contradict the claim about precision variable?
# Try another simulation:
simulate.one.instance.pre<-function(beta.true){
  eta=c(3,4);
  n=100;
  X=matrix(rnorm(2*n),ncol=2)
  Z=matrix(rnorm(2*n),ncol=2)
  noise = rnorm(n)/2;
  y.XZ=X%%beta.true+Z%%eta + noise;

  beta.full = fit.linear.model(covariate=cbind(X,Z),outcome=y.XZ)
  beta.under= fit.linear.model(covariate=X,outcome=y.XZ)

  return(cbind(beta.full[2:3]-beta.true,beta.under[2:3]-beta.true))
}

beta.true=c(1,2)

set.seed(1);
N.sim=10000;
pre.sim=replicate(N.sim,simulate.one.instance.pre(beta.true))

apply(pre.sim[1,,],1,mean) # almost zero bias

### Overfitting
# Just modify the code slightly..
simulate.one.instance<-function(beta.true,X,Z){
  eta=c(3,4);
  n= dim(X)[1];
  noise = rnorm(n)/2;
  y.X=X%%beta.true+ noise;

  beta.over= fit.linear.model(covariate=cbind(X,Z),outcome=y.X)
  beta.full= fit.linear.model(covariate=X,outcome=y.X)

  return(cbind(beta.full[2:3]-beta.true,beta.over[2:3]-beta.true))
}

```

```

beta.true=c(1,2)
n=100;
X=matrix(rnorm(2*n),ncol=2)
Z=matrix(rnorm(2*n),ncol=2)

set.seed(1);
N.sim=10000;
over.sim=replicate(N.sim,simulate.one.instance(beta.true,X,Z))
# Fitting the true model gives

apply(over.sim[1,,],1,mean) # No bias

```

7.5 Sampling distribution and inference

Reading materials: Slides 153 - 165 in STA108_LinearRegression_S20.pdf.

We will skip the part for multivariate normal distribution. You will learn more about this in STA 135: Multivariate Data Analysis.

7.5.1 ANOVA

```

attach(dat.advertising)
full.model=lm(sales~TV+radio+newspaper);
reduced.model=lm(sales~TV);
anova(reduced.model,full.model)

RSS.Reduced<-sum((summary(reduced.model)$residuals)^2);
RSS.Full<-sum((summary(full.model)$residuals)^2);
df.Reduced <- summary(reduced.model)$df[2];
df.Full <- summary(full.model)$df[2];

tF<-((RSS.Reduced-RSS.Full)/(df.Reduced-df.Full))/(RSS.Full/df.Full);
tF

```

7.5.2 Wald test

Suppose that we want to test the null hypothesis

$$H_0 : \beta_1 - 2\beta_2 = 0 \text{ v.s. } H_a : \beta_1 - 2\beta_2 \neq 0,$$

where β_1 is the regression coefficient for TV and β_2 is the regression coefficient for radio.

```

full.model=lm(sales~TV+radio+newspaper);
covariance.LSE=vcov(full.model); # Calculate the covariance for the least squares estim

```

```

coef.LSE=coef(full.model);
covariance.test = covariance.LSE[2:3,2:3]; # Extract the submatrix of covariance for be
coef.test=coef.LSE[2:3];
R=matrix(0,1,2);
R[1,1]=1;R[1,2]=-2;

test.stat=t(R%%coef.test)%%solve(R%%covariance.test%%t(R))%%(R%%coef.test);
p.value=1-pchisq(test.stat,df=1);

```

We can verify that the Wald test yields similar results with the F-test using `anova`.

$$H_0 : \beta_1 = \beta_2 = 0 \text{ v.s. } H_a : \beta_1 \neq 0, \beta_2 \neq 0.$$

```

full.model=lm(sales~TV+radio+newspaper);
covariance.LSE=vcov(full.model); # Calculate the covariance for the least squares estim
coef.LSE=coef(full.model);
covariance.test = covariance.LSE[2:3,2:3]; # Extract the submatrix of covariance for be
coef.test=coef.LSE[2:3];
R=matrix(0,2,2);
R[1,1]=1;R[2,2]=1;

test.stat=t(R%%coef.test)%%solve(R%%covariance.test%%t(R))%%(R%%coef.test);
p_value=1-pchisq(test.stat,df=2);

reduced_model=lm(sales~newspaper);
anova(reduced.model,full.model)

```


Chapter 8

Model selection

8.1 Criteria

Reading materials: Slides 166 - 181 in STA108_LinearRegression_S20.pdf.

8.2 Selection procedure

Reading materials: Slides 182 - 189 in STA108_LinearRegression_S20.pdf.

Appendix A

R basics

Reading materials: Sections 4, 6, and 8 in R for data science by Garrett Golemund and Hadley Wickham (optional).

This section reviews some basic functions in R that will be useful in this class.

A.1 Summary statistics

A.2 Data structures

A.2.1 String

A.2.2 Data frame

A.3 List

A.4 Functions in R

A.5 Miscellaneous

Appendix B

Linear algebra

Reading materials: STA108_LinearAlgebra_S20.pdf

B.1 Vector

B.2 Matrix

B.3 Other operations on vectors and matrices

B.3.1 Array

Appendix C

Simulation and visualization

Reading materials: Chapters 3, 5, and 7 in R for data science by Garrett Golemund and Hadley Wickham (optional).

C.1 Simulation and visualization: univariate

C.2 Simulation and visualization: multivariate