

36350-A HW9

Joong Ho Choi

TOTAL POINTS

108 / 120

QUESTION 1

1 Q1 20 / 20

- ✓ - **0 pts** Correct
- **2 pts** does not sample 1000 values
- **6 pts** incorrect setup for bivariate normal
- **4 pts** no histogram
- **2 pts** no vertical line for population
- **4 pts** question links incorrectly
- **2 pts** question partially linked
- **20 pts** blank

QUESTION 2

2 Q2 20 / 20

- ✓ - **0 pts** Correct
- **2 pts** does not sample 1000 values
- **6 pts** incorrect setup for trivariate normal
- **4 pts** no plot using ggplot
- **4 pts** incorrect sample means
- **4 pts** question links incorrectly
- **20 pts** blank

QUESTION 3

3 Q3 8 / 20

- **0 pts** Correct
- ✓ - **4 pts** incorrect Sigma.cond setup
- ✓ - **4 pts** incorrect mu.cond setup
- **2 pts** not setting $x_2 = 1$
- **3 pts** no plot using ggplot
- ✓ - **2 pts** Missing/incorrect conditional correlation coefficient
- ✓ - **2 pts** Missing/incorrect sample correlation matrix.
- **1 pts** Wrong plot (should be a scatterplot)
- **4 pts** question links incorrectly
- **20 pts** blank

QUESTION 4

4 Q4 20 / 20

- ✓ - **0 pts** Correct
- **4 pts** incorrect setup for bivariate normal model 1
- **4 pts** incorrect setup for bivariate normal model 2
- **2 pts** does not sample 100 samples from each model
- **3 pts** no plot with different color
- **5 pts** not using glm for classification
- **4 pts** missing/incorrect % of misclassification
- **4 pts** question links incorrectly
- **20 pts** blank

QUESTION 5

5 Q5 20 / 20

- ✓ - **0 pts** Correct
- **5 pts** incorrect # of PCs to retain
- **4 pts** missing/incorrect proportion of variance explained by PCs
- **6 pts** missing/incorrect conclusion of original variables' importance
- **4 pts** question links incorrectly
- **20 pts** blank

QUESTION 6

6 Q6 20 / 20

- ✓ - **0 pts** Correct
- **2 pts** incorrect setup for lm.orig
- **2 pts** missing/incorrect R^2 for lm.orig
- **2 pts** missing/incorrect SSE for lm.orig
- **4 pts** incorrect/missing setup for lm.pc
- **2 pts** missing/incorrect R^2 for lm.pc
- **2 pts** missing/incorrect SSE for lm.pc
- **4 pts** missing conclusion about whether R^2 and SSE are close
- **4 pts** question links incorrectly
- **20 pts** blank

QUESTION 7

7 Late Penalty 0 / 0

✓ - 0 pts Correct

HW: Week 9

36-350 – Statistical Computing

Week 9 – Spring 2021

Name: Joong Ho Choi

Andrew ID: joonghoc

You must submit **your own** HW as a PDF file on Gradescope.

HW Length Cap Instructions

- If the question requires you to print a data frame in your solution e.g. `q1_out_df`, you must first apply `head(q1_out_df, 30)` and `dim(q1_out_df)` in the final knitted pdf output for such a data frame.
- Please note that this only applies if you are knitting the Rmd to a pdf, for Gradescope submission purposes.
- If you are using the data frame output for visualization purposes (for example), use the entire data frame in your exploration
- The **maximum allowable length** of knitted pdf HW submission is **30 pages**. Submissions exceeding this length *will not be graded* by the TAs. All pages must be tagged as usual for the required questions per the usual policy
- For any concerns about HW length for submission, please reach out on Piazza during office hours

Question 1

(20 points)

Display the sampling distribution for $R_{1,2}$, the off-diagonal element of a two-dimensional sample correlation matrix for a bivariate normal, given that $\mu_1 = \mu_2 = 2$, $\sigma_1 = 1$, $\sigma_2 = 2$, $\rho_{1,2} = -0.5$, and $n = 100$. Sample 1000 values and display them in a histogram; include a vertical line for the population value (-0.5). (Reminder: if I don't specify how exactly to visualize something, i.e., if I don't specify base R versus `ggplot`, then you can choose whichever.) Note that the final distribution that you see will not be normal.

```
library(emdbook)
suppressMessages(library(MASS))
library(ggplot2)
mu=c(2,2)
sigma.1=1
sigma.2=2
rho.12=-0.5
n=100
Sigma=matrix(c(sigma.1^2,rho.12*sigma.1*sigma.2,rho.12*sigma.1*sigma.2,sigma.2^2),nrow=2)
set.seed(101)
```

```

res=c()
for (ii in 1:1000){
  data=mvrnorm(n,mu,Sigma)
  data1=cor(data)
  j=data1[1,2]
  res=c(res,j)
}
length(res) #1000

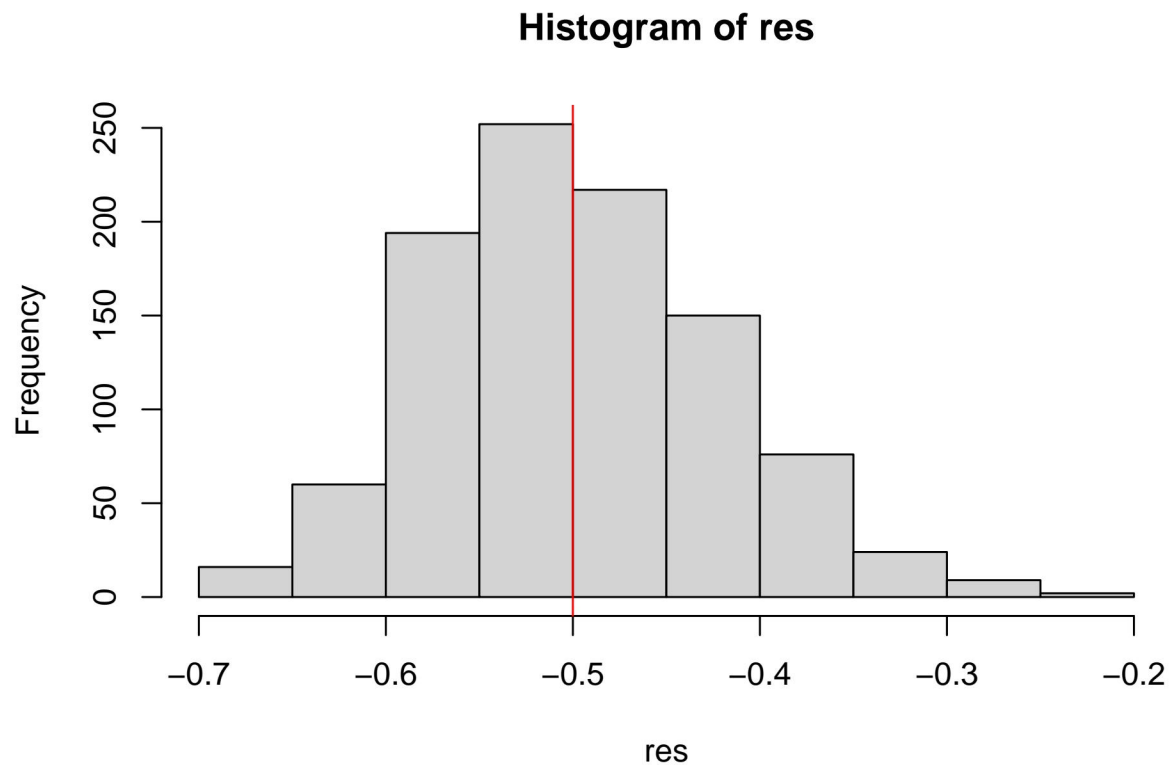
```

```
## [1] 1000
```

```

hist(res)
abline(v = -0.5 , col="red")

```



Question 2

(20 points)

Assume you have a trivariate multivariate normal with $\mu = \{2, 3, 4\}$ and $\sigma = \{1, 2, 1\}$, and $\rho_{1,2} = 0.4$, $\rho_{1,3} = 0.7$, and $\rho_{2,3} = -0.2$. Sample 1000 data from the marginal distribution for (x_1, x_3) , and display them via `ggplot`. Compute the sample means along each marginal axis: they should be approximately 2 and 4.

1 Q1 20 / 20

✓ - 0 pts Correct

- 2 pts does not sample 1000 values
- 6 pts incorrect setup for bivariate normal
- 4 pts no histogram
- 2 pts no vertical line for population
- 4 pts question links incorrectly
- 2 pts question partially linked
- 20 pts blank

```

res=c()
for (ii in 1:1000){
  data=mvrnorm(n,mu,Sigma)
  data1=cor(data)
  j=data1[1,2]
  res=c(res,j)
}
length(res) #1000

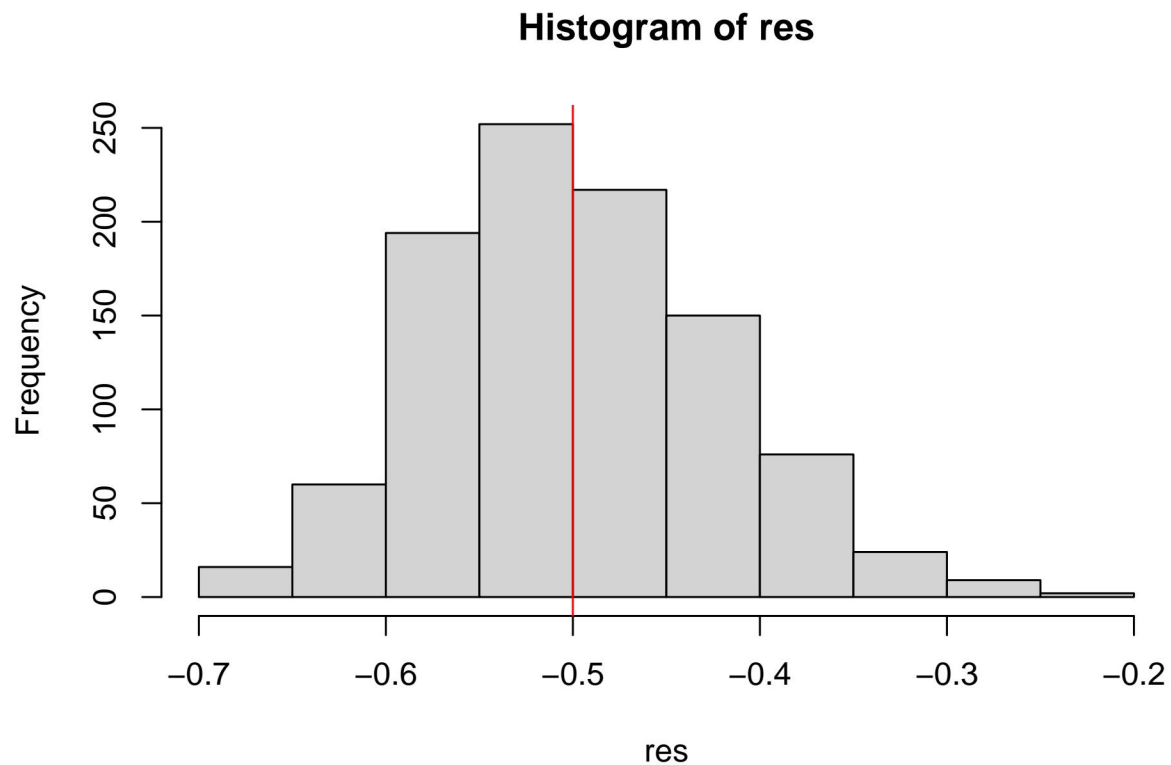
```

```
## [1] 1000
```

```

hist(res)
abline(v = -0.5 , col="red")

```



Question 2

(20 points)

Assume you have a trivariate multivariate normal with $\mu = \{2, 3, 4\}$ and $\sigma = \{1, 2, 1\}$, and $\rho_{1,2} = 0.4$, $\rho_{1,3} = 0.7$, and $\rho_{2,3} = -0.2$. Sample 1000 data from the marginal distribution for (x_1, x_3) , and display them via `ggplot`. Compute the sample means along each marginal axis: they should be approximately 2 and 4.

```

mu=c(2,3,4)
sigma.1=1
sigma.2=2
sigma.3=1

rho.12=-0.4
rho.13=0.7
rho.23=-0.2

Sigma=matrix(c(sigma.1^2,rho.12*sigma.1*sigma.2,rho.13*sigma.1*sigma.3,
               rho.12*sigma.1*sigma.2,sigma.2^2,rho.23*sigma.2*sigma.3,
               rho.13*sigma.1*sigma.3,rho.23*sigma.2*sigma.3,sigma.3^2),nrow=3)

p1=c(1)
p3=c(3)
res1=c()

x=seq(-1,6,by=0.01)

(mu.marg=mu[c(1,3)])

```

```
## [1] 2 4
```

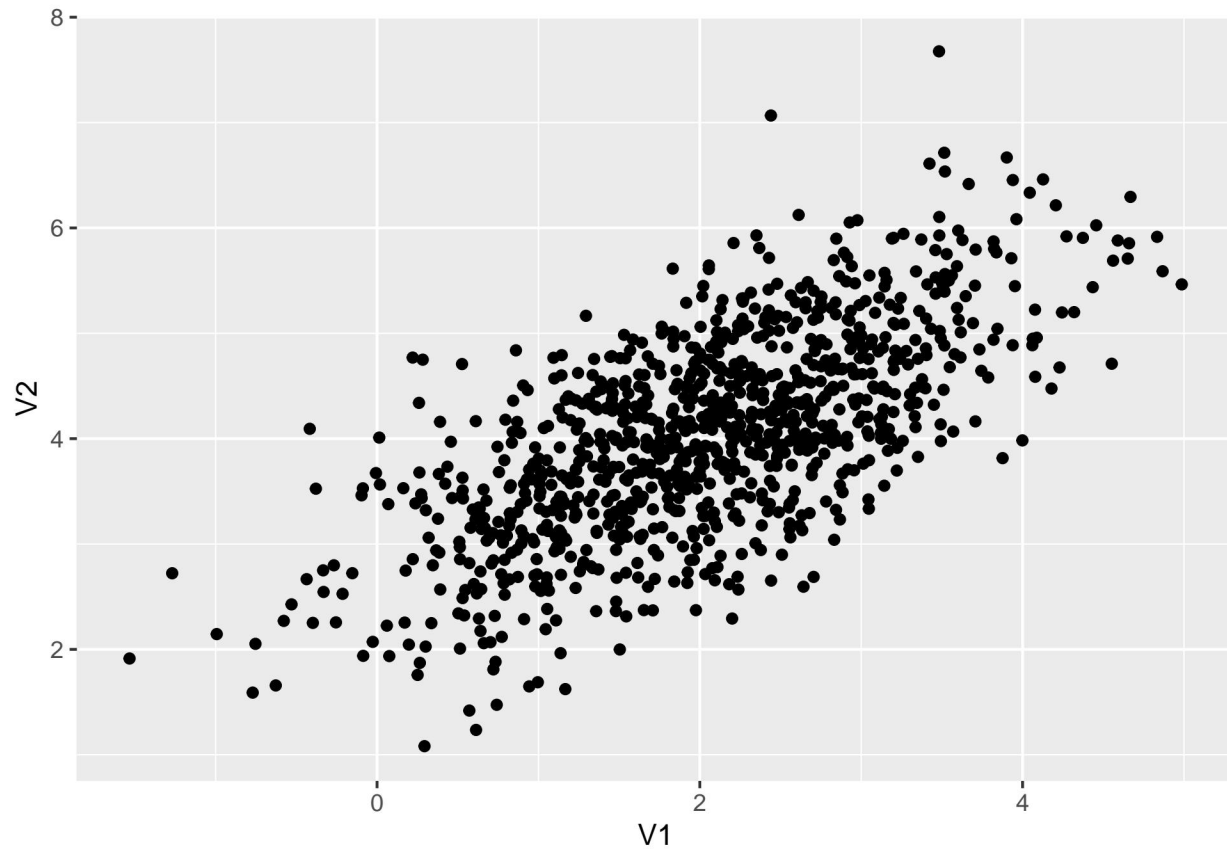
```

Sigma.marg=Sigma[c(1,3),c(1,3)]

df1<-as.data.frame(matrix(nrow=1000,ncol=2))
set.seed(123)
for (ii in 1:1000){
  fx=mvrnorm(1,mu.marg,Sigma.marg)
  df1[ii,1]=fx[1]
  df1[ii,2]=fx[2]
}

ggplot(df1,aes(x=V1,y=V2))+geom_point()

```



Question 3

(20 points)

Repeat Q2, except here you should sample from the conditional distribution $f(x_1, x_3 | x_2 = 1)$. Compute the conditional correlation coefficient ρ_{cond} between the data along the x_1 and x_3 axes, given $x_2 = 1$, and display the sample correlation matrix.

```
#conditional correlation coefficient part
x.2=1 #condition: x_2=1
k=c(1,3)
d.minus.k=c(2)
Sigma.kk=Sigma[k,k]
Sigma.kd=Sigma[k,d.minus.k]
Sigma.dk=Sigma[d.minus.k,k]
Sigma.dd=Sigma[d.minus.k,d.minus.k]

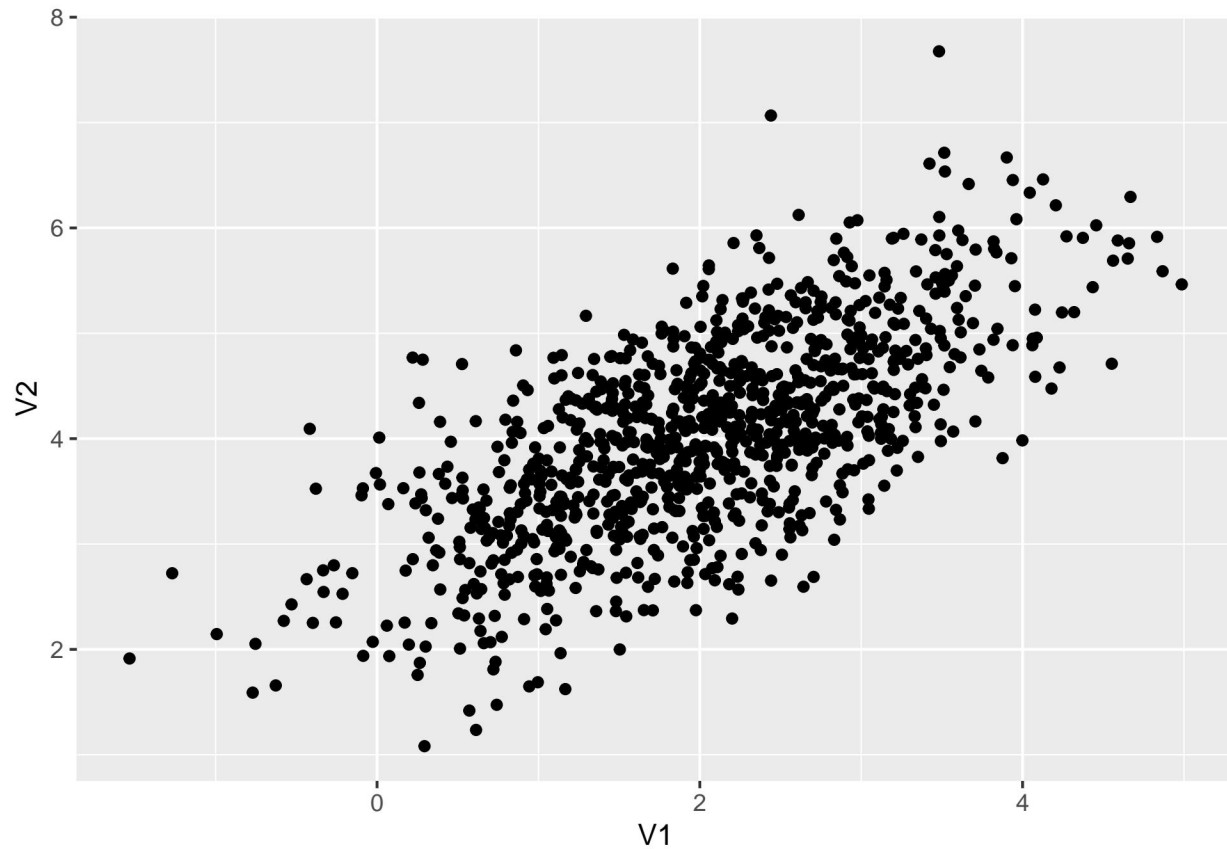
mu.cond=mu[c(1,3)]+Sigma.kd %*% solve(Sigma.dd) %*% matrix(c(1)-mu[2])
mu.cond

##      [,1]
## [1,]  2.4
## [2,]  4.2
```


2 Q2 20 / 20

✓ - 0 pts Correct

- 2 pts does not sample 1000 values
- 6 pts incorrect setup for trivariate normal
- 4 pts no plot using ggplot
- 4 pts incorrect sample means
- 4 pts question links incorrectly
- 20 pts blank



Question 3

(20 points)

Repeat Q2, except here you should sample from the conditional distribution $f(x_1, x_3 | x_2 = 1)$. Compute the conditional correlation coefficient ρ_{cond} between the data along the x_1 and x_3 axes, given $x_2 = 1$, and display the sample correlation matrix.

```
#conditional correlation coefficient part
x.2=1 #condition: x_2=1
k=c(1,3)
d.minus.k=c(2)
Sigma.kk=Sigma[k,k]
Sigma.kd=Sigma[k,d.minus.k]
Sigma.dk=Sigma[d.minus.k,k]
Sigma.dd=Sigma[d.minus.k,d.minus.k]

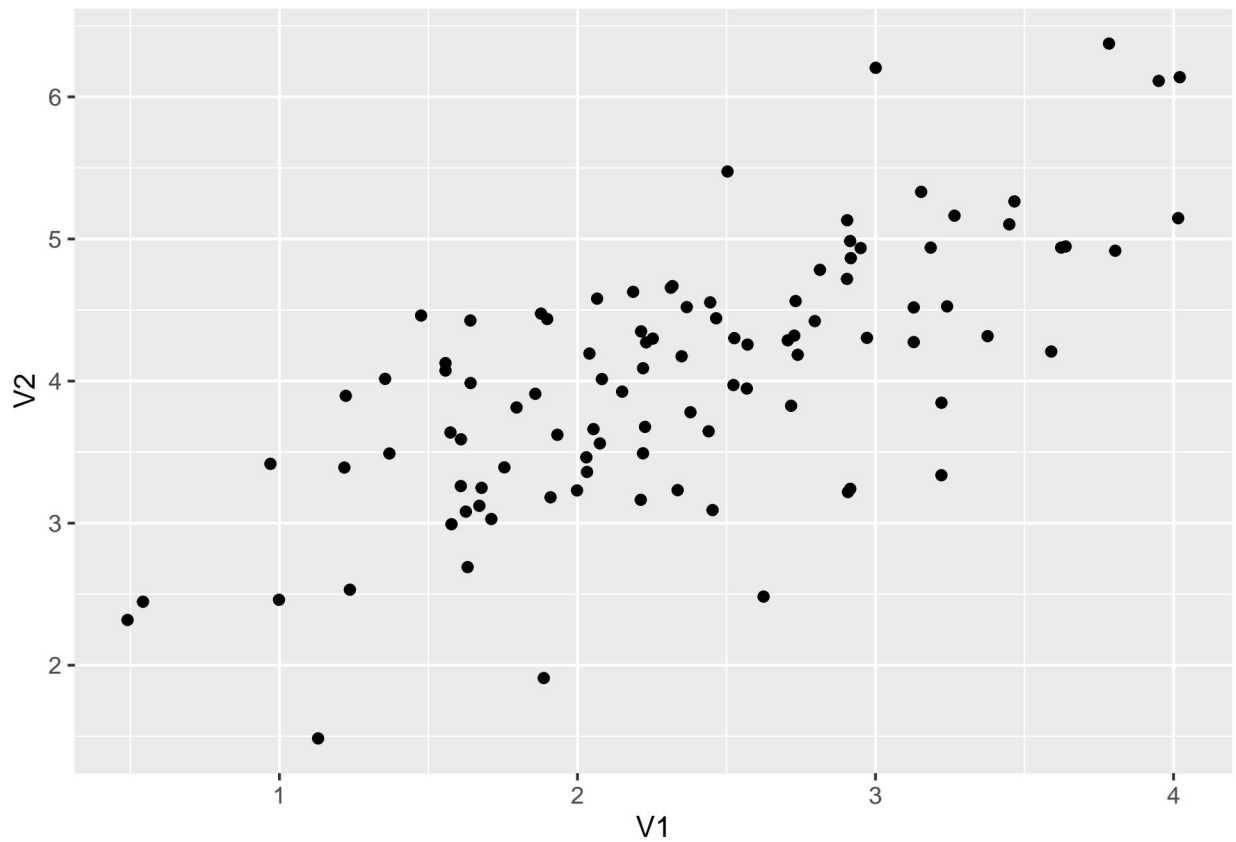
mu.cond=mu[c(1,3)]+Sigma.kd %*% solve(Sigma.dd) %*% matrix(c(1)-mu[2])
mu.cond

##      [,1]
## [1,]  2.4
## [2,]  4.2
```

```
Sigma.cond=Sigma.kk-Sigma.kd %*% solve(Sigma.dd) %*% Sigma.dk
Sigma.cond
```

```
##      [,1] [,2]
## [1,] 0.84 0.62
## [2,] 0.62 0.96
```

```
data<-mvrnorm(100,mu.cond,Sigma.cond)
df2<-as.data.frame(data)
ggplot(df2,aes(x=V1,y=V2))+geom_point()
```



```
res<-cor(data)
res  #sample correlation matrix
```

```
##      [,1] [,2]
## [1,] 1.0000000 0.7002904
## [2,] 0.7002904 1.0000000
```

```
Sigma.cond[2]/sqrt(Sigma.cond[1]*Sigma.cond[4]) #conditional correlation coefficient
```

```
## [1] 0.6904249
```

3 Q3 8 / 20

- 0 pts Correct
- ✓ - 4 pts incorrect Sigma.cond setup
- ✓ - 4 pts incorrect mu.cond setup
 - 2 pts not setting $x_2 = 1$
 - 3 pts no plot using ggplot
- ✓ - 2 pts Missing/incorrect conditional correlation coefficient
- ✓ - 2 pts Missing/incorrect sample correlation matrix.
 - 1 pts Wrong plot (should be a scatterplot)
 - 4 pts question links incorrectly
 - 20 pts blank

Question 4

(20 points)

Assume that you have a mixture model: you have 100 data sampled from a bivariate normal with $\mu = \{1, 1\}$ and $\sigma = \{1.2, 1.2\}$, with $\rho = 0.4$, and another 100 data sampled from a bivariate normal with $\mu = \{3, 3\}$, $\sigma = \{1, 1\}$, and $\rho = -0.6$.

Plot your sampled data with separate colors for each component of the mixture.

Then perform logistic regression to try to classify each sampled point as being from component 1 or component 2, and output the proportion of times you misclassify a point. (Don't worry about breaking your data up into training and testing sets, as this is a simple academic exercise; just use all 200 points to train your classifier, then output the training misclassification error.)

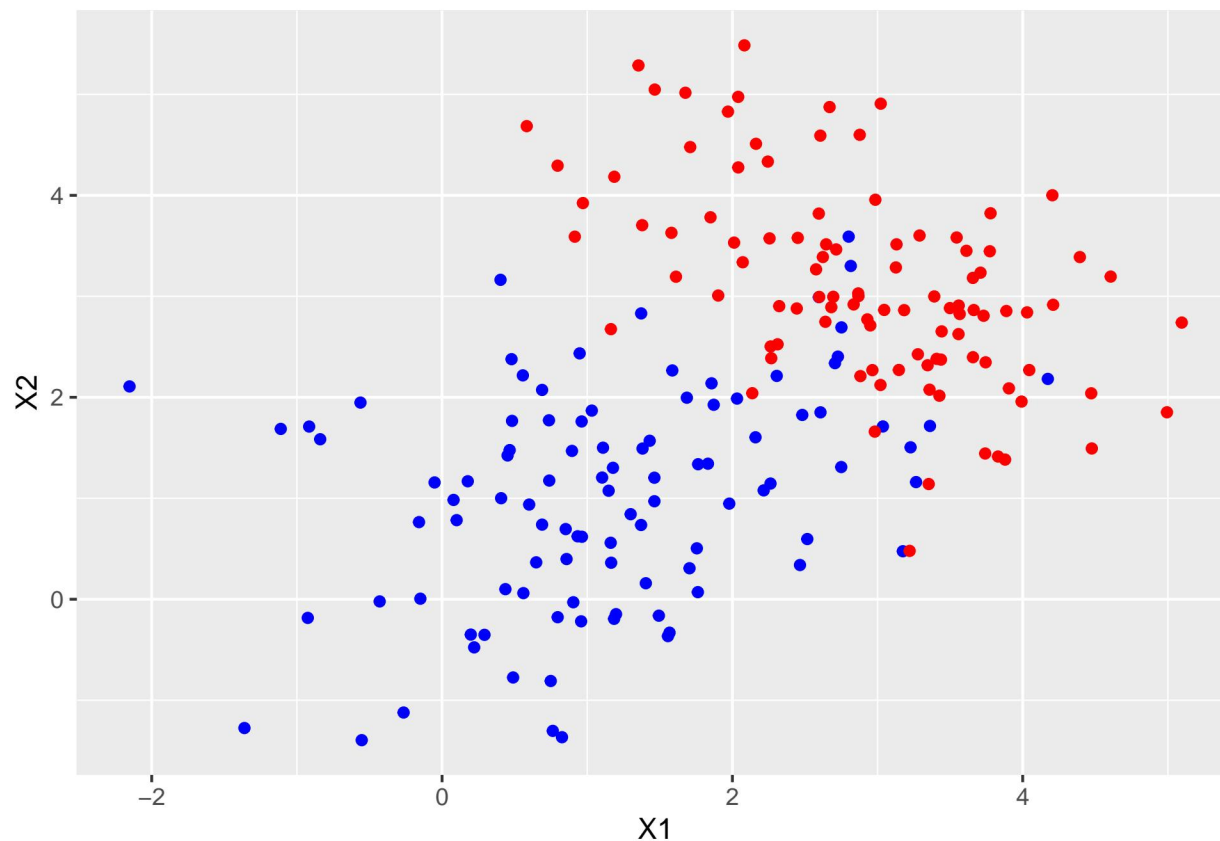
How to train your logistic classifier and get the misclassification rate?

- Assuming you already have a data frame with sampled x_1 and x_2 values in the first and second columns, add a third column with the labeled class. (Name this column `class`.) Use 0 for the first class, and 1 for the second.
- Use `glm()` with model formula `class~.`, your data frame, and the argument `family=binomial`.
- Use `predict()` with the output of `glm()` and with the argument `type="response"`. This will generate 200 predictions between 0 and 1.
- Round off all predictions to 0 or 1.
- Create a `table()` with the arguments being your rounded-off predictions and the labeled classes.
- Compute the proportion of table elements that are “off-diagonal” (upper-right and lower-left). Done.

```
set.seed(123)
mu1=c(1,1)
Sigma1=matrix(c(1.2^2,0.4*1.2*1.2,0.4*1.2*1.2,1.2^2),nrow=2)
data1=mvrnorm(100,mu1,Sigma1)
df1<-data.frame(data1)

mu2=c(3,3)
Sigma2=matrix(c(1^2,-0.6*1*1,-0.6*1*1,1^2),nrow=2)
data2=mvrnorm(100,mu2,Sigma2)
df2<-data.frame(data2)

ggplot()+geom_point(data=df1,aes(x=X1,y=X2),color="blue")+geom_point(data=df2,aes(x=X1,y=X2),color="red")
```



```
df1$class=0
df2$class=1

df3<-rbind(df1,df2)

model<-glm(class~.,data=df3,family="binomial")
pred<-predict(model,type="response")

pred_rounded<-ifelse(pred>=0.5,1,0)
k<-table(pred_rounded,df3$class)
k
```

```
##
## pred_rounded 0 1
##              0 91 5
##              1 9 95
```

```
(k[2]+k[3])/(k[1]+k[2]+k[3]+k[4]) #proportion of table elements that are "off-diagonal"
```

```
## [1] 0.07
```

In the following code chunk, we input seven measurements for each of 5000 asteroids. The data frame is `df`. There is another variable, `q`, that is also loaded and which we will utilize later.

```
load(url("http://www.stat.cmu.edu/~mfarag/350/HW_09_PCA.Rdata"))
names(df)
```

```
## [1] "diameter" "albedo" "a" "e" "i" "per_y" "H"
```

Question 5

(20 points)

Perform PCA on the data frame `df`. Use the rule-of-thumb in the notes to determine the number of principal components (or PCs) to retain, and for those PCs, indicate the mapping from PC to original variables. (Also, display the proportion of variance explained by the PCs you retain.) Are any of the original variables “unimportant” within the context of the retained PCs?

```
p=prcomp(df,scale=TRUE)
a=p$sdev[1]^2/sum(p$sdev^2)
a #0.3354716
```

```
## [1] 0.3354716
```

```
b=p$sdev[2]^2/sum(p$sdev^2)
b #0.2615384
```

```
## [1] 0.2615384
```

```
c=p$sdev[3]^2/sum(p$sdev^2)
c #0.145932
```

```
## [1] 0.145932
```

```
d=p$sdev[4]^2/sum(p$sdev^2)
d #0.1234097
```

```
## [1] 0.1234097
```

```
e=p$sdev[5]^2/sum(p$sdev^2)
e #0.1076712
```

```
## [1] 0.1076712
```

```
f=p$sdev[6]^2/sum(p$sdev^2)
f #0.02251361
```

```
## [1] 0.02251361
```

4 Q4 20 / 20

✓ - 0 pts Correct

- 4 pts incorrect setup for bivariate normal model 1
- 4 pts incorrect setup for bivariate normal model 2
- 2 pts does not sample 100 samples from each model
- 3 pts no plot with different color
- 5 pts not using glm for classification
- 4 pts missing/incorrect % of misclassification
- 4 pts question links incorrectly
- 20 pts blank


```
load(url("http://www.stat.cmu.edu/~mfarag/350/HW_09_PCA.Rdata"))
names(df)
```

```
## [1] "diameter" "albedo" "a" "e" "i" "per_y" "H"
```

Question 5

(20 points)

Perform PCA on the data frame `df`. Use the rule-of-thumb in the notes to determine the number of principal components (or PCs) to retain, and for those PCs, indicate the mapping from PC to original variables. (Also, display the proportion of variance explained by the PCs you retain.) Are any of the original variables “unimportant” within the context of the retained PCs?

```
p=prcomp(df,scale=TRUE)
a=p$sdev[1]^2/sum(p$sdev^2)
a #0.3354716
```

```
## [1] 0.3354716
```

```
b=p$sdev[2]^2/sum(p$sdev^2)
b #0.2615384
```

```
## [1] 0.2615384
```

```
c=p$sdev[3]^2/sum(p$sdev^2)
c #0.145932
```

```
## [1] 0.145932
```

```
d=p$sdev[4]^2/sum(p$sdev^2)
d #0.1234097
```

```
## [1] 0.1234097
```

```
e=p$sdev[5]^2/sum(p$sdev^2)
e #0.1076712
```

```
## [1] 0.1076712
```

```
f=p$sdev[6]^2/sum(p$sdev^2)
f #0.02251361
```

```
## [1] 0.02251361
```

```
g=p$sdev[7]^2/sum(p$sdev^2)
g #0.003463479
```

```
## [1] 0.003463479
```

```
#Using the rule-of-thumb, we retain 5 principal components. PC1,2,3,4,5
p$rotation[,1:5] #The mapping from PC to original variables.for PC1-5
```

```
##           PC1           PC2           PC3           PC4           PC5
## diameter -0.10427784  0.64018993 -0.381807654 -0.06374223  0.08430670
## albedo   0.02926049  0.27356501  0.891476176 -0.02430516 -0.21790935
## a        -0.61875917 -0.08069235  0.091831781 -0.14969568  0.25344191
## e        -0.30222711  0.04700618 -0.006640299  0.94592379 -0.10576940
## i        -0.37669959 -0.08308185 -0.183319322 -0.21654799 -0.87665888
## per_y    -0.60414428 -0.08470000  0.131810599 -0.16680797  0.31781953
## H        0.08469313 -0.70179878 -0.005997209  0.05871658  0.02116336
```

```
pcs_var = p$sdev^2
pcs_var_pct = round(pcs_var/sum(pcs_var)*100,2)
pcs_var_pct[1:5] #Proportion of variance explained by the PCs you retain.)
```

```
## [1] 33.55 26.15 14.59 12.34 10.77
```

None of the original variables is ‘unimportant’ within the context of the retained PCs.

Question 6

(20 points)

Something that one can do with principal components is regression: after all, all you’ve done is transform your data to a new coordinate system.

Below, linearly regress the variable `q` upon all the variables in `df`, and print the adjusted R^2 and the sum of squared errors for the model.

Then repeat linear regression, except now regress the variable `q` upon only the retained PCs. Again, print out the adjusted R^2 and the sum of squared errors. Are the second value close to the first? (They often will be, but don’t have to be.)

(Hint: look at the names of the list elements that are output by the `summary` of your linear regression fits, as one of those list elements may help you with extracting the adjusted R^2 . As far as the sum of squared errors, you need to simply compute the `sum()` of the difference between the observed values of `q` and the predicted values from `predict()`.)

```
fit<-lm(q~diameter+albedo+a+e+i+per_y+H,data=df)
summary(fit)
```

```
##
## Call:
## lm(formula = q ~ diameter + albedo + a + e + i + per_y + H, data = df)
##
```

5 Q5 20 / 20

✓ - 0 pts Correct

- 5 pts incorrect # of PCs to retain
- 4 pts missing/incorrect proportion of variance explained by PCs
- 6 pts missing/incorrect conclusion of original variables' importance
- 4 pts question links incorrectly
- 20 pts blank

```
g=p$sdev[7]^2/sum(p$sdev^2)
g #0.003463479
```

```
## [1] 0.003463479
```

```
#Using the rule-of-thumb, we retain 5 principal components. PC1,2,3,4,5
p$rotation[,1:5] #The mapping from PC to original variables.for PC1-5
```

```
##           PC1           PC2           PC3           PC4           PC5
## diameter -0.10427784  0.64018993 -0.381807654 -0.06374223  0.08430670
## albedo   0.02926049  0.27356501  0.891476176 -0.02430516 -0.21790935
## a        -0.61875917 -0.08069235  0.091831781 -0.14969568  0.25344191
## e        -0.30222711  0.04700618 -0.006640299  0.94592379 -0.10576940
## i        -0.37669959 -0.08308185 -0.183319322 -0.21654799 -0.87665888
## per_y    -0.60414428 -0.08470000  0.131810599 -0.16680797  0.31781953
## H         0.08469313 -0.70179878 -0.005997209  0.05871658  0.02116336
```

```
pcs_var = p$sdev^2
pcs_var_pct = round(pcs_var/sum(pcs_var)*100,2)
pcs_var_pct[1:5] #Proportion of variance explained by the PCs you retain.)
```

```
## [1] 33.55 26.15 14.59 12.34 10.77
```

None of the original variables is ‘unimportant’ within the context of the retained PCs.

Question 6

(20 points)

Something that one can do with principal components is regression: after all, all you’ve done is transform your data to a new coordinate system.

Below, linearly regress the variable `q` upon all the variables in `df`, and print the adjusted R^2 and the sum of squared errors for the model.

Then repeat linear regression, except now regress the variable `q` upon only the retained PCs. Again, print out the adjusted R^2 and the sum of squared errors. Are the second value close to the first? (They often will be, but don’t have to be.)

(Hint: look at the names of the list elements that are output by the `summary` of your linear regression fits, as one of those list elements may help you with extracting the adjusted R^2 . As far as the sum of squared errors, you need to simply compute the `sum()` of the difference between the observed values of `q` and the predicted values from `predict()`.)

```
fit<-lm(q~diameter+albedo+a+e+i+per_y+H,data=df)
summary(fit)
```

```
##
## Call:
## lm(formula = q ~ diameter + albedo + a + e + i + per_y + H, data = df)
##
```

```
## Residuals:
##      Min       1Q   Median       3Q      Max
## -26.3444  -0.2649  -0.0679   0.3428  15.7684
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.207e-01  1.704e-01   0.708   0.479
## diameter     7.129e-03  4.616e-04  15.444 <2e-16 ***
## albedo       2.198e-02  1.773e-01   0.124   0.901
## a            6.331e-01  6.135e-03 103.191 <2e-16 ***
## e           -6.280e+00  1.291e-01 -48.629 <2e-16 ***
## i            6.241e-05  1.282e-03   0.049   0.961
## per_y       -3.239e-02  3.352e-04 -96.647 <2e-16 ***
## H            1.327e-01  1.285e-02  10.326 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.7033 on 4992 degrees of freedom
## Multiple R-squared:  0.7394, Adjusted R-squared:  0.739
## F-statistic: 2023 on 7 and 4992 DF,  p-value: < 2.2e-16
```

```
summary(fit)$adj.r.squared # 0.7390329
```

```
## [1] 0.7390329
```

```
sum(summary(fit)$residual^2) # 2469.368
```

```
## [1] 2469.368
```

```
fit1<-lm(q~p$x[,1]+p$x[,2]+p$x[,3]+p$x[,4]+p$x[,5])
summary(fit1)$adj.r.squared #0.1842183
```

```
## [1] 0.1842183
```

```
sum(summary(fit1)$residual^2) #7722.323
```

```
## [1] 7722.323
```

```
fit1<-lm(q~p$x[,1]+p$x[,2]+p$x[,3]+p$x[,4]+p$x[,5])
```

The second value is not close to the first.

6 Q6 20 / 20

✓ - 0 pts Correct

- 2 pts incorrect setup for $lm.orig$
- 2 pts missing/incorrect R^2 for $lm.orig$
- 2 pts missing/incorrect SSE for $lm.orig$
- 4 pts incorrect/missing setup for $lm.pc$
- 2 pts missing/incorrect R^2 for $lm.pc$
- 2 pts missing/incorrect SSE for $lm.pc$
- 4 pts missing conclusion about whether R^2 and SSE are close
- 4 pts question links incorrectly
- 20 pts blank

7 Late Penalty 0 / 0

✓ - 0 pts Correct