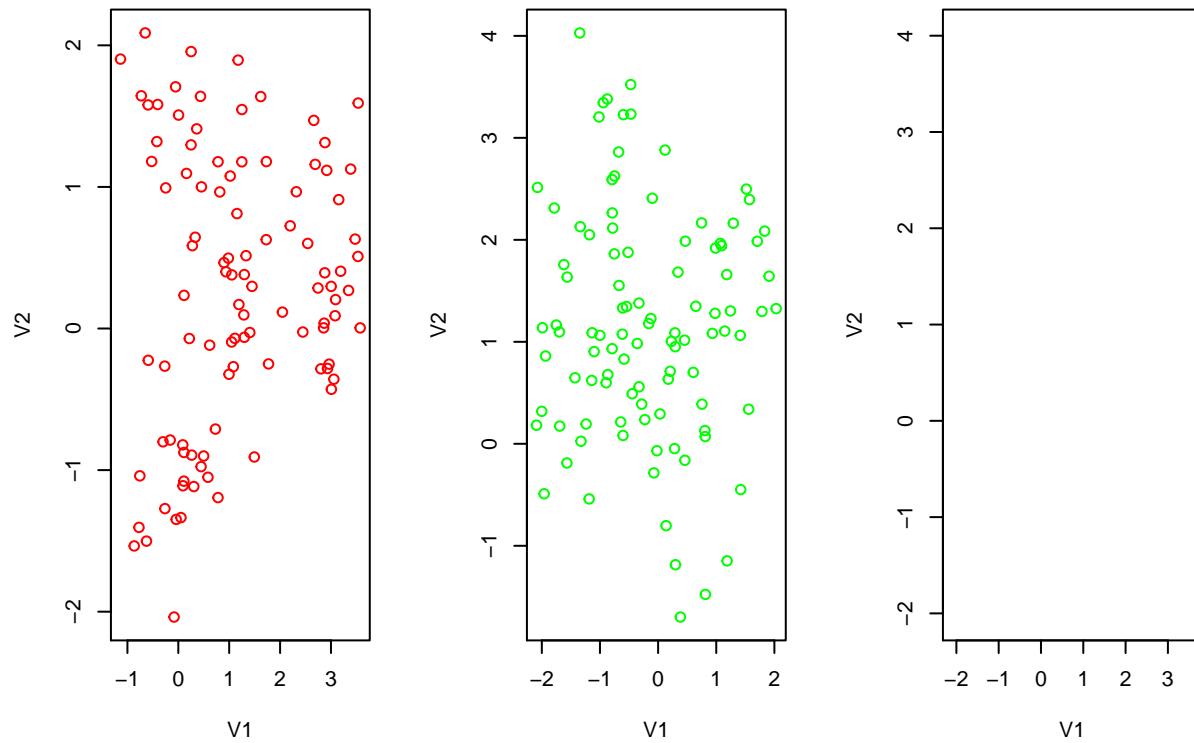
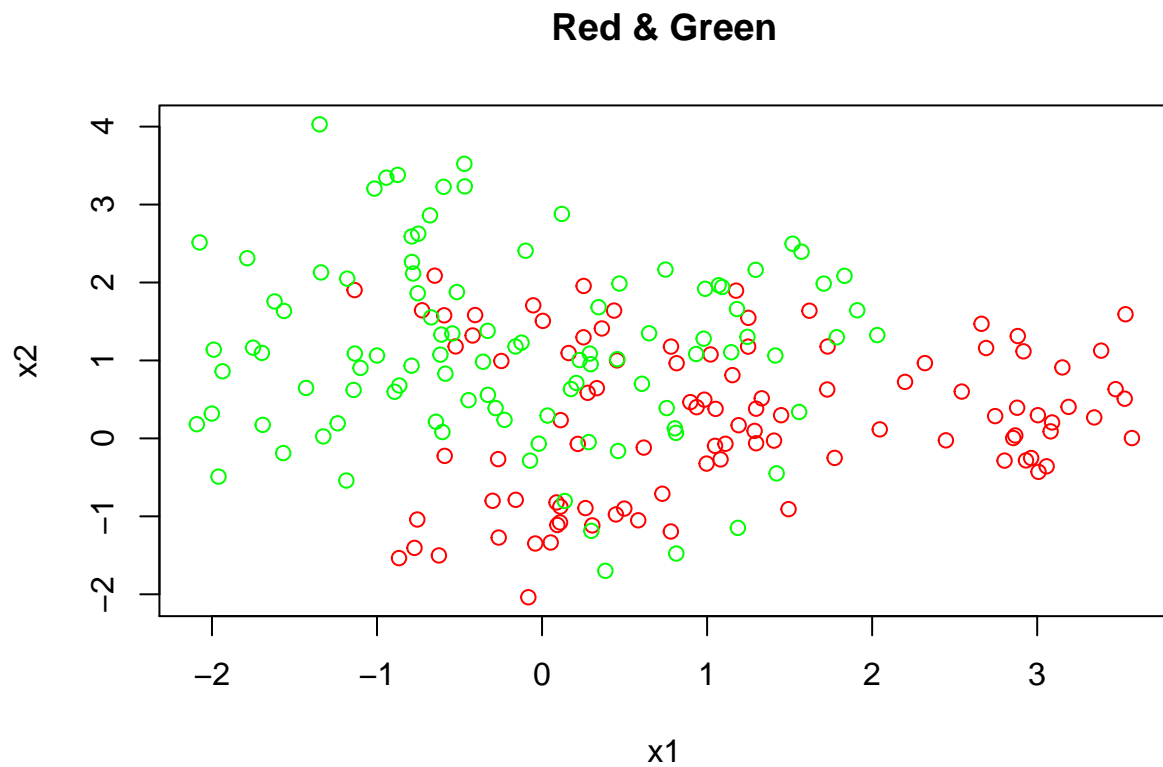


1

```
plot(X,type="n")
```



```
par(mfrow=c(1,1))
plot(X,type="n",main='Red & Green',xlab='x1',ylab='x2') # type="n":nothing
points(redpoints,col="red")
points(greenpoints,col="green")
```



```
##### Q3. #####
```

```
is.matrix(X)
```

```
## [1] FALSE
```

```
is.data.frame(X)
```

```
## [1] TRUE
```

```
X=as.matrix(X)
```

```
is.matrix(X)
```

```
## [1] TRUE
```

```
# Regression 1
```

```
X1=cbind(rep(1,nrow(X)),X) # design matrix : 200 by 3  
dim(X1)
```

```
## [1] 200 3
```

```
beta.hat=solve(t(X1)%*%X1)%*%t(X1)%*%y  
beta.hat
```

```
## [1]
```

```
## 0.5235718
```

```
## V1 0.1601129
```

```
## V2 -0.1415877
```

```

y.hat=X1%*%beta.hat # 200 by 1

result <- ifelse(y.hat>0.5,1,0)
result[,1]

##      [1] 1 1 1 0 0 1 1 0 1 1 0 1 1 1 1 0 1 1 1 1 1 1 0 0 1 1 1 1 1 1
##     [38] 1 1 1 1 1 1 1 1 0 1 1 1 1 1 0 0 0 1 1 1 1 0 1 1 1 1 1 1 1 0 1 1 1
##     [75] 1 0 1 1 1 0 1 1 1 0 1 1 1 1 0 0 1 1 1 1 0 1 1 1 1 0 0 0 1 0 0 0 1 0 1
##    [112] 0 0 1 0 0 0 0 0 0 0 1 1 0 0 0 0 0 0 1 0 0 0 0 0 0 1 1 0 0 0 0 1 1 0 0 0
##    [149] 0 0 0 1 0 0 0 0 0 0 1 1 0 0 0 1 0 1 0 0 0 0 1 1 0 1 0 0 0 0 0 1 0 1 0 0
##    [186] 0 0 0 0 0 0 1 0 0 1 0 0 0 1 0

sum(result[,1]==y)

## [1] 155

te=1-sum(result[,1]==y)/length(y); te

## [1] 0.225

# Regression 2
fit=lm(y~X[,1]+X[,2])
summary(fit)

##
## Call:
## lm(formula = y ~ X[, 1] + X[, 2])
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.87609 -0.33002  0.00249  0.30485  0.92773
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.52357    0.03653  14.332 < 2e-16 ***
## X[, 1]       0.16011    0.02124   7.539 1.69e-12 ***
## X[, 2]      -0.14159    0.02545  -5.563 8.59e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.408 on 197 degrees of freedom
## Multiple R-squared:  0.344, Adjusted R-squared:  0.3374
## F-statistic: 51.66 on 2 and 197 DF, p-value: < 2.2e-16

summary(fit)$coef

##              Estimate Std. Error  t value    Pr(>|t|)
## (Intercept)  0.5235718 0.03653068 14.332387 2.199775e-32
## X[, 1]       0.1601129 0.02123825  7.538895 1.692278e-12
## X[, 2]      -0.1415877 0.02545322 -5.562664 8.587034e-08

summary(fit)$coef[,1]

## (Intercept)      X[, 1]      X[, 2]
##    0.5235718    0.1601129   -0.1415877

beta.hat2 = fit$coefficients
y.hat2 = fit$fitted.values # y.hat2 = fitted.values(fit)

```

```

g.hat=as.numeric(y.hat>0.5) # T/F로 반환해서 numeric으로 1/0
g.hat

##      [1] 1 1 1 0 0 1 1 0 1 1 0 1 1 1 1 0 1 1 1 1 1 1 0 0 1 1 1 1 1 1
##     [38] 1 1 1 1 1 1 1 1 0 1 1 1 1 1 0 0 0 1 1 1 1 0 1 1 1 1 1 1 1 0 1 1 1
##     [75] 1 0 1 1 1 0 1 1 1 0 1 1 1 1 0 0 1 1 1 1 0 1 1 1 1 0 0 0 1 0 0 0 1 0 1
##    [112] 0 0 1 0 0 0 0 0 0 0 1 1 0 0 0 0 0 0 1 0 0 0 0 0 0 1 1 0 0 0 0 0 1 0 0 0
##    [149] 0 0 0 1 0 0 0 0 0 0 1 1 0 0 0 1 0 1 0 0 0 0 1 1 0 1 0 0 0 0 0 0 1 0 1 0 0
##    [186] 0 0 0 0 0 0 1 0 0 1 0 0 0 1 0

z=sum(g.hat==y)
training.error=1-z/200 # 오분류율을 찾는 거니까!
training.error # 0.225

## [1] 0.225

##### Q4. #####

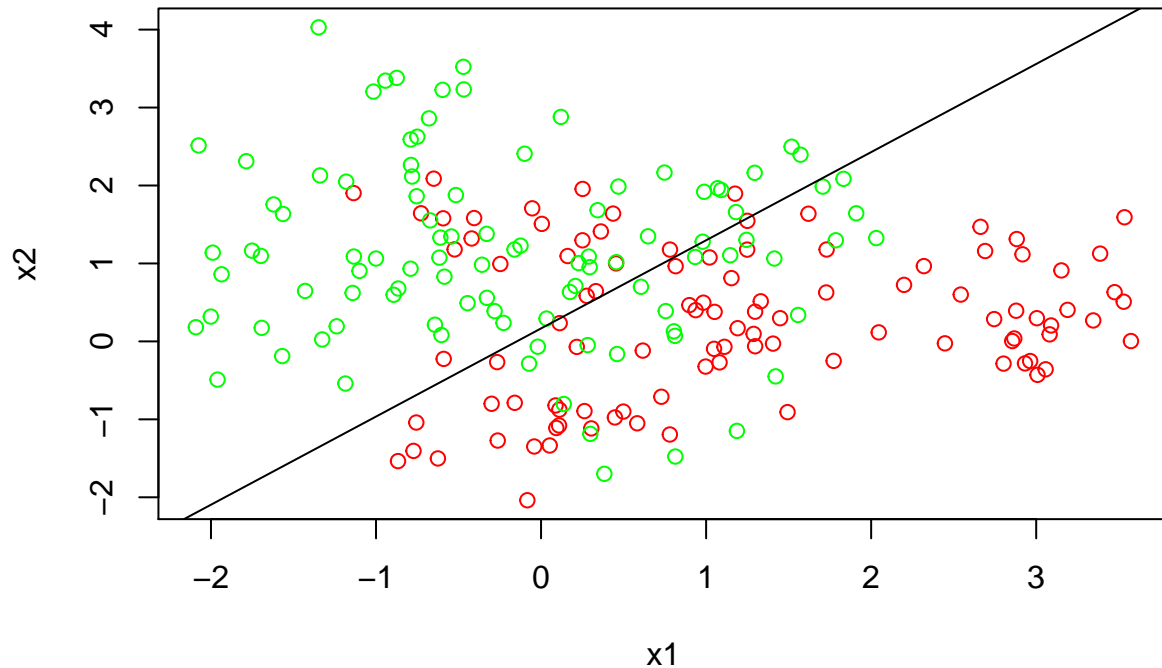
plot(X,type="n",main='Red & Green',xlab='x1',ylab='x2') # type="n":nothing
points(redpoints,col="red")
points(greenpoints,col="green")
beta.hat

##           [,1]
##      0.5235718
## V1  0.1601129
## V2 -0.1415877

abline((.5-beta.hat[1])/beta.hat[3],-beta.hat[2]/beta.hat[3])

```

## Red & Green



```
#abline((.5-beta.hat[1,1])/beta.hat[3,1],-beta.hat[2,1]/beta.hat[3,1])

##### Q5. #####

X0=rbind(redtestpoints,greentestpoints) # testset은 각 1000 by 2
X0=as.matrix(X0) # X0는 2000 by 2

y0=c(rep(1,1000),rep(0,1000)) # 마찬가지로 y를 1~1000(red)=1, 1001~2000(green)=0
X01=cbind(rep(1,2000),X0) # design matrix(2000 by 3)

# 추정한 모수들을 이용하여 y=xb, 즉 testset에서의 예측된 색(y)을 구함, testerror
y0.hat=X01%*%beta.hat
g0.hat=as.numeric(y0.hat>0.5)

test.error=1-sum(g0.hat==y0)/2000
test.error # 0.245 > 0.225(training error)
```

```
## [1] 0.245
```

```
# testerror is always larger than training error !!!
```

```
##### Q6. #####
```

```
# Training observations 사이의 거리를 계산해야 함.
# Euclidean distance matrix
D=matrix(0,200,200) # 200 by 200 matrix all values are zero.
```



```
# Euclidean distance matrix
```



```

D0=matrix(0,2000,200) # row = testset, col = trainingset
dim(D0)

## [1] 2000 200

# (x1,x1)~(x1,x200)~...~(x2000,x1)~(x2000,x200)
for (i in 1:2000) for (j in 1:200) D0[i,j]=sqrt(sum((X0[i,]-X[j,])^2))

k=7
g0.hat=rep(0,2000)
for (i in 1:2000) g0.hat[i]=(mean(y[order(D0[i,])[1:k]])>0.5)
head(g0.hat,100)

## [1] 1 0 1 1 0 1 1 1 0 1 1 1 0 0 1 1 1 1 0 1 1 1 0 1 0 0 0 1 1 1 0 1 0 1 0 1 1
## [38] 1 0 0 0 1 1 1 0 1 1 1 1 0 1 1 1 1 0 1 1 0 0 0 1 0 1 1 1 0 1 1 1 1 0 0 1 1
## [75] 1 1 1 0 1 1 0 1 1 1 0 1 1 0 0 0 1 0 1 0 0 1 1 1 1 1

test.error=1-sum(g0.hat==y0)/2000
test.error

## [1] 0.2965

k=3
g0.hat=rep(0,2000)
for (i in 1:2000) g0.hat[i]=(mean(y[order(D0[i,])[1:k]])>0.5)
head(g0.hat,100)

## [1] 1 0 1 1 0 1 1 1 0 1 1 1 0 1 1 1 1 1 1 1 0 1 0 0 1 1 1 0 1 1 0 1 0 1 1
## [38] 1 1 0 0 1 0 1 0 1 1 1 1 1 1 1 1 1 1 0 0 1 0 0 0 1 0 1 0 1 0 1 1 1 1 0 1 1 0
## [75] 1 1 1 1 0 1 0 0 1 1 0 1 1 0 0 0 1 0 0 0 0 1 1 1 1 0

test.error=1-sum(g0.hat==y0)/2000
test.error

## [1] 0.328

k=1
g0.hat=rep(0,2000)
for (i in 1:2000) g0.hat[i]=(mean(y[order(D0[i,])[1:k]])>0.5)
head(g0.hat, 100)

## [1] 1 1 1 1 0 1 0 1 0 1 1 1 0 0 1 0 1 1 0 1 1 1 0 1 0 0 0 0 1 1 1 1 0 0 1 0 1
## [38] 1 0 0 0 1 1 0 0 1 1 1 1 1 1 1 1 0 1 0 1 1 0 0 0 1 0 1 0 0 0 1 1 1 1 0 1 1 1
## [75] 1 1 1 1 1 1 0 1 1 1 0 1 1 0 0 0 1 0 1 0 0 1 1 1 0 1

test.error=1-sum(g0.hat==y0)/2000
test.error

## [1] 0.3225

g2 <- NULL
for(k in c(1,3,7)){
  for(i in 1:2000){
    g2[i] <- mean(y[order(D0[i,])[1:k]])>0.5
  }
  print(1-sum(g2==y0)/2000)
}

## [1] 0.3225

```

```
## [1] 0.328  
## [1] 0.2965
```