# Statistical Learning for Biological and Information Systems Problem Set #3

Chenye Yang `cy2540@columbia.edu`

October 22, 2019

## P1

### (a)

**Ans:**
Assume $A = $ "receive an A", $\bar{A} = $ "not receive an A". Thus:

$$\log \frac{p_A(X)}{p_{\bar{A}}(X)} = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$
$$= -6 + 0.05 X_1 + X_2$$

For $X_1 = 40, \ X_2 = 3.5$,

$$\log \frac{p_A(X)}{1 - p_{\bar{A}}(X)} = -6 + 0.05 \times 40 + 3.5 = -0.5$$

$\Rightarrow$

$$p_A(X) = 0.38$$

Therefore, the probability that a student who studies for 40 hours and has an undergrad GPA of 3.5 gets an A in the class is 0.38

### (b)

**Ans:**
When $p_A(X) = 0.5$, $\log \frac{p_A(X)}{p_{\bar{A}}(X)} = 0$

$$\begin{cases} -6 + 0.05 X_1 + X_2 = 0 \\ X_2 = 3.5 \end{cases}$$

$\Rightarrow$

$$X_1 = 50$$

Therefore, the student in part (a) need to study 50 hours to have a 50% chance of getting an A in the class.

Chenye Yang cy2540 ⁣1

# P2

**Ans:**

Given:

$$Y = \text{Yes} \quad \bar{X} = 10 \quad \hat{\sigma}^2 = 36 \quad \pi_{\text{Yes}} = 0.8 \quad X \sim N(10, 36)$$
$$Y = \text{No} \quad \bar{X} = 0 \quad \hat{\sigma}^2 = 36 \quad \pi_{\text{No}} = 0.2 \quad X \sim N(0, 36)$$

Thus,

$$
\begin{aligned}
P(Y = \text{Yes}|X = 4) &= \frac{\pi_{\text{Yes}} f_{\text{Yes}}(X)}{\pi_{\text{Yes}} f_{\text{Yes}}(X) + \pi_{\text{No}} f_{\text{No}}(X)} \\
&= \frac{0.8 \times \frac{1}{6 \times \sqrt{2\pi}} e^{-\frac{(4-10)^2}{2 \times 36}}}{0.8 \times \frac{1}{6 \times \sqrt{2\pi}} e^{-\frac{(4-10)^2}{2 \times 36}} + 0.2 \times \frac{1}{6 \times \sqrt{2\pi}} e^{-\frac{(4-0)^2}{2 \times 36}}} \\
&= \frac{4 \times e^{-1/2}}{4 \times e^{-1/2} + e^{-2/9}} \\
&= 0.752
\end{aligned}
$$

Therefore, the probability that a company will issue a dividend this year given that its percentage profit was $X = 4$ last year is 0.752

## P3

**Ans:**

Assume $\boldsymbol{X_1} = [1, 0.0]^T$, $\boldsymbol{X_2} = [1, 0.2]^T$, $\boldsymbol{X_3} = [1, 0.4]^T$, $\boldsymbol{X_4} = [1, 0.6]^T$, $\boldsymbol{X_5} = [1, 0.8]^T$, $\boldsymbol{X_6} = [1, 1.0]^T$ and $y_1 = 0$, $y_2 = 0$, $y_3 = 0$, $y_4 = 1$, $y_5 = 0$, $y_6 = 1$ and $\boldsymbol{\beta} = [\beta_0, \beta_1]^T$.

The log-likelihood function, $\ell(\boldsymbol{\beta}) = \ell(\beta_0, \beta_1)$:

$$
\begin{aligned}
\ell(\boldsymbol{\beta}) &= \sum_{i=1}^{6} \left[ y_i \log p\left(\boldsymbol{X_i}; \boldsymbol{\beta}\right) + (1 - y_i) \log \left(1 - p\left(\boldsymbol{X_i}; \boldsymbol{\beta}\right)\right) \right] \\
&= \sum_{i=1}^{6} \left[ y_i \boldsymbol{\beta}^T \boldsymbol{X_i} - \log \left(1 + e^{\boldsymbol{\beta}^T \boldsymbol{X_i}}\right) \right] \\
&= -\log(1 + e^{\boldsymbol{\beta}^T \boldsymbol{X_1}}) - \log(1 + e^{\boldsymbol{\beta}^T \boldsymbol{X_2}}) - \log(1 + e^{\boldsymbol{\beta}^T \boldsymbol{X_3}}) + \boldsymbol{\beta}^T \boldsymbol{X_4} - \log(1 + e^{\boldsymbol{\beta}^T \boldsymbol{X_4}}) \\
&\quad - \log(1 + e^{\boldsymbol{\beta}^T \boldsymbol{X_5}}) + \boldsymbol{\beta}^T \boldsymbol{X_6} - \log(1 + e^{\boldsymbol{\beta}^T \boldsymbol{X_6}}) \\
&= -\log(1 + e^{\beta_0}) - \log(1 + e^{\beta_0 + \beta_1 \times 0.2}) - \log(1 + e^{\beta_0 + \beta_1 \times 0.4}) + (\beta_0 + \beta_1 \times 0.6) \\
&\quad - \log(1 + e^{\beta_0 + \beta_1 \times 0.6}) - \log(1 + e^{\beta_0 + \beta_1 \times 0.8}) + (\beta_0 + \beta_1 \times 1.0) - \log(1 + e^{\beta_0 + \beta_1 \times 1.0})
\end{aligned}
\tag{1}
$$

To maximize the log-likelihood, we set its derivatives to zero:

$$
\frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \sum_{i=1}^{6} \boldsymbol{X_i} \left[y_i - p\left(\boldsymbol{X_i}; \boldsymbol{\beta}\right)\right] = 0
\tag{2}
$$

In practice, we use the Newton–Raphson algorithm to solve the maximizing problem:

$$
\begin{aligned}
\frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} &= \sum_{i=1}^{6} \boldsymbol{X_i} \left[y_i - p\left(\boldsymbol{X_i}; \boldsymbol{\beta}\right)\right] \\
\frac{\partial^2 \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} &= -\sum_{i=1}^{6} \boldsymbol{X_i} \boldsymbol{X_i}^T p\left(\boldsymbol{X_i}; \boldsymbol{\beta}\right) \left(1 - p\left(\boldsymbol{X_i}; \boldsymbol{\beta}\right)\right) \\
\boldsymbol{\beta}^{\text{new}} &= \boldsymbol{\beta}^{\text{old}} - \left(\frac{\partial^2 \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T}\right)^{-1} \frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}}
\end{aligned}
\tag{3}
$$

**Code:**

```
1  beta_old = c(0,0)
2  beta_new = c(0,0)
3  # initiate X
4  X = matrix(data = rep(0,12), nrow = 2, ncol = 6)
5  for (i in 1:6){
6      X[,i] = c(1,(i-1)/5)
7  }
8  # initiate y
9  y = c(0, 0, 0, 1, 0, 1)
10 # initiate derivatives
11 l_b = c(0, 0)
12 for (i in 1:6){
```

```
13     l_b = l_b + X[,i]*(y[i]-exp(X[,i] %*% beta_old)/(1+exp(X[,i] %*%
           beta_old)))
14  }
15  # initiate second derivatives
16  l_bb = matrix(data = rep(0,4), nrow = 2, ncol = 2)
17  for (i in 1:6){
18     l_bb = l_bb - (X[,i] %*% t(X[,i]))*((exp(X[,i] %*% beta_old)/(1+exp(
           X[,i] %*% beta_old)))*(1-(exp(X[,i] %*% beta_old)/(1+exp(X[,i] %*
           % beta_old))))))[1]
19  }
20  # iteration
21  for (i in 1:10){
22     beta_new = t(t(beta_old)) - solve(l_bb) %*% t(t(l_b))
23     beta_old = t(t(beta_new))
24  }
25  t(t(beta_new))
```

**Result:**

```
1  > t(t(beta_new))
2              [,1]
3  [1,]  -23.80952
4  [2,]   34.28571
```

The estimated coefficients of this logistic regression problem is $\hat{\beta}_0 = -23.81$, $\hat{\beta}_1 = 34.29$.

## P4

**Ans:**

From given,

$$\begin{aligned}
\text{cov}[\boldsymbol{Y}] &= \text{cov}[\boldsymbol{A}\boldsymbol{X}] \\
&= \boldsymbol{A}\,\text{cov}[\boldsymbol{X}]\boldsymbol{A}^T \\
&= \boldsymbol{A}\boldsymbol{\Sigma}\boldsymbol{A}^T = \boldsymbol{I}
\end{aligned} \tag{4}$$

Because $\boldsymbol{\Sigma}$ is symmetric, assume $\boldsymbol{\Sigma}$ has non-degenerate eigenvalues $\lambda_1,\ \lambda_2$ and corresponding linearly independent eigenvectors $v_1,\ v_2$. Define:

$$\boldsymbol{Q} = \left[\begin{array}{cc} \boldsymbol{v_1} & \boldsymbol{v_2} \end{array}\right] = \left[\begin{array}{cc} v_1^1 & v_2^1 \\ v_1^2 & v_2^2 \end{array}\right]$$

$$\boldsymbol{\Lambda} = \left[\begin{array}{cc} \lambda_1 & 0 \\ 0 & \lambda_2 \end{array}\right] \tag{5}$$

Thus we have,

$$\boldsymbol{\Sigma} = \boldsymbol{Q}\boldsymbol{\Lambda}\boldsymbol{Q}^T = \boldsymbol{Q}\boldsymbol{\Lambda}\boldsymbol{Q}^{-1} \tag{6}$$

Therefore,

$$\boldsymbol{A}\boldsymbol{Q}\boldsymbol{\Lambda}\boldsymbol{Q}^T\boldsymbol{A}^T = \boldsymbol{I}$$

$$(\boldsymbol{A}\boldsymbol{Q})\boldsymbol{\Lambda}(\boldsymbol{A}\boldsymbol{Q})^T = \boldsymbol{I} \tag{7}$$

$$\Rightarrow$$

$$\boldsymbol{A}\boldsymbol{Q} = \left(\boldsymbol{\Lambda}^{-1}\right)^{\frac{1}{2}}$$

$$\boldsymbol{A} = \left(\boldsymbol{\Lambda}^{-1}\right)^{\frac{1}{2}}\boldsymbol{Q}^{-1} \tag{8}$$

Using Matlab program to calculate the eigenvalues and eigenvectors of $\boldsymbol{\Lambda}$:

**Code:**

```
1  # Matlab
2  syms sigma_1 sigma_2 rho
3  A = [sigma_1^2, rho*sigma_1*sigma_2;
4       rho*sigma_1*sigma_2, sigma_2^2];
5  [eigenvector, eigenvalue] = eig(A)
```

**Result:**

$$\boldsymbol{Q} = \left[\begin{array}{cc} \boldsymbol{v_1} & \boldsymbol{v_2} \end{array}\right]$$

$$\boldsymbol{v_1} = \left[\begin{array}{c} (\sigma_1^2/2 - (4\rho^2\sigma_1^2\sigma_2^2 + \sigma_1^4 - 2\sigma_1^2\sigma_2^2 + \sigma_2^4)^{1/2}/2 + \sigma_2^2/2)/(\rho\sigma_1\sigma_2) - \sigma_2/(\rho\sigma_1) \\ 1 \end{array}\right]$$

$$\boldsymbol{v_2} = \left[\begin{array}{c} ((4\rho^2\sigma_1^2\sigma_2^2 + \sigma_1^4 - 2\sigma_1^2\sigma_2^2 + \sigma_2^4)^{1/2}/2 + \sigma_1^2/2 + \sigma_2^2/2)/(\rho\sigma_1\sigma_2) - \sigma_2/(\rho\sigma_1) \\ 1 \end{array}\right]$$

$$\boldsymbol{\Lambda} = \left[\begin{array}{cc} \lambda_1 & 0 \\ 0 & \lambda_2 \end{array}\right]$$

$$\lambda_1 = \sigma_1^2/2 - (4*\rho^2*\sigma_1^2*\sigma_2^2 + \sigma_1^4 - 2*\sigma_1^2*\sigma_2^2 + \sigma_2^4)^{1/2}/2 + \sigma_2^2/2$$

$$\lambda_2 = (4*\rho^2*\sigma_1^2*\sigma_2^2 + \sigma_1^4 - 2*\sigma_1^2*\sigma_2^2 + \sigma_2^4)^{1/2}/2 + \sigma_1^2/2 + \sigma_2^2/2$$

(9)

Therefore, use Matlab program to calculate $\boldsymbol{A}$:

**Code:**

```
1  A = (inv(eigenvalue))^(1/2)*inv(eigenvector)
```

**Result:**

$$\boldsymbol{A} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

$a_{11} = -(2^{1/2}\rho\sigma_1\sigma_2(1/(\sigma_1^2 - (4\rho^2\sigma_1^2\sigma_2^2 + \sigma_1^4 - 2\sigma_1^2\sigma_2^2 + \sigma_2^4)^{1/2} + \sigma_2^2))^{1/2})/(4\rho^2\sigma_1^2\sigma_2^2 + \sigma_1^4 - 2\sigma_1^2\sigma_2^2 + \sigma_2^4)^{1/2}$

$a_{12} = (2^{1/2}(1/(\sigma_1^2 - (4\rho^2\sigma_1^2\sigma_2^2 + \sigma_1^4 - 2\sigma_1^2\sigma_2^2 + \sigma_2^4)^{1/2} + \sigma_2^2))^{1/2}((4\rho^2\sigma_1^2\sigma_2^2 + \sigma_1^4 - 2\sigma_1^2\sigma_2^2 + \sigma_2^4)^{1/2} + \sigma_1^2 - \sigma_2^2))$
$\qquad /(2(4\rho^2\sigma_1^2\sigma_2^2 + \sigma_1^4 - 2\sigma_1^2\sigma_2^2 + \sigma_2^4)^{1/2})$

$a_{21} = (2^{1/2}\rho\sigma_1\sigma_2(1/((4\rho^2\sigma_1^2\sigma_2^2 + \sigma_1^4 - 2\sigma_1^2\sigma_2^2 + \sigma_2^4)^{1/2} + \sigma_1^2 + \sigma_2^2))^{1/2})/(4\rho^2\sigma_1^2\sigma_2^2 + \sigma_1^4 - 2\sigma_1^2\sigma_2^2 + \sigma_2^4)^{1/2}$

$a_{22} = (2^{1/2}(1/((4\rho^2\sigma_1^2\sigma_2^2 + \sigma_1^4 - 2\sigma_1^2\sigma_2^2 + \sigma_2^4)^{1/2} + \sigma_1^2 + \sigma_2^2))^{1/2}((4\rho^2\sigma_1^2\sigma_2^2 + \sigma_1^4 - 2\sigma_1^2\sigma_2^2 + \sigma_2^4)^{1/2} - \sigma_1^2 + \sigma_2^2))$
$\qquad /(2(4\rho^2\sigma_1^2\sigma_2^2 + \sigma_1^4 - 2\sigma_1^2\sigma_2^2 + \sigma_2^4)^{1/2})$

$$(10)$$

# P5

**Ans:**

Because the variance of each population is the same $(\sigma^2)$, we have:

$$\hat{\sigma}_k^2 = \frac{1}{n_k - 1} \sum_{i:y_i=k} (x_i - \hat{\mu}_k)^2 \tag{11}$$

$$\mathbb{E}\left[\hat{\sigma}_k^2\right] = \sigma^2$$

Thus,

$$\mathbb{E}\left[\hat{\sigma}^2\right] = \mathbb{E}\left[\sum_{k=1}^{K} \alpha_k \hat{\sigma}_k^2\right] = \sum_{k=1}^{K} \alpha_k \mathbb{E}\left[\hat{\sigma}_k^2\right]$$

$$= \sigma^2 \sum_{k=1}^{K} \alpha_k = \sigma^2 \tag{12}$$

Because $x_i$ obeys Gaussian distribution, thus

$$\frac{(n_k - 1)\,\hat{\sigma}_k^2}{\sigma^2} \sim \mathcal{X}_{n_k-1}^2$$

Therefore,

$$\mathrm{var}\left[\frac{(n_k - 1)\,\hat{\sigma}_k^2}{\sigma^2}\right] = 2\,(n_k - 1)$$

$$\frac{(n_k - 1)^2}{\sigma^4}\,\mathrm{var}\left[\hat{\sigma}_k^2\right] = 2\,(n_k - 1) \tag{13}$$

$$\mathrm{var}\left[\hat{\sigma}_k^2\right] = \frac{2\sigma^4}{n_k - 1}$$

Considering,

$$\mathrm{var}\left[\hat{\sigma}^2\right] = \mathrm{var}\left[\sum_{k=1}^{K} \alpha_k \hat{\sigma}_k^2\right]$$

$$= \sum_{k=1}^{K} \alpha_k^2\,\mathrm{var}\left[\hat{\sigma}_k^2\right]$$

$$= \sum_{k=1}^{K} \alpha_k^2 \frac{2\sigma^4}{n_k - 1}$$

$$= 2\sigma^4 \sum_{k=1}^{K} \frac{\alpha_k^2}{n_k - 1} \sum_{i=1}^{K} \alpha_i \tag{14}$$

$$= 2\sigma^4 \sum_{k=1}^{K} \frac{\alpha_k^2}{n_k - 1} \sum_{i=1}^{K} \frac{n_i - 1}{n - K}$$

$$= \frac{2\sigma^4}{n - K} \sum_{k=1}^{K} \frac{\alpha_k^2}{n_k - 1} \sum_{i=1}^{K} (n_i - 1)$$

According to Cauchy-Schwarz inequality,

$$\sum_{k=1}^{K} \frac{\alpha_k^2}{n_k - 1} \sum_{k=1}^{K} (n_k - 1) \geq \left( \sum_{k=1}^{K} \frac{\alpha_k}{\sqrt{n_k - 1}} \times \sqrt{n_k - 1} \right)^2$$

$$= \left( \sum_{k=1}^{K} \alpha_k \right) \tag{15}$$

$$= 1$$

Therefore,

$$\mathrm{var}\left[ \hat{\sigma}^2 \right] \geq \frac{2\sigma^4}{n - K} \tag{16}$$

The equality happens if and only if:

$$\frac{\alpha_1}{n_1 - 1} = \frac{\alpha_2}{n_2 - 1} = ... = \frac{\alpha_K}{n_K - 1} \tag{17}$$

which can be satisfied by $\alpha_k = (n_k - 1)/(n - K)$, i.e.

$$\frac{\alpha_1}{n_1 - 1} = \frac{\alpha_2}{n_2 - 1} = ... = \frac{\alpha_K}{n_K - 1} = \frac{1}{n - K} \tag{18}$$

Therefore, $\alpha_k = (n_k - 1)/(n - K)$ minimizes the variance of $\sigma^2$ under the Gaussian assumption.

# P6

**Ans:**

Majority vote approach:

$X$ is classified to *Red*: 6 estimates, $X$ is classified to *Green*: 4 estimates. Thus, $X$ is classified to *Red*.

Average probability:

$$\begin{aligned}
\bar{\mathbb{P}}[\text{Class is Red}|X] &= \frac{\sum^n \mathbb{P}[\text{Class is Red}|X]}{n} \\
&= \frac{0.1 + 0.15 + 0.2 + 0.2 + 0.55 + 0.6 + 0.6 + 0.65 + 0.7 + 0.75}{10} \\
&= 0.45
\end{aligned}$$

Because $\bar{\mathbb{P}}[\text{Class is Red}|X] = 0.45 < 0.5$, $X$ is classified to *Green*.

# P7

## (a)

**Code:**

```
1  # (a)
2  rm(list=ls())
3  set.seed(1000)
4  load('/Users/yangchenye/Downloads/OJ.rda')
5  dim(OJ)
6  OJ=na.omit(OJ) # remove incomplete cases
7  dim(OJ)
8  names(OJ)
9
10 # 800: the sample size
11 smp_size = floor(800)
12 train_ind = sample(seq_len(nrow(OJ)), size = smp_size)
13 # select train data and test data
14 train = subset(OJ[train_ind, ])
15 test = subset(OJ[-train_ind, ])
```

## (b)

**Code:**

```
1  # (b)
2  library(tree)
3  tree.fit = tree(train$Purchase~., data=train)
4  summary(tree.fit)
```

**Result:**

```
1  > summary(tree.fit)
2
3  Classification tree:
4  tree(formula = train$Purchase ~ ., data = train)
5  Variables actually used in tree construction:
6  [1] "LoyalCH"      "PriceDiff"    "SalePriceMM"
7  Number of terminal nodes:  8
8  Residual mean deviance:  0.7486 = 592.9 / 792
9  Misclassification error rate: 0.16 = 128 / 800
```

**Ans:**

The training error rate is 128/800=0.16
The tree has 8 terminal nodes.

## (c)

**Code:**

```
1  # (c)
2  tree.fit
```

**Result:**

```
 1  > tree.fit
 2  node), split, n, deviance, yval, (yprob)
 3          * denotes terminal node
 4
 5  1) root 800 1066.00 CH ( 0.61500 0.38500 )
 6  2) LoyalCH < 0.5036 353    422.60 MM ( 0.28612 0.71388 )
 7  4) LoyalCH < 0.276142 170    131.00 MM ( 0.12941 0.87059 )
 8  8) LoyalCH < 0.035047 57     10.07 MM ( 0.01754 0.98246 ) *
 9  9) LoyalCH > 0.035047 113    108.50 MM ( 0.18584 0.81416 ) *
10  5) LoyalCH > 0.276142 183    250.30 MM ( 0.43169 0.56831 )
11  10) PriceDiff < 0.05 78     79.16 MM ( 0.20513 0.79487 ) *
12  11) PriceDiff > 0.05 105    141.30 CH ( 0.60000 0.40000 ) *
13  3) LoyalCH > 0.5036 447    337.30 CH ( 0.87472 0.12528 )
14  6) LoyalCH < 0.764572 187    206.40 CH ( 0.75936 0.24064 )
15  12) SalePriceMM < 2.125 120    156.60 CH ( 0.64167 0.35833 )
16  24) PriceDiff < −0.35 16     17.99 MM ( 0.25000 0.75000 ) *
17  25) PriceDiff > −0.35 104    126.70 CH ( 0.70192 0.29808 ) *
18  13) SalePriceMM > 2.125 67     17.99 CH ( 0.97015 0.02985 ) *
19  7) LoyalCH > 0.764572 260     91.11 CH ( 0.95769 0.04231 ) *
```

**Ans:**

For the terminal node:

$$8) \quad LoyalCH < 0.035047 \quad 57 \quad 10.07 \quad MM \quad (0.01754 \, 0.98246)*$$

"8)" is the node number.

"$LoyalCH < 0.035047$" is a two-column matrix of the labels for the left and right splits at the node.

"57" is the number of cases reaching that node ($LoyalCH < 0.035047$).

"10.07" is the deviance of the node.

"$MM$" is the fitted value at the node (the mean for regression trees, a majority class for classification trees).

"(0.01754 0.98246)" is a matrix of fitted probabilities for each response level.

## (d)

**Code:**

```
1  # (d)
2  plot(tree.fit)
3  text(tree.fit, pretty=0)
```
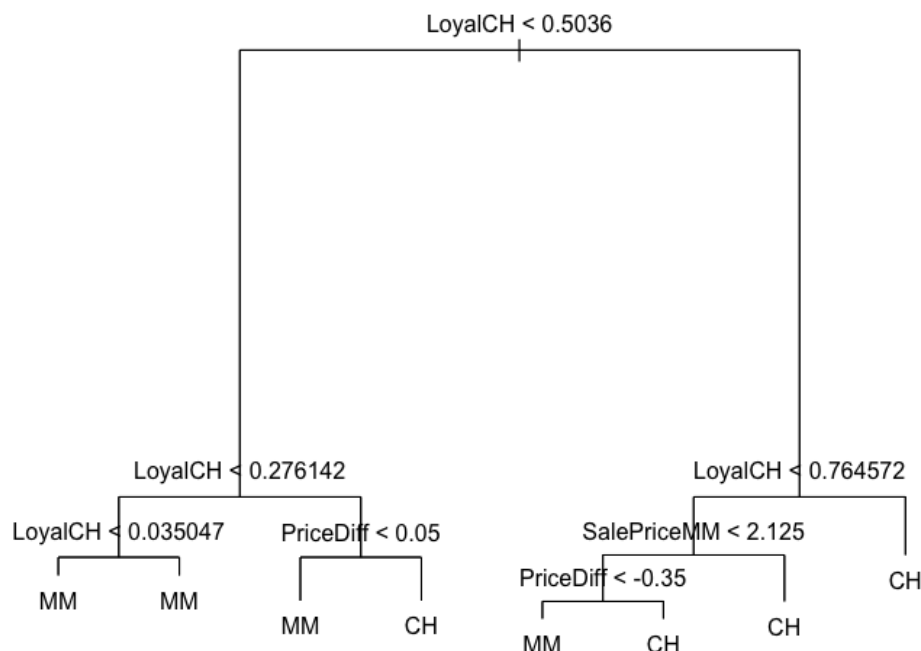
**Result:**



Figure 1: Regression tree with 8 terminal nodes

**Ans:**
For customers with $LoyalCH < 0.276142$, they will choose Minute Maid Orange Juice.
For customers with $0.276142 < LoyalCH < 0.5036$ & $PriceDiff < 0.05$, they will choose Minute Maid Orange Juice.
For customers with $0.5036 < LoyalCH < 0.764572$ & $SalePriceMM < 2.125$ & $PriceDiff < -0.35$, they will choose Minute Maid Orange Juice.
Otherwise, customers will choose Citrus Hill Orange Juice.

## (e)

**Code:**

```
# (e)
tree.pred = predict(tree.fit, test, type="class")
with(test, table(tree.pred, Purchase))
```

**Result:**

```
> with(test, table(tree.pred, Purchase))
            Purchase
tree.pred   CH   MM
       CH  150   38
       MM   11   71
```

**Ans:**
The test error rate is $\frac{11+38}{270} = 0.18$

## (f)

**Code:**

```
1  # (f)
2  cv.fit = cv.tree(tree.fit)
3  cv.fit
4  cv.fit$size[which.min(cv.fit$dev)]
```

**Result:**

```
1  > cv.fit
2  $size
3  [1] 8 7 6 5 4 3 2 1
4
5  $dev
6  [1]  673.8589  677.7294  675.2341  729.4281  727.2145  772.9214
        773.5642
7  [8]  1067.3829
8
9  $k
10 [1]       -Inf   11.87503   12.41171   29.77434   31.80546   39.82936
        41.38321
11 [8]  306.37571
12
13 $method
14 [1] "deviance"
15
16 attr(,"class")
17 [1] "prune"          "tree.sequence"
18 > cv.fit$size[which.min(cv.fit$dev)]
19 [1] 8
```

**Ans:**
The optimal tree size is 8.

## (g)

**Code:**

```
1  # (g)
2  par(mfrow = c(1, 2))
3  # default plot
4  plot(cv.fit)
5  # better plot
6  plot(cv.fit$size, cv.fit$dev / smp_size, type = "b",
7        xlab = "Tree Size", ylab = "CV Misclassification Rate")
```
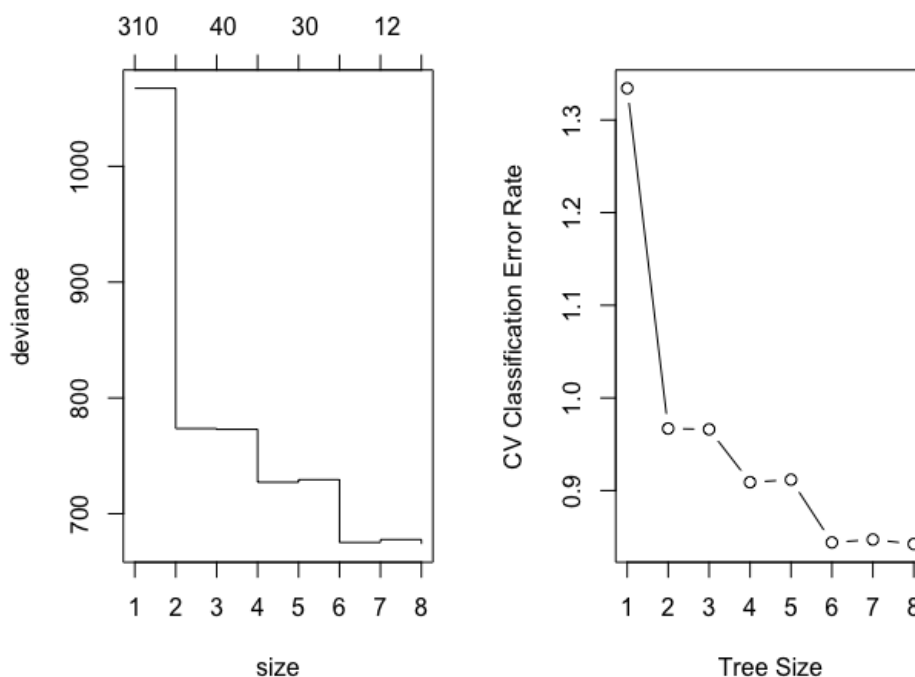
**Result:**

Figure 2: Error rate to tree size

## (h)

**Code:**

```
1  # (h)
2  cv.best = cv.fit$size[which.min(cv.fit$dev / smp_size)] #
      misclassification rate of each tree
3  cv.best
```

**Result:**

```
1  > cv.best
2  [1] 8
```

**Ans:**

The tree size corresponds to the lowest cross-validated classification error rate is 8.

## (i)

**Code:**

```
1  # (i)
2  prune.fit = prune.tree(tree.fit, best = 5)
3  plot(prune.fit)
4  text(prune.fit, pretty=0)
```
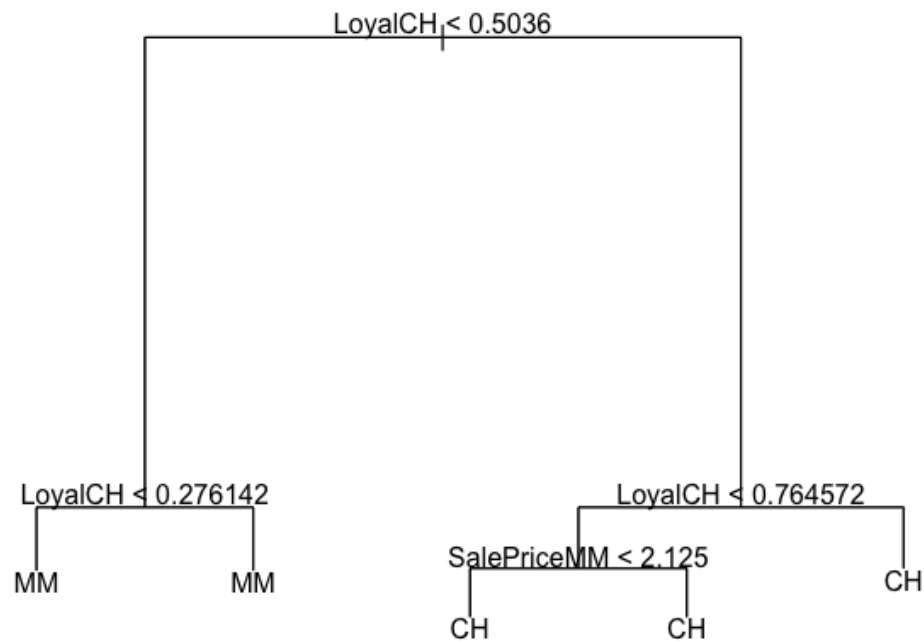
**Result:**

Figure 3: Pruned classification tree

## (j)

**Code:**

```
1  # (j)
2  summary(prune.fit)
```

**Result:**

```
1  > summary(prune.fit)
2
3   Classification tree:
4   snip.tree(tree = tree.fit, nodes = c(12L, 4L, 5L))
5   Variables actually used in tree construction:
6   [1] "LoyalCH"      "SalePriceMM"
7   Number of terminal nodes:   5
8   Residual mean deviance:   0.8138 = 646.9 / 795
9   Misclassification error rate: 0.1962 = 157 / 800
```

**Ans:**
Training error rate of unpruned tree: 128/800=0.16
Training error rate of pruned tree: 157/800=0.1962
The training error rate of pruned tree is higher.

## (k)

**Code:**

```
1  # (k)
```

```
2  tree.pred = predict(prune.fit, test, type="class")
3  with(test, table(tree.pred, Purchase))
```

**Result:**

```
1  > with(test, table(tree.pred, Purchase))
2              Purchase
3  tree.pred   CH   MM
4         CH  129   25
5         MM   32   84
```

**Ans:**

Test error rate of unpruned tree: $(11 + 38)/270 = 0.18$

Test error rate of pruned tree: $(32 + 25)/270 = 0.21$

The test error rate of pruned tree is higher.