

## Homework 3

---

Yi Yang (yy3089)

November 2, 2021

1 P1

1.1 a

$Y = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2$ . The possibility that  $P(Y = \text{true})$  is:

$$\begin{aligned} P(Y = \text{true}) &= \frac{e^{(\hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2)}}{1 + e^{(\hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2)}} \\ &= \frac{e^{(-6 + 0.05 * 40 + 1 * 3.5)}}{1 + e^{(-6 + 0.05 * 40 + 1 * 3.5)}} \\ &= 0.378 \end{aligned} \tag{1.1}$$

1.2 b

We can write the equation that:

$$\begin{aligned} P(Y = \text{true}) &= \frac{e^{(-6 + 0.05 * X_2 + 1 * 3.5)}}{1 + e^{(-6 + 0.05 * X_2 + 1 * 3.5)}} \\ &= 0.5 \end{aligned} \tag{1.2}$$

2 P2

$$\begin{aligned}
 P(Y = \text{div} \mid X = 4) &= \frac{P(Y = \text{div}, X = 4)}{P(X = 4)} \\
 &= \frac{P(X = 4 \mid Y = \text{div}) P(Y = \text{div})}{P(X = 4)} \\
 &= \frac{P(\text{div}) * p(4 \mid \text{div})}{P(\text{div}) * p(4 \mid \text{div}) + P(\text{not-div}) * p(4 \mid \text{not-div})} \\
 &= \frac{0.8 * \frac{1}{\sqrt{2\pi\hat{\sigma}^2}} e^{-\frac{(4-\mu_{\text{div}})^2}{2\hat{\sigma}^2}}}{0.8 * \frac{1}{\sqrt{2\pi\hat{\sigma}^2}} e^{-\frac{(4-\mu_{\text{div}})^2}{2\hat{\sigma}^2}} + 0.2 * \frac{1}{\sqrt{2\pi\hat{\sigma}^2}} e^{-\frac{(4-\mu_{\text{not-div}})^2}{2\hat{\sigma}^2}}} \\
 &= \frac{0.8 * 0.04032845}{0.8 * 0.04032845 + 0.2 * 0.05324133} \\
 &= 0.752
 \end{aligned} \tag{2.1}$$

### 3 P3

$P(Y = \text{true}) = \frac{e^{(\beta_0 + \beta_1 x_i)}}{1 + e^{(\beta_0 + \beta_1 x_i)}}$ ,  $P(Y = \text{false}) = \frac{1}{1 + e^{(\beta_0 + \beta_1 x_i)}}$ . Assume that  $k_i = 1$  if  $y_i = \text{true}$ , else  $k_i = 0$

$$\mathcal{L}(Y | \beta) = \prod_{i=1}^6 \left( \frac{e^{(\beta_0 + \beta_1 x_i)}}{1 + e^{(\beta_0 + \beta_1 x_i)}} \right)^{k_i} \left( \frac{1}{1 + e^{(\beta_0 + \beta_1 x_i)}} \right)^{1-k_i} \quad (3.1)$$

To find the optimal value of likelihood function, the question is equal to find the optimal value for:

$$\ln(\mathcal{L}(Y | \beta)) = \sum_{i=1}^6 k_i (\beta_0 + \beta_1 x_i) - \ln(1 + e^{(\beta_0 + \beta_1 x_i)}) \quad (3.2)$$

Which is equal to solve:

$$\begin{aligned} \frac{\partial l(\beta_0, \beta_1)}{\partial \beta_0} &= \sum_{i=1}^n \left( k_i - \frac{e^{(\beta_0 + \beta_1 x_i)}}{1 + e^{(\beta_0 + \beta_1 x_i)}} \right) = 0 \\ \frac{\partial l(\beta_0, \beta_1)}{\partial \beta_1} &= \sum_{i=1}^n \left( x_i k_i - \frac{x_i e^{(\beta_0 + \beta_1 x_i)}}{1 + e^{(\beta_0 + \beta_1 x_i)}} \right) = 0 \end{aligned} \quad (3.3)$$

According to Newton-Raphson Algorithm, a single Newton Update is:

$$\beta^{\text{new}} = \beta^{\text{old}} - \left( \frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} \right)^{-1} \frac{\partial \ell(\beta)}{\partial \beta} \quad (3.4)$$

The matrix form of this problem is:

$$\begin{aligned} \mathbf{z} &= \mathbf{X}\beta^{\text{old}} + \mathbf{W}^{-1}(\mathbf{y} - \mathbf{p}) \\ \beta^{\text{new}} &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z} \end{aligned} \quad (3.5)$$

```

1 library(matlib)
2
3 x = c(0.0, 0.2, 0.4, 0.6, 0.8, 1.0)
4 X = matrix(cbind(1, x), ncol = 2)
5 Y = matrix(c(0, 0, 0, 1, 0, 1), ncol = 1)
6
7 beta = matrix(c(0, 0), ncol = 1)
8
9 cat("Initial_Beta_values:", "beta0=", beta[1], ",",
10     "beta1=", beta[2], "\n")
11 # Initial Beta values: beta0 = 0, beta1 = 0
12
13 maxIter = 10
14
15 for(i in 1:maxIter){
16     p = exp(X %*% beta) / (1 + exp(X %*% beta))
17     w = as.vector(p * (1 - p))
18     W = diag(w)
19     z = X %*% beta + inv(W) %*% (Y - p)
20     beta = inv(t(X) %*% W %*% X) %*% t(X) %*% W %*% z

```

```

21     cat("Beta values at iter #:", i, "beta0=", beta[1], ",",
22         "beta1=", beta[2], "\n")
23 }
24 # Beta values at iter #: 1 beta0 = -2.380952 , beta1 = 3.428571
25 # Beta values at iter #: 2 beta0 = -3.522775 , beta1 = 4.966947
26 # Beta values at iter #: 3 beta0 = -4.022333 , beta1 = 5.624766
27 # Beta values at iter #: 4 beta0 = -4.096585 , beta1 = 5.721513
28 # Beta values at iter #: 5 beta0 = -4.09797 , beta1 = 5.723308
29 # Beta values at iter #: 6 beta0 = -4.09797 , beta1 = 5.723309
30 # Beta values at iter #: 7 beta0 = -4.09797 , beta1 = 5.723309
31 # Beta values at iter #: 8 beta0 = -4.09797 , beta1 = 5.723309
32 # Beta values at iter #: 9 beta0 = -4.09797 , beta1 = 5.723309
33 # Beta values at iter #: 10 beta0 = -4.09797 , beta1 = 5.723309

```

#### 4 P4

$\text{Cov}[\mathbf{Y}]$  can be written as:

$$\begin{aligned}
 \text{Cov}[\mathbf{Y}] &= \text{Cov}[\mathbf{A}\mathbf{X}] \\
 &= \mathbb{E}[(\mathbf{A}\mathbf{X} - \mathbb{E}[\mathbf{A}\mathbf{X}])(\mathbf{A}\mathbf{X} - \mathbb{E}[\mathbf{A}\mathbf{X}])^T] \\
 &= \mathbf{A}\mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^T]\mathbf{A}^T \\
 &= \mathbf{A}\text{Cov}[\mathbf{X}]\mathbf{A}^T \\
 &= \mathbf{A}\Sigma\mathbf{A}^T
 \end{aligned} \tag{4.1}$$

To get the eigenvalues of the matrix, we need to compute that:

$$\begin{aligned}
 |\lambda E - \Sigma| &= (\lambda - \sigma_1^2)(\lambda - \sigma_2^2) - \rho^2\sigma_1^2\sigma_2^2 \\
 &= \lambda^2 - (\sigma_1^2 + \sigma_2^2)\lambda + (1 - \rho^2)(\sigma_1^2 + \sigma_2^2) \\
 \lambda &= \frac{\sigma_1^2 + \sigma_2^2 \pm \sqrt{(\sigma_1^2 + \sigma_2^2)^2 - 4(1 - \rho^2)\sigma_1^2\sigma_2^2}}{2}
 \end{aligned} \tag{4.2}$$

Thus, matrix  $\Sigma$  and  $\mathbf{A}\Sigma\mathbf{A}^T$  can be written as:

$$\begin{aligned}
 \Sigma &= V\text{diag}(\lambda)V^T \\
 \mathbf{A}\Sigma\mathbf{A}^T &= \mathbf{A}V\text{diag}(\lambda)V^T\mathbf{A}^T \\
 &= \mathbf{A}V\text{diag}(\lambda)^{1/2}((\lambda)^{1/2})^T\mathbf{A}^T
 \end{aligned} \tag{4.3}$$

Where  $V$  is the eigenvector of matrix  $\Sigma$ . if we let  $\mathbf{A} = (V\text{diag}(\lambda)^{1/2})^{-1}$ , we can get that  $\text{Cov}[\mathbf{Y}] = \mathbf{A}\Sigma\mathbf{A}^T = E$

## 5 P5

As  $\hat{\sigma}^2$  is the unbiased estimation of  $\sigma^2$ , we can get:

$$\begin{aligned} E[\hat{\sigma}^2] &= E\left[\sum_{k=1}^K \alpha_k \hat{\sigma}_k^2\right] = \sum_{k=1}^K \alpha_k E[\hat{\sigma}_k^2] \\ &= \sigma^2 \end{aligned} \quad (5.1)$$

Also, according to Gaussian assumption:

$$\begin{aligned} \frac{(n_k - 1)\hat{\sigma}_k^2}{\sigma^2} &\sim \chi_{n_k-1}^2 \\ \text{Var}\left[\frac{(n_k - 1)\hat{\sigma}_k^2}{\sigma^2}\right] &= \text{Var}[\chi_{n_k-1}^2] \\ \frac{(n_k - 1)^2}{\sigma^4} \text{Var}[\hat{\sigma}_k^2] &= 2(n_k - 1) \\ \text{Var}[\hat{\sigma}_k^2] &= \frac{2\sigma^4}{n_k - 1} \end{aligned} \quad (5.2)$$

Thus:

$$\begin{aligned} \text{Var}[\hat{\sigma}^2] &= \text{Var}\left[\sum_{k=1}^K \alpha_k \hat{\sigma}_k^2\right] \\ &= \sum_{k=1}^K \alpha_k^2 \text{Var}[\hat{\sigma}_k^2] \\ &= \sum_{k=1}^K \alpha_k^2 \frac{2\sigma^4}{n_k - 1} \\ &= 2\sigma^4 \sum_{k=1}^K \frac{\alpha_k^2}{n_k - 1} \end{aligned} \quad (5.3)$$

To find the value of  $\alpha_k$  that minimizes  $\text{Var}[\hat{\sigma}^2]$ , the problem is equal to find the  $\alpha_k$  that minimizes  $\sum_{k=1}^K \frac{\alpha_k^2}{n_k - 1}$ . According to the Lagrange Multiplier Method, we can set up the equations to find the optimal value:

$$\mathcal{L}(\alpha_k) = \sum_{k=1}^K \frac{\alpha_k^2}{n_k - 1} + \lambda \left( \sum_{k=1}^K \alpha_k - 1 \right) \quad (5.4)$$

For each  $\alpha_k$ , we need to guarantee that:

$$\begin{aligned} \frac{\partial \mathcal{L}(\alpha_k)}{\partial \alpha_k} &= \frac{2\alpha_k}{n_k - 1} + \lambda = 0 \\ \alpha_k &= -\frac{\lambda}{2} * (n_k - 1) \end{aligned} \quad (5.5)$$

Also, note the constrain that  $\sum_{k=1}^K \alpha_k = 1$ . We can calculate that  $\alpha_k = \frac{n_k - 1}{n - K}$ .

## 6 P6

### 6.1 Majority Vote

```
1 samples = c(0.1, 0.15, 0.2, 0.2, 0.55, 0.6, 0.6, 0.65, 0.7, 0.75)
2
3 red = 0
4 green = 0
5
6 for(i in 1:length(samples)) {
7   if(samples[i] >= 0.5) {
8     red = red + 1
9   }
10  else {
11    green = green + 1
12  }
13 }
14 red
15 # [1] 6
16 green
17 # [1] 4
```

### 6.2 Average Probability

```
1 samples = c(0.1, 0.15, 0.2, 0.2, 0.55, 0.6, 0.6, 0.65, 0.7, 0.75)
2
3 mean(samples) >= 0.5
4 # False
```

### 6.3 Conclusion

The result of Majority Vote indicate that the answer should be red. The result of Average Probability indicate that the answer should be green.

## 7 P7

### 7.1 a

```
1 library(ISLR)
2 library(tree)
3 attach(OJ)
4 set.seed(1000)
5
6 train = sample(dim(OJ)[1], 800)
7 OJ.train = OJ[train, ]
8 OJ.test = OJ[-train, ]
```

### 7.2 b

```
1 OJ.tree = tree(Purchase ~ ., data = OJ.train)
2 summary(OJ.tree)
3 # Classification tree:
4 # tree(formula = Purchase ~ ., data = OJ.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
```

The tree uses LoyalCH, PriceDiff and salePriceMM. The training error rate is 0.16. The tree has 8 terminal nodes.

### 7.3 c

```
1 OJ.tree
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 ) *
```

I want to take terminal node 7 as an example. The variable got in this terminal is LoyalCH. Value of this variable is 0.764572. There are 260 nodes in the sub-trees of this node. Deviance of these nodes is 91.11. The predicted y-value for Purchase of this node is CH. There are 0.95769 of all



nodes in the sub-trees giving a prediction of CH for Purchase while 0.04231 of the nodes giving a prediction of MM. The star marks this node as a terminal node.

#### 7.4 d

```
1 plot(OJ.tree)
2 text(OJ.tree)
3 title("Tree")
```

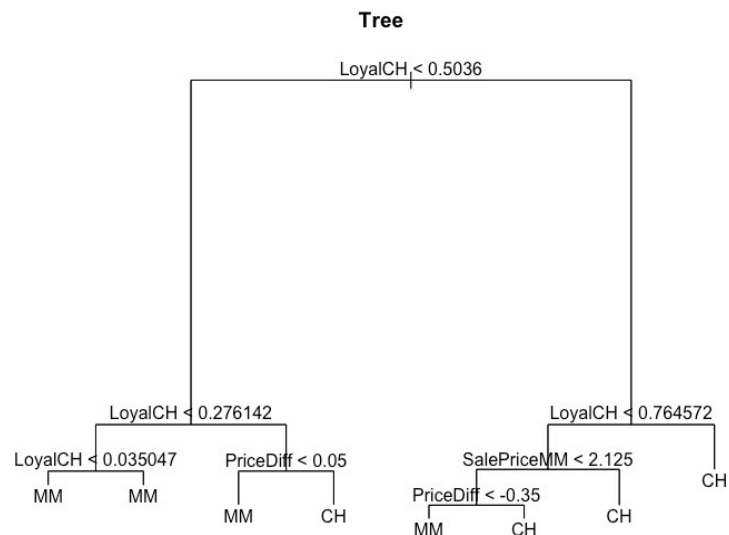


Figure 7.1: Screenshot for 'Tree Model'

#### 7.5 e

```
1 OJ.pred = predict(OJ.tree, OJ.test, type = "class")
2 table(OJ.test$Purchase, OJ.pred)
3 #      OJ.pred
4 #      CH  MM
5 #  CH 150  11
6 #  MM  38  71
7 OJ.error = sum(OJ.test$Purchase != OJ.pred)/dim(OJ.test)[1]
8 OJ.error
9 # [1] 0.1814815
```

#### 7.6 f

```
1 OJ.cvtree = cv.tree(OJ.tree, FUN = prune.tree)
2 OJ.cvtree
3 # $size
4 # [1] 8 7 6 5 4 3 2 1
5 #
6 # $dev
```

```

7 # [1] 663.1416 670.8785 669.0620 715.5108 737.4116 777.8179
  770.9825 1068.1253
8 #
9 # $k
10 # [1] -Inf 11.87503 12.41171 29.77434 31.80546 39.82936
   41.38321 306.37571
11 #
12 # $method
13 # [1] "deviance"
14 #
15 # attr("class")
16 # [1] "prune" "tree.sequence"

```

The optimal size is 8.

### 7.7 g

```

1 plot(OJ.cvtree$size, OJ.cvtree$dev, pch = 20, type = "o", col = "blue",
2       xlab = "Tree Size", ylab = "Dev")
3 title("Dev-Size Curve")

```

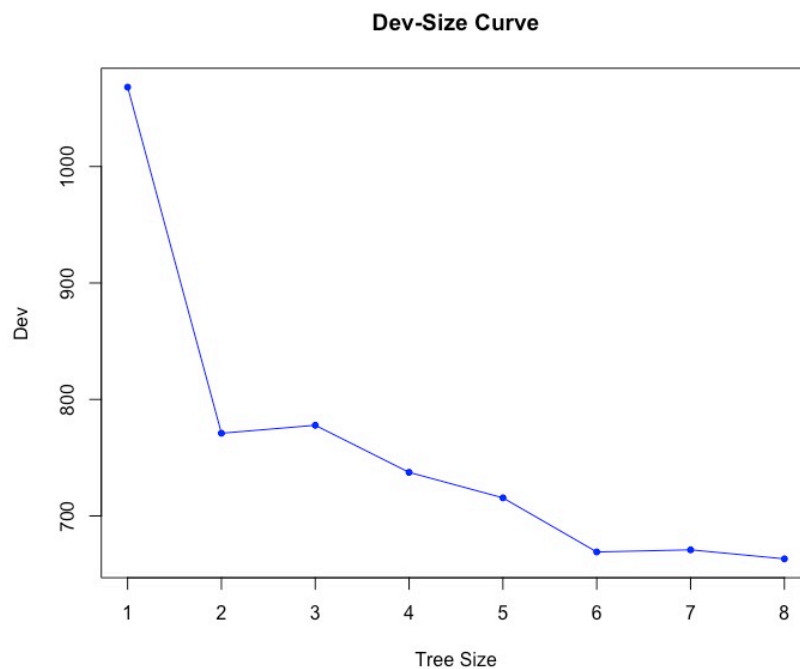


Figure 7.2: Screenshot for 'Dev-Size Curve'

### 7.8 h

```

1 OJ.cvtree$size[which.min(OJ.cvtree$dev)]
2 # [1] 8

```

When the size is 8, we can get the lowest cross-validation error.

### 7.9 i

```

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

```

## 7.10 j

```

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

```

The pruned tree selects LoyalCH, PriceDiff and SalePriceMM. It has 8 terminal nodes. The test error does not change.

## 7.11 k

```

1 OJ.prunedpred = predict(OJ.pruned, OJ.test, type = "class")
2 table(OJ.test$Purchase, OJ.prunedpred)
3 #      OJ.prunedpred
4 #      CH  MM
5 # CH 150  11
6 # MM  38  71
7 OJ.prunederror = sum(OJ.test$Purchase != OJ.prunedpred)/dim(OJ.test)[1]
8 OJ.prunederror
9 # [1] 0.1814815

```

The pruned test error is the same with the original test error.