

```
1 ### init.R
2
3 library(MASS)
4 mush0 = read.csv("mydata.csv")
5 mush = mush0[,-17]
6
7
8 ### util.R
9
10 cmp = function(data, cIndex) {
11   print(paste("colname:", names(data)[cIndex], " [edible,poisonous]"))
12   print(summary(data[which(data$class=="e"),cIndex]))
13   print(summary(data[which(data$class=="p"),cIndex]))
14 }
15
16
17 ### myclass.R
18
19 #myclass = rep("?", nrow(mush))
20
21 #for (i in 1:nrow(mush))
22 #{
23 # if (mush[i,]$odor %in% c("a","l") | mush[i,]$spore.print.color %in% c("b",
24 # {
25 #   myclass[i] = "e"
26 # }
27 # else if (mush[i,]$odor %in% c("c","f","m","p","s","y") | mush[i,]$spore.pr
28 # {
29 #   myclass[i] = "p"
30 # }
31 #}
32
33 myclass = ifelse(mush$odor %in% c("a","l") | mush$spore.print.color %in% c("
34   ifelse(mush$odor %in% c("c","f","m","p","s","y") | mush$spore.print.color
35     ifelse(mush$spore.print.color != "w", "e",
36       ifelse(mush$gill.size=="b", "e",
37         ifelse(mush$gill.spacing=="c" | mush$stalk.surface.above.ring=="k",
38           ifelse(mush$population=="c", "p",
39             "?"
40           )
41         )
42       )
43     )
44   )
45 )
```

```
46 myclass = factor(myclass)
47 foo = mush[which(myclass=="?"),]
48
49 print("-----")
50 for (i in (1:ncol(mush))[-c(1,17)])
51 {
52   cmp(foo,i)
53 }
54
55 print(table(myclass, mush$class))
56
57 ### br.R
58
59 start = Sys.time()
60 print(start)
61
62 # Constants
63 nCol = ncol(mush)
64 index = 1:nrow(mush)
65 k = 5
66
67 # Feature value summary table setup.
68 setsize = 1:nCol
69 for (i in 1:nCol)
70 {
71   setsize[i] = length(unique(mush[,i]))
72 }
73 maxFcmt = max(setsize)
74
75 # Results structures.
76 preds = data.frame(mush$class,mush$class,mush$class,mush$class,mush$class,rep(NA,nCol))
77 names(preds) = c("real", "nbc", "lda", "qda", "log.reg", "fold")
78
79 # Set of test features. (not including class)
80 testFeatureSet = 2:nCol
81
82 for (fold in 0:(k-1))
83 {
84   # Partition Data for NBC
85   test = mush[which(index %% k == fold),]
86   train = mush[which(index %% k != fold),]
87   traine = train[which(train$class == "e"),]
88   trainp = train[which(train$class == "p"),]
89
90   # Initialize results structures.
```

```
91 test.preds = preds[which(index %% k == fold),]
92 test.preds$fold = fold
93 test.preds$real = test$class
94
95 pe = log(nrow(traine)) # P(edible)
96 pp = log(nrow(trainp)) # P(poisonous)
97
98 # Create concise matrix of log(P(feature value|class)) estimates
99 pfe = matrix(nrow=nCol, ncol=maxFcnt)
100 pfp = matrix(nrow=nCol, ncol=maxFcnt)
101 eFeatValLogit = matrix(nrow=nCol, ncol=maxFcnt)
102 for (c in 1:nCol)
103 {
104   l = levels(mush[1,c])
105   for (v in 1:length(l))
106   {
107     valueTotal = max(log(sum(train[,c] == l[v])), 0) * 100
108     pfe[c,v] = max(log(sum(traine[,c] == l[v])), -valueTotal)
109     pfp[c,v] = max(log(sum(trainp[,c] == l[v])), -valueTotal)
110   }
111 }
112 # Logit(ish) Estimate for numerical methods
113 eFeatValLogit = pfe - pfp # Matrix subtraction
114
115 # Create probabalistic data table.
116 mush.eprob = data.frame(mush[,1], rep(list(rep(-1, nrow(mush))), nCol))
117 names(mush.eprob) = names(mush)
118 for (c in 2:nCol)
119 {
120   # Fancy vector method for setting a column of values at a time.
121   mush.eprob[,c] = eFeatValLogit[(as.integer(mush[,c])-1) * nCol + c]
122 }
123 test.eprob = mush.eprob[which(index %% k == fold),]
124 train.eprob = mush.eprob[which(index %% k != fold),]
125
126 # NBC
127 res = rep("?", nrow(test))
128 for (ti in 1:nrow(test))
129 {
130   pi.e = pe + sum(pfe[(as.integer(test[ti, testFeatureSet])-1)*nCol + testF
131   pi.p = pp + sum(pfp[(as.integer(test[ti, testFeatureSet])-1)*nCol + testF
132
133   res[ti] = ifelse(pi.e > pi.p, "e", "p")
134 }
135 test.preds$ NBC = factor(res, levels=levels(mush[1,1]))
```

```
136
137 f = as.formula(paste("class~",paste(names(mush)[testFeatureSet], collapse="
138 # LDA
139 lda.fit = lda(f, data=train.eprob)
140 test.preds$lda = predict(lda.fit, test.eprob)$class
141
142 # QDA
143 qda.fit = qda(f, data=train.eprob)
144 test.preds$qda = predict(qda.fit, test.eprob)$class
145
146 # Logistic Regression
147 glm.fit = glm(f, family=binomial, data=train.eprob)
148 glm.prob = predict(glm.fit, test.eprob, type="response")
149 test.preds$log.reg = factor(ifelse(glm.prob < .5, "e", "p"), levels=levels
150
151 preds[which(index %% k == fold),] = test.preds
152
153 }
154 print(table(preds$nb, preds$real))
155 print(table(preds$lda, preds$real))
156 print(table(preds$qda, preds$real))
157 print(table(preds$log.reg, preds$real))
158 print(Sys.time() - start)
159
160
161 ### nb.R
162
163 mylog = log
164
165 index = 1:nrow(mush)
166 k = 5
167 caution = .5
168 cautionOff = 0
169 cautionRange = .01
170 setsize = 1:ncol(mush)
171 for (i in 1:ncol(mush))
172 {
173   setsize[i] = length(unique(mush[,i]))
174 }
175 maxFcmt = max(setsize)
176
177 pred = factor(c("e", "p"))
178 real = factor(c("e", "p"))
179 confuse = table(pred, real)
180 confuse[,] = 0
```

```
181
182 coef = matrix(nrow=nrow(mush), ncol=2)
183 pprob = rep(0, nrow(mush))
184 eprob = pprob
185 realClass = rep('?', nrow(mush))
186 coefi = 0
187
188 for (fold in 0:(k-1))
189 {
190   test = mush[which(index %% k == fold),]
191   train = mush[which(index %% k != fold),]
192   traine = train[which(train$class == "e"),]
193   trainp = train[which(train$class == "p"),]
194
195   pe = mylog(nrow(traine))# * 2
196   pp = mylog(nrow(trainp))# * 2
197
198   pfe = matrix(nrow=ncol(mush), ncol=maxFcnt)
199   pfp = pfe
200   for (c in 1:ncol(mush))
201   {
202     l = levels(mush[,c])
203     for (v in 1:length(l))
204     {
205       emptyWeight = 0#-max(mylog(sum(train[,c] == l[v])), 0)
206       pfe[c,v] = max(mylog(sum(traine[,c] == l[v])), emptyWeight)
207       pfp[c,v] = max(mylog(sum(trainp[,c] == l[v])), emptyWeight)
208     }
209   }
210
211   res = rep("?", nrow(test))
212   for (ti in 1:nrow(test))
213   {
214     pi.e = pe
215     pi.p = pp
216     #print(paste(pi.e, pi.p))
217     #for (fi in 2:ncol(mush))
218     for (fi in c(6,21))
219     {
220       pi.e = pi.e + pfe[fi,as.integer(test[ti,fi])]
221       pi.p = pi.p + pfp[fi,as.integer(test[ti,fi])]
222     }
223
224     res[ti] = #ifelse(
225       #abs(pi.e - pi.p) < cautionRange
```

```
226     #abs(.5 - pi.e / (pi.e + pi.p)) < cautionRange
227     #, "?",
228     ifelse(pi.e * (1-caution) > pi.p * caution, "e", "p")
229 #)
230
231 #coef[coefi,] = c(pi.e / (pi.e + pi.p), ifelse(res[ti] == test[ti,1], 1,
232 pprob[coefi] = pi.e
233 eprob[coefi] = pi.p
234 realClass[coefi] = test[ti,1]
235 coefi = coefi + 1
236 }
237 resf = factor(res)#, levels=c(1,2,3), labels=c("e","p","?"))
238
239 curTab = table(res, test$class)
240 print(curTab)
241 confuse = confuse + curTab
242 }
243 print(confuse)
244
245
246
```