

```

1  ###決定木 (rpart:2進木解析) #####
2  ##wine data
3  white.wine <- ("https://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/winequality-white.csv")
4  dat.white <- read.table(white.wine, header = T, sep = ";")
5  str(dat.white)
6  colnames(dat.white)
7  red.wine <- ("https://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/winequality-red.csv")
8  dat.red <- read.table(red.wine, header = T, sep = ";")
9  str(dat.red)
10 colnames(dat.red)
11 anyNA(dat.white);anyNA(dat.red)
12 library(caret)
13 set.seed(123)
14 dat.w <- createDataPartition( dat.white$quality ,
15                               p = 0.75 , list = F )
16 train.w <- dat.white[ dat.w , ]
17 test.w <- dat.white[ -dat.w , ]
18 set.seed(123)
19 dat.r <- createDataPartition( dat.red$quality ,
20                               p = 0.75 , list = F )
21 train.r <- dat.red[ dat.r , ]
22 test.r <- dat.red[ -dat.r , ]
23 #回帰木
24 ##分類木;library(e1071)でチューニング
25 #minsplit (nodeの最少数),minbucket (leafの最少数,Default;minsplit/3)
26 #cp (複雑パラメータ)
27 #maxdepth (木の深さの最大数),xval (number of cross-validations)
28 options(scipen=50)
29 library(e1071)
30 tu.r <- tune.rpart(quality ~ ., data=train.r,
31                  cp = c(1e-3,2e-3,3e-3,4e-3,5e-3,6e-3,7e-3,9e-3))
32 summary(tu.r)
33 plot(tu.r)
34 dep.r <- tune.rpart(quality ~ ., data=train.r,
35                   maxdepth = 3:10)
36 summary(dep.r)
37 plot(dep.r)
38
39 #####
40 library(rpart)
41 library(partykit)
42 wr.out0 <- rpart(quality ~ ., data=train.r)
43 printcp(wr.out0)
44 sum(residuals(wr.out0)^2) #残差二乗和
45
46 #
47 cotr <- rpart.control(cp = 0.003 , maxdepth = 5 , xval=20)
48 wr.out <- rpart(quality ~ ., data=train.r ,
49               method = "anova",
50               parms = list(split = "gini"),#information:エントロピー
51               control = cotr)
52 #method:anova:回帰木,poisson:生起,class:分類木,exp:生存;指定した方がいい
53 names(wr.out) #出力内容
54 as.party(wr.out) #Fitと分割内容
55 wr.out$cptable
56 wr.out$variable.importance #重要度
57 plot(as.party(wr.out))
58 printcp(wr.out)
59 plotcp(wr.out)#CP.Graf
60 sum(residuals(wr.out)^2) #残差二乗和
61
62 #回帰推定による残差二乗和と比較
63 out.lm <- step(lm( quality ~ . ,
64                 data=train.r))#回帰推定
65 summary(out.lm)
66 sum(residuals(out.lm)^2) #回帰推定による残差二乗和
67 boxplot(data.frame(Rpart=residuals(wr.out),
68                   LM=residuals(out.lm))) #残差二乗和の比較
69 grid()
70 ##Test dataで検証#####
71 #predict (予測モデル,data)
72 p.lm.r <- predict(out.lm, test.r) #重回帰
73 p.rpart.r <- predict(wr.out, test.r)#回帰木
74 p.rpart.r0 <- predict(wr.out0, test.r)#回帰木 (デフォルト)
75 #RSME
76 sqrt( sum((test.r$quality - p.lm.r)^2) / 10788 ) #重回帰
77 sqrt( sum((test.r$quality - p.rpart.r)^2) / 10788 ) #回帰木
78 sqrt( sum((test.r$quality - p.rpart.r0)^2) / 10788 ) #回帰木 (デフォルト)
79 #MAE
80 mean(abs(test.r$quality - p.lm.r)) #重回帰
81 mean(abs(test.r$quality - p.rpart.r)) #回帰木
82 mean(abs(test.r$quality - p.rpart.r0)) #回帰木 (デフォルト)
83 #cor
84 cor(p.lm.r , test.r$quality) #重回帰
85 cor(p.rpart.r , test.r$quality) #回帰木
86 cor(p.rpart.r0 , test.r$quality) #回帰木 (デフォルト)
87 ##分類木#####
88 #第1列から48列:変数名に用いた文字列がメールに使用された頻度
89 #num857のようにnum***はその数値***が現れた頻度
90 #49列から54列:記号;、(、[、!、$、#の使用頻度

```

```

91 #55列から57列:用いられた大文字の平均、大文字が連続使用された
92 #最も長い文字列の文字数、用いられた大文字の総数
93 ##train.test data作成
94 library(kernlab);data(spam)
95 library(caret)
96 anyNA(spam)
97 set.seed(123)
98 dat.s <- createDataPartition( spam$type , p = 0.7 , list = F )
99 train.spam <- spam[ dat.s ,]
100 test.spam <- spam[-dat.s ,]
101 ##分類木;library(e1071)でチューニング
102 #minsplit(nodeの最少数),minbucket(leafの最少数,Default;minsplit/3)
103 #cp(複雑パラメータ)
104 #maxdepth(木の深さの最大数),xval(number of cross-validations)
105 options(scipen=50)
106 library(e1071)
107 tun <- tune.rpart(type ~ ., data=train.spam,
108                 cp = c(1e-4,1e-3,3e-3,5e-3,7e-3,1e-2))
109 summary(tun)
110 plot(tun)
111
112 msp <- tune.rpart(type ~ ., data=train.spam,
113                 minsplit=seq(10,50,10))
114 summary(msp)
115 plot(msp)
116
117 dep <- tune.rpart(type ~ ., data=train.spam,
118                 maxdepth = seq(3,10,1))
119 summary(dep)
120 plot(dep)
121
122 #####
123 library(rpart)
124 library(partykit)
125 cont <- rpart.control(minsplit = 30, minbucket = round(30/3),#minsplit/3
126                     cp = 0.001 , maxdepth = 6 ,xval=20)
127 spam.out <- rpart(type ~ ., data=train.spam , method = "class",
128                 parms = list(split = "gini"),#information:イントロピー
129                 control = cont)
130 names(spam.out) #出力内容
131 as.party(spam.out) #Fitと分割内容
132 spam.out$sctable
133 spam.out$variable.importance #重要度
134 plot(as.party(spam.out))
135 pd <- predict(spam.out , type = "class")
136 tb <- xtabs( ~ type + pd , data = train.spam )
137 tb
138 round((tb[1,1] + tb[2,2])/length(train.spam$type),3)*100
139 #confusion matrix
140 pd <- predict(spam.out , type = "class")
141 library(caret)
142 confusionMatrix(pd , train.spam$type )
143 ##Test dataで検証#####
144 #predict(予測モデル,data)
145 p.rpart.t <- predict(spam.out, test.spam , type = "class")
146 tb.t <- xtabs( ~ type + p.rpart.t , data = test.spam )
147 tb.t
148 round((tb.t[1,1] + tb.t[2,2])/length(test.spam$type),3)*100
149 library(caret)
150 confusionMatrix(p.rpart.t, test.spam$type )
151 ##判別分析#####
152 library(MASS)
153 lda.out <- lda(type ~ ., data=train.spam )
154 summary(lda.out)
155 predict.lda <- predict(lda.out)
156 tab <- table(train.spam$type , predict.lda$class)
157 round(((tab[1,1]+ tab[2,2])/sum(tab)*100),2)
158 predict.out.test <- predict(lda.out , newdata = test.spam)
159 tab.test <- table(test.spam$type , predict.out.test$class)
160 #rn.ics
161 round(((tab.test[1,1] + tab.test[2,2])/sum(tab.test)*100),2) #判別
162 round((tb.t[1,1] + tb.t[2,2])/length(test.spam$type),3)*100 #分類木
163
164 ##多クラス分類
165 ##手書き数 data
166 test.su <- ("https://archive.ics.uci.edu/ml/machine-learning-databases/optdigits//optdigits.tes")
167 test.su <- read.table(test.su , sep = ",")
168 train.su <- ("https://archive.ics.uci.edu/ml/machine-learning-databases/optdigits//optdigits.tra")
169 train.su <- read.table(train.su , sep = ",")
170 str(train.su)
171 colnames(train.su)
172 str(test.su)
173 colnames(test.su)
174 table(test.su$V65)
175 table(train.su$V65)
176 #y:factor変換 confusion matrix出力のため(factor変換:解析には不要)
177 train.su$V65 <- factor(train.su$V65)
178 str(train.su)
179 test.su$V65 <- factor(test.su$V65)
180 str(test.su)

```

```

181 colnames(train.su)
182 colnames(test.su)
183 table(test.su$V65)
184 table(train.su$V65)
185 ##
186 ##分類木;library(e1071)でチューニング
187 options(scipen=50)
188 library(e1071)
189 tue <- tune.rpart(V65 ~ ., data=train.su ,
190                  cp = c(1e-4,4e-4,7e-4,1e-3,3e-3,5e-3,7e-3,1e-2))
191 summary(tue)
192 plot(tue)
193 readline(tue)
194 spl <- tune.rpart(V65 ~ ., data=train.su ,
195                  minsplit=seq(10,80,20))
196 summary(spl)
197 plot(spl)
198 readline(spl)
199 mdp <- tune.rpart(V65 ~ ., data=train.su ,
200                  maxdepth = 3:10)
201 summary(mdp)
202 plot(mdp)
203 readline(mdp)
204 #####
205 library(rpart)
206 library(partykit)
207 cotr <- rpart.control(minsplit = 10, minbucket = round(10/3),#minsplit/3
208                      cp = 0.0004 , maxdepth = 8 ,xval=20)
209 suchi.out <- rpart(V65 ~ ., data=train.su , method = "class",
210                  parms = list(split = "gini"),#information:イントロピー
211                  control = cotr)
212 suchi.out
213 names(suchi.out) #出力内容
214 as.party(suchi.out) #Fitと分割内容
215 suchi.out$sctable
216 suchi.out$variable.importance #重要度
217 plot(as.party(suchi.out))
218 printcp(suchi.out)
219 plotcp(suchi.out) #CP.Graf
220 #クロス集計:的中率
221 est.su <- predict(suchi.out , type = "class")
222 est.su2 <- xtabs(~ V65 + est.su , data = train.su )
223 est.su2
224 #round(sum(diag(est.su2))/sum(table(train.su$V65)),4)*100
225 round(sum(diag(est.su2)) /sum(est.su2),4)*100
226 #confusion matrix
227 library(caret)
228 confusionMatrix(est.su, train.su$V65 )
229 ##Test dataで検証#####
230 #predict(予測モデル,data)
231 p.rpart.su <- predict(suchi.out, test.su , type = "class")
232 est.su2.2 <- xtabs(~ V65 + p.rpart.su , data = test.su )
233 est.su2.2
234 round(sum(diag(est.su2.2)) /sum(est.su2.2),4)*100
235 #confusion matrix
236 library(caret)
237 confusionMatrix(p.rpart.su, test.su$V65 )
238 #
239 options(scipen=0)
240

```