

P8106_yiminchen_secondaryanalysis

Yimin Chen (yc4195), Yang Yi (yy3307), Qingyue Zhuo (qz2493)

Contents

Import and data manipulation	1
Data visualization	2
Model training	2

Import and data manipulation

```
# Load recovery.RData environment
load("./recovery.Rdata")

dat %>% na.omit()

# dat1 draw a random sample of 2000 participants Uni:3307
set.seed(3307)

dat1 = dat[sample(1:10000, 2000),]

dat1 =
  dat1[, -1] %>%
  mutate(
    recovery_time = as.factor(
      case_when(recovery_time <= 30 ~ "long", recovery_time > 30 ~ "short")
    ),
    gender = as.factor(gender),
    race = as.factor(race),
    smoking = as.factor(smoking),
    hypertension = as.factor(hypertension),
    diabetes = as.factor(diabetes),
    vaccine = as.factor(vaccine),
    severity = as.factor(severity),
    study = as.factor(
      case_when(study == "A" ~ 1, study == "B" ~ 2, study == "C" ~ 3)
    )
  )

# dat2 draw a random sample of 2000 participants Uni:2493
set.seed(2493)

dat2 = dat[sample(1:10000, 2000),]

dat2 =
  dat2[, -1] %>%
```

```

mutate(
  recovery_time = as.factor(
    case_when(recovery_time <= 30 ~ "long", recovery_time > 30 ~ "short")
  ),
  gender = as.factor(gender),
  race = as.factor(race),
  smoking = as.factor(smoking),
  hypertension = as.factor(hypertension),
  diabetes = as.factor(diabetes),
  vaccine = as.factor(vaccine),
  severity = as.factor(severity),
  study = as.factor(
    case_when(study == "A" ~ 1, study == "B" ~ 2, study == "C" ~ 3)
  )
)

# Merged dataset with unique observation
covid_dat = rbind(dat1, dat2) %>%
  unique()

covid_dat2 = model.matrix(recovery_time ~ ., covid_dat)[, -1]

# Partition dataset into two parts: training data (70%) and test data (30%)
rowTrain = createDataPartition(y = covid_dat$recovery_time, p = 0.7, list = FALSE)

trainData = covid_dat[rowTrain, ]
testData = covid_dat[-rowTrain, ]

# matrix of predictors
x1 = covid_dat2[rowTrain,]
# vector of response
y1 = covid_dat$recovery_time[rowTrain]
# matrix of predictors
x2 = covid_dat2[-rowTrain,]
# vector of response
y2 = covid_dat$recovery_time[-rowTrain]

ctrl1 = trainControl(method = "repeatedcv", number = 10, repeats = 5)
ctrl2 = trainControl(method = "cv",
  classProbs = TRUE,
  summaryFunction = twoClassSummary)

```

Data visualization

Model training

classification - classification tree: L11 - glm + penalized logistic regression L8 - GAM L8 - MARS L8 - QDA
 L9 - LDA L9 - Naive Bayes L9 - random forest L12 - boosting L12 - support vector machines L13