

Midterm

Jeffrey Zhuohui Liang

3/7/2021

Method

Table 1: Data summary

Name	hrt_data
Number of rows	920
Number of columns	15
Column type frequency:	
factor	9
numeric	6
Group variables	None

Variable type: factor

skim_variable	n_missing	complete_rate	ordered	n_unique	top_counts
diagnosis_heart_disease	0	1.00	FALSE	2	pre: 509, abs: 411
location	0	1.00	FALSE	4	cle: 303, hun: 294, va: 200, swi: 123
sex	0	1.00	FALSE	2	mal: 726, fem: 194
chest_pain_type	0	1.00	FALSE	4	asy: 496, non: 204, aty: 174, typ: 46
fasting_blood_sugar	90	0.90	FALSE	2	fas: 692, fas: 138
resting_ecg	2	1.00	FALSE	3	nor: 551, lef: 188, ST-: 179
exercise_induced_angina	55	0.94	FALSE	2	no: 528, yes: 337
peak_exercise_st_segment	309	0.66	FALSE	3	Fla: 345, Up-: 203, Dow: 63
thalassemia	486	0.47	FALSE	3	nor: 196, rev: 192, fix: 46

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100
age	0	1.00	53.51	9.42	28.0	47	54.0	60.0	77.0
resting_blood_pressure	59	0.94	132.13	19.07	0.0	120	130.0	140.0	200.0
serum_cholesterol	30	0.97	199.13	110.78	0.0	175	223.0	268.0	603.0
max_heart_rate_achieved	55	0.94	137.55	25.93	60.0	120	140.0	157.0	202.0
st_depression_exercise	62	0.93	0.88	1.09	-2.6	0	0.5	1.5	6.2
num_major_vessels_flouro	611	0.34	0.68	0.94	0.0	0	0.0	1.0	3.0

```

# split into training set
train_index = createDataPartition(hrt_data$diagnosis_heart_disease,list = F)

Y_tr = hrt_data$diagnosis_heart_disease[train_index]

Y_ts = hrt_data$diagnosis_heart_disease[-train_index]

options(na.action = "na.pass")
X_tr = model.matrix(diagnosis_heart_disease ~., hrt_data,na.action = "na.pass")[train_index,-1]

X_ts = model.matrix(diagnosis_heart_disease ~., hrt_data,na.action = "na.pass")[-train_index,-1]

TRC = caret::trainControl(method = "repeatedcv",repeats=5,
                           summaryFunction = twoClassSummary,
                           classProbs = T)

PPS = c("knnImpute", "center", "scale")

```

Modeling

```

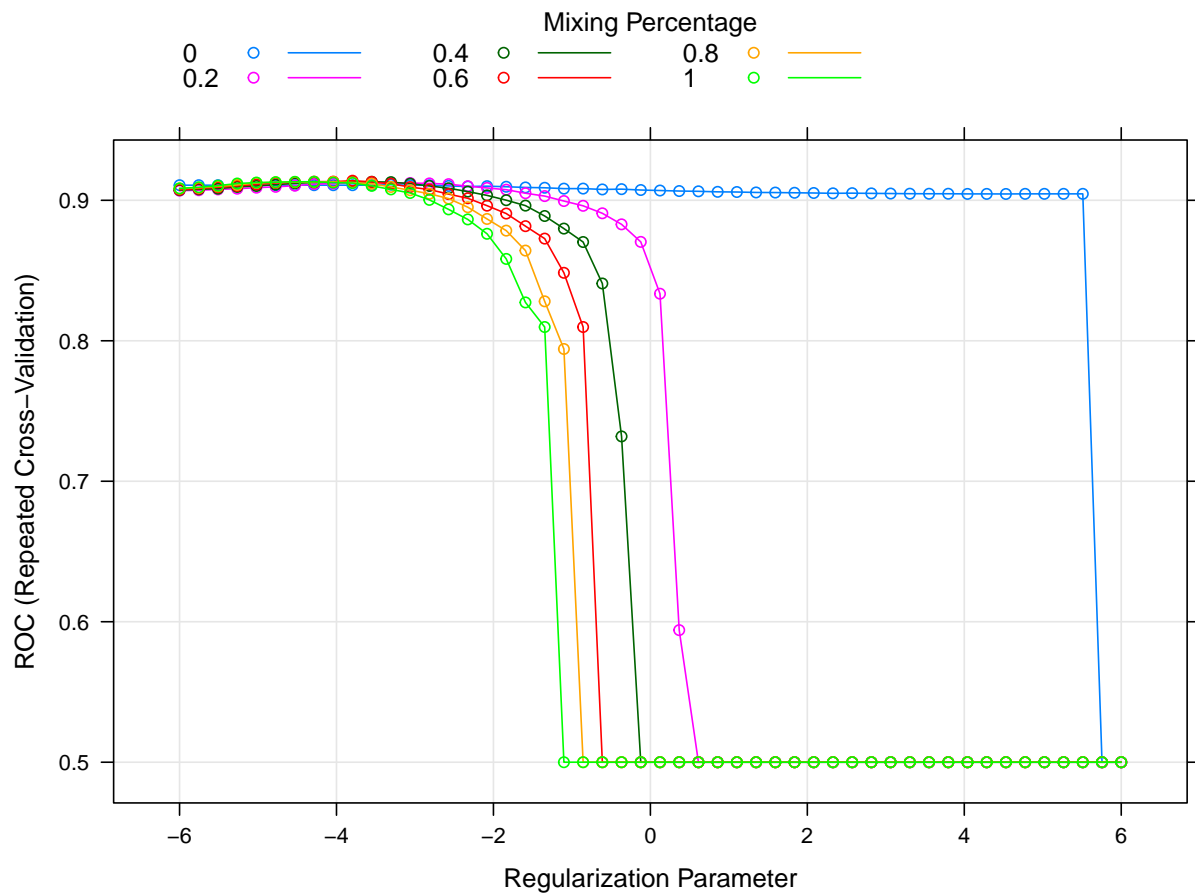
set.seed(123123)
cl = parallel::makePSOCKcluster(5)
doParallel::registerDoParallel(cl)

logistic_model =
  train(
    X_tr,
    Y_tr,
    method = "glmnet",
    tuneGrid = expand.grid(alpha = seq(0,1,length=6),
                           lambda = exp(seq(
                               6, to = -6, length = 50
                           ))),
    family = "binomial",
    preProcess = PPS,
    metric = "ROC",
    trControl = TRC
  )

stopCluster(cl)

plot(logistic_model,xTrans = function(x) log(x))

```



```
coef(logistic_model$finalModel,logistic_model$bestTune$lambda) %>%
  as.vector() %>%
  tibble(term = c("Intercept",colnames(X_tr)),
         coefficient = .) %>%
  knitr::kable(caption = "Coefficient of Lasso Logistic Regression")
```

Table 4: Coefficient of Lasso Logistic Regression

term	coefficient
Intercept	0.427
locationhun	0.000
locationswi	0.584
locationva	0.229
age	0.000
sexmale	0.224
chest_pain_typeatypical angina	-0.496
chest_pain_typenon-angina pain	-0.384
chest_pain_typetypical angina	0.000
resting_blood_pressure	0.000
serum_cholesterol	0.000
fasting_blood_sugarfasting blood sugar > 120 mg/dl	0.000
resting_ecgnormal	-0.025
resting_ecgST-T wave abnormality	0.109

term	coefficient
max_heart_rate_achieved	0.000
exercise_induced_anginayes	0.328
st_depression_exercise	0.223
peak_exercise_st_segmentFlat	0.092
peak_exercise_st_segmentUp-sloaping	-0.114
num_major_vessels_flouro	0.748
thalassemianormal	-0.613
thalassemiareversible defect	0.149

```

set.seed(123123)
cl = parallel::makePSOCKcluster(5)
doParallel::registerDoParallel(cl)

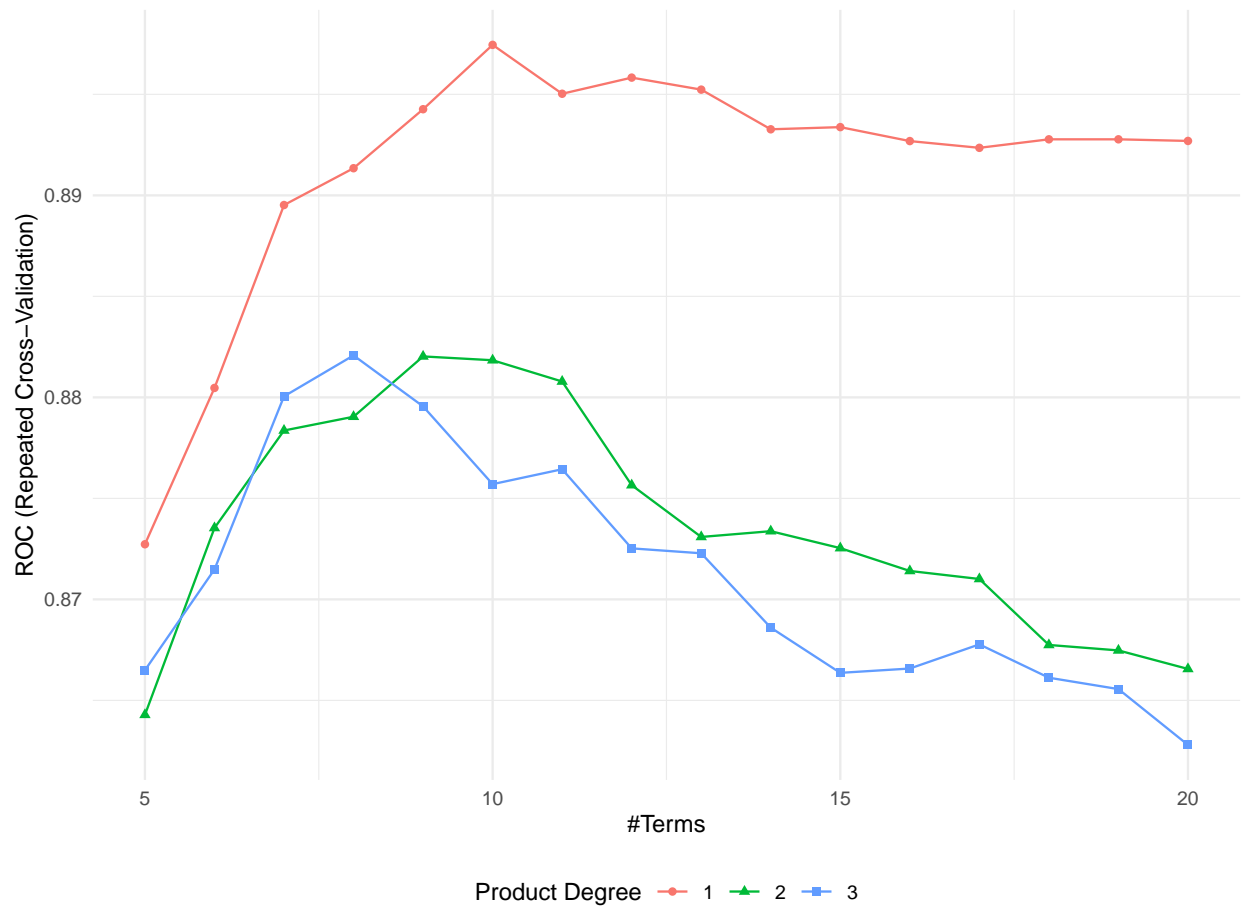
mars_model =
  train(X_tr,
        Y_tr,
        method = "earth",
        tuneGrid = expand.grid(degree = 1:3,
                               nprune = 5:20),

        preProcess = PPS,
        trControl = TRC,
        metric = "ROC")

stopCluster(cl)

ggplot(mars_model)

```

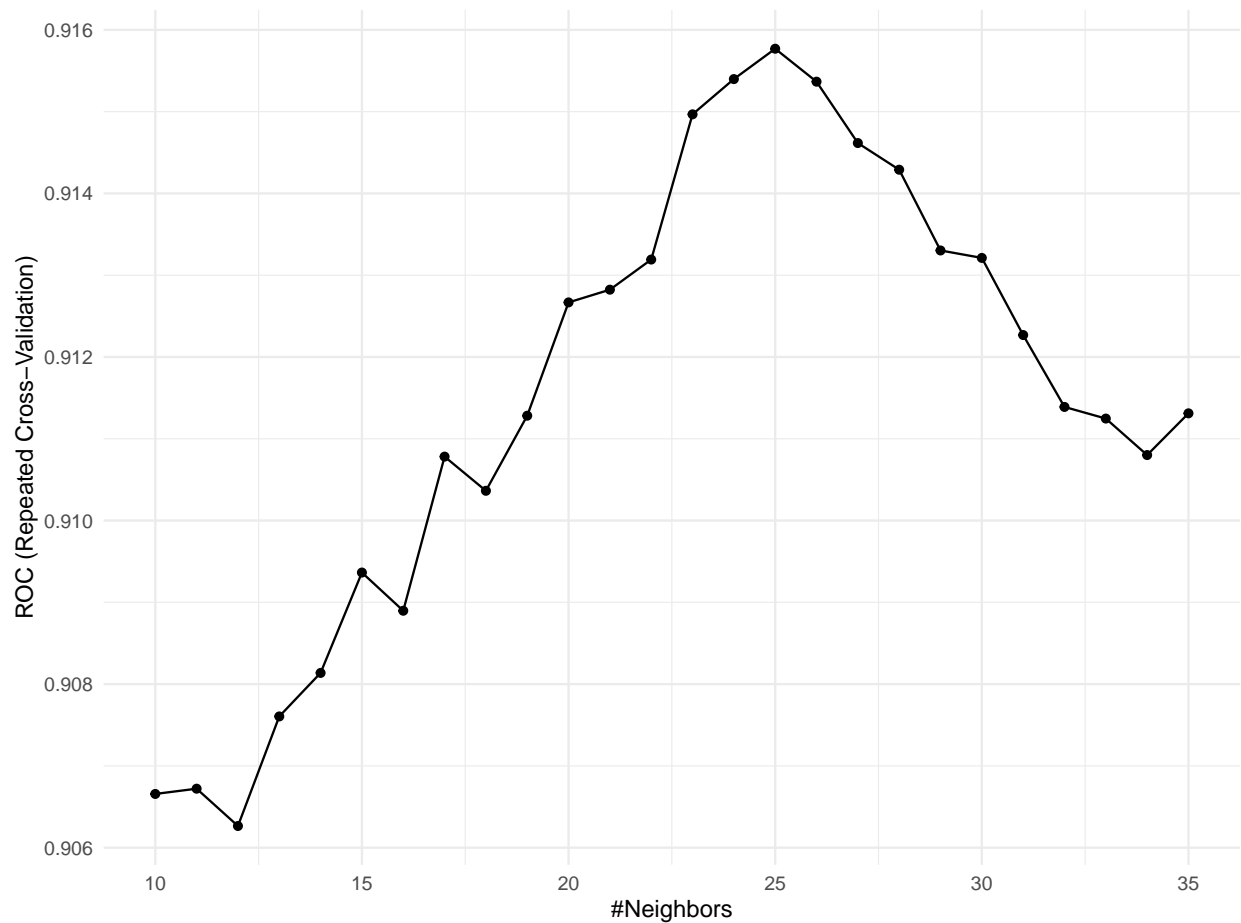


```
set.seed(123123)
cl = parallel::makePSOCKcluster(5)
doParallel::registerDoParallel(cl)

knn_model =
  train(X_tr,
        Y_tr,
        method = "knn",
        tuneGrid = expand.grid(k = 10:35),
        preProcess = PPS,
        trControl = TRC,
        metric = "ROC")

stopCluster(cl)

ggplot(knn_model)
```



```
set.seed(123123)
cl = parallel::makePSOCKcluster(5)
doParallel::registerDoParallel(cl)

lda_model = train(
  X_tr,
  Y_tr,
  method = "lda",
  preProcess = PPS,
  trControl = TRC,
  metric = "ROC"
)

stopCluster(cl)
```

```
set.seed(123123)
cl = parallel::makePSOCKcluster(5)
doParallel::registerDoParallel(cl)

qda_model = train(
  X_tr,
  Y_tr,
  method = "qda",
```

```

preProcess = PPS,
trControl = TRC,
metric = "ROC"
)

stopCluster(cl)

set.seed(123123)
cl = parallel::makePSOCKcluster(5)
doParallel::registerDoParallel(cl)

nb_model =
  train(
    X_tr,
    Y_tr,
    method = "nb",
    tuneGrid = expand.grid(
      usekernel = c(T, F),
      fL = 1,
      adjust = seq(.2, 3, by = .2)
    ),
    preProcess = PPS,
    trControl = TRC,
    metric = "ROC"
  )
stopCluster(cl)

```

Performance comparison

```

rsmp = resamples(
  list(
    logistic = logistic_model,
    MARS = mars_model,
    knn = knn_model,
    lda = lda_model,
    qda = qda_model
  ),
  metric = c("ROC", "Kappa")
)

summary(rsmp)

```

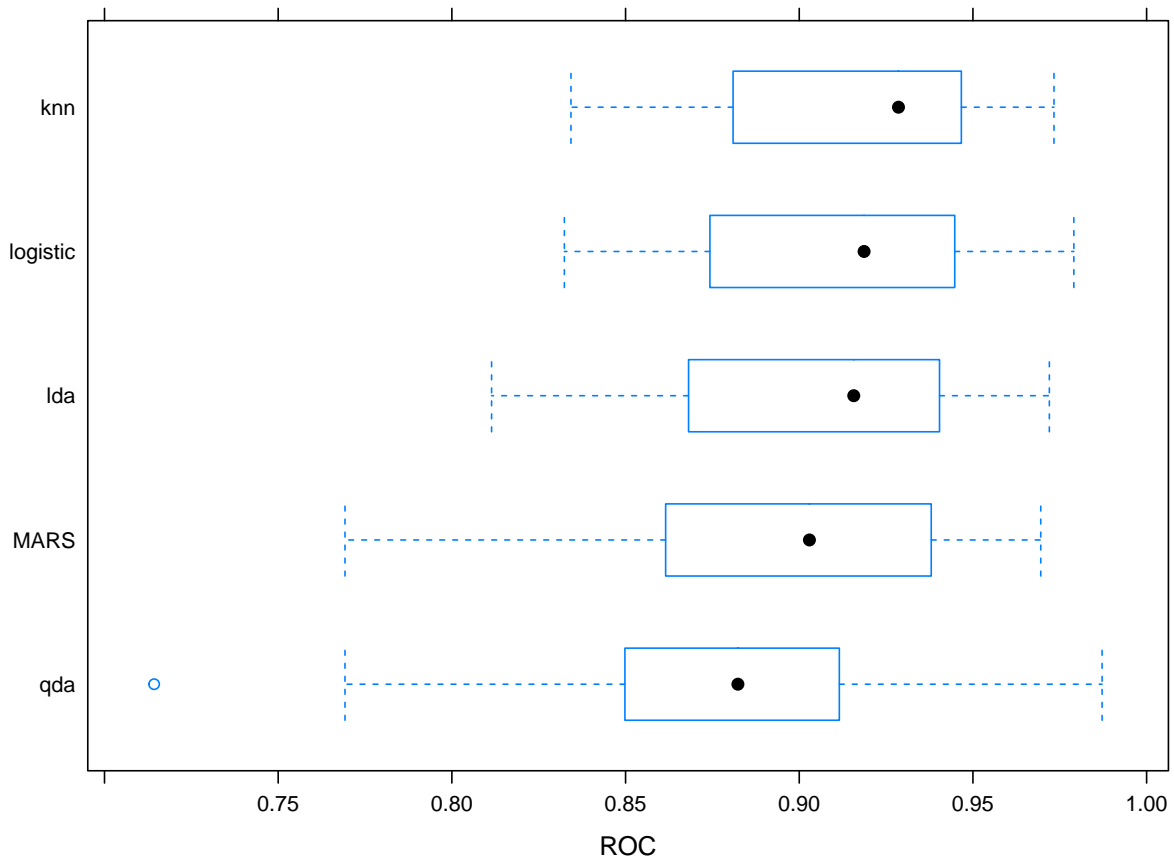
```

##
## Call:
## summary.resamples(object = rsmp)
##
## Models: logistic, MARS, knn, lda, qda
## Number of resamples: 50
##
## ROC
##           Min. 1st Qu. Median  Mean 3rd Qu.  Max. NA's

```

```
## logistic 0.832 0.876 0.919 0.911 0.944 0.979 0
## MARS 0.769 0.862 0.903 0.897 0.938 0.970 0
## knn 0.834 0.881 0.929 0.916 0.947 0.973 0
## lda 0.811 0.869 0.916 0.908 0.940 0.972 0
## qda 0.714 0.851 0.882 0.879 0.910 0.987 0
##
## Sens
##      Min. 1st Qu. Median Mean 3rd Qu. Max. NA's
## logistic 0.619 0.753 0.805 0.807 0.857 0.952 0
## MARS 0.619 0.750 0.800 0.792 0.850 0.952 0
## knn 0.619 0.750 0.800 0.800 0.857 0.952 0
## lda 0.619 0.762 0.810 0.814 0.857 0.952 0
## qda 0.571 0.750 0.800 0.791 0.840 0.952 0
##
## Spec
##      Min. 1st Qu. Median Mean 3rd Qu. Max. NA's
## logistic 0.680 0.810 0.846 0.855 0.885 1 0
## MARS 0.600 0.808 0.880 0.856 0.920 1 0
## knn 0.760 0.855 0.885 0.889 0.920 1 0
## lda 0.680 0.808 0.863 0.858 0.920 1 0
## qda 0.692 0.808 0.846 0.855 0.911 1 0
```

```
bwplot(rsmp,metric = "ROC")
```




```

ROC =
  expand.grid(
    test_X = list(X_ts),
    test_Y = list(Y_ts),
    model = list(logistic_model, mars_model, knn_model, lda_model, qda_model)
  ) %>%
  mutate(
    pred = map2(model, test_X, ~ predict(.x, newdata = .y, type = "prob")[, 2]),
    roc = map2(test_Y, pred, ~ pROC::roc(.x, .y))
  ) %>%
  pull(roc)

auc = c()

for (i in 1:5){
  auc = append(auc, ROC[[i]]$auc[1])
  plot(ROC[[i]], col = i, add = T * (i>1), legacy.axes = T * (i==1))
}

model_name =
  c("lasso logistic", "MARS", "KNN", "LDA", "QDA")

legend("bottomright",
  legend = paste0(model_name, "~", round(auc, 3)), col=1:5, lwd=2)

```

