

Diabetes Prediction model

DS II Final team

Contents

Load Data	2
EDA	3
Summary statistics	3
Density plots (numerical covariates)	7
Bar plots (categorical covariates)	8
Partition plots	9
Models	12
training data	12
Nonlinear models	13

```
library(RNHANES)
library(tidyverse)
library(summarytools)
library(leaps)
library(readr)
library(caret)
library(ggplot2)
library(patchwork)
library(mgcv)
library(nlme)
library(dplyr)
library(plyr)
library(AppliedPredictiveModeling)
library(dplyr)
library(scales)
library(pROC)
#library(MASS)
#library(klaR)
```

Load Data

```
data_files <- nhanes_load_data(file_name = "DIQ_H", year = "2013-2014")

data_files <- data_files %>%
  left_join(nhanes_load_data("HDL_H", "2013-2014"), by = "SEQN") %>%
  left_join(nhanes_load_data("INS_H", "2013-2014"), by = "SEQN") %>%
  left_join(nhanes_load_data("TRIGLY_H", "2013-2014"), by = "SEQN") %>%
  left_join(nhanes_load_data("DEMO_H", "2013-2014"), by = "SEQN") %>%
  left_join(nhanes_load_data("BMX_H", "2013-2014"), by = "SEQN") %>%
  left_join(nhanes_load_data("OGTT_H", "2013-2014"), by = "SEQN") %>%
  left_join(nhanes_load_data("BPX_H", "2013-2014"), by = "SEQN") %>%
  left_join(nhanes_load_data("PAQ_H", "2013-2014"), by = "SEQN") %>%
  left_join(nhanes_load_data("DPQ_H", "2013-2014"), by = "SEQN") %>%
  left_join(nhanes_load_data("SLQ_H", "2013-2014"), by = "SEQN")

raw_data <- data_files %>%
  select(SEQN, RIAGENDR, RIDAGEYR, RIDRETH3, BMXBMI, LBDHDD, LBDLDL, LBXTR, LBXIN, LBXGLT, BPXSY1, BPXD1)

raw_data <- raw_data[raw_data$DIQ010 != 3 & raw_data$DIQ010 != 7 & raw_data$DIQ010 != 9, ] %>% mutate(
  drop_na(DIQ010)

colnames(raw_data) <- c("ID", "gender", "age", "race", "bmi", "hdl", "ldl", "triglyceride", "insulin",
  contrasts(raw_data$diabetes)

2 1 0 2 1

levels(raw_data$diabetes)[1] <- "yes"
levels(raw_data$diabetes)[2] <- "no"
contrasts(raw_data$diabetes)
```

no

yes 0 no 1

```
write.csv(raw_data, "final_data.csv")
```

EDA

Summary statistics

```
st_options(plain.ascii = FALSE,
           style = "rmarkdown",
           dfSummary.silent = TRUE,
           footnote = NA,
           subtitle.emphasis = FALSE)

dfSummary(raw_data[, -1], valid.col = FALSE)
```

Data Frame Summary

raw_data
Dimensions: 9578 x 18
Duplicates: 319

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Missing
1	gender [factor]	1. 1	4706 (49.1%)	IIIIIIII	0
		2. 2	4872 (50.9%)	IIIIIIIIII	(0.0%)
2	age [numeric]	Mean (sd) : 32.4 (23.9)	80 distinct values	:	0
		min < med < max:		::	(0.0%)
		1 < 28 < 80		:::	
		IQR (CV) : 41 (0.7)		:::::::::::	
3	race [factor]	1. 1	1616 (16.9%)	III	0
		2. 2	893 (9.3%)	I	(0.0%)
		3. 3	3449 (36.0%)	IIIIIIII	
		4. 4	2148 (22.4%)	IIII	
		5. 6	1033 (10.8%)	II	
		6. 7	439 (4.6%)		
4	bmi [numeric]	Mean (sd) : 25.6 (7.9)	436 distinct	:	706
		min < med < max:	values	:::	(7.4%)
		12.1 < 24.6 < 82.9		:::	
		IQR (CV) : 10.4 (0.3)		:::.	
5	hdl [numeric]	Mean (sd) : 53.2 (15.2)	116 distinct	:	2128
		min < med < max:	values	:	(22.2%)
		10 < 51 < 173		:::	
		IQR (CV) : 19 (0.3)		:::	
				:::.	

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Missing
15	married [factor]	1. 1	2866 (51.3%)	IIIIIIIIII	3986
		2. 2	419 (7.5%)	I	(41.6%)
		3. 3	637 (11.4%)	II	
		4. 4	170 (3.0%)		
		5. 5	1096 (19.6%)	III	
		6. 6	401 (7.2%)	I	
		7. 77	2 (0.0%)		
		8. 99	1 (0.0%)		
16	depression [numeric]	Mean (sd) : 0.4 (0.8)	0 : 3955 (75.5%)	IIIIIIIIIIIIII	4343
		min < med < max:	1 : 876 (16.7%)	III	(45.3%)
		0 < 0 < 9	2 : 205 (3.9%)		
		IQR (CV) : 0 (2.1)	3 : 194 (3.7%)		
			7 : 2 (0.0%)		
17	sleep [numeric]	Mean (sd) : 7 (3.2)	12 distinct values	:	3300
		min < med < max:		:	(34.5%)
		2 < 7 < 99		:	
		IQR (CV) : 2 (0.5)		:	
				:	
18	diabetes [factor]	1. yes	737 (7.7%)	I	0
		2. no	8841 (92.3%)	IIIIIIIIIIIIIIII	(0.0%)

```
raw_data <- raw_data[-c(7:10)]
dfSummary(raw_data[,-1], valid.col = FALSE)
```

Data Frame Summary

raw_data
Dimensions: 9578 x 14
Duplicates: 319

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Missing
1	gender [factor]	1. 1	4706 (49.1%)	IIIIIIII	0
		2. 2	4872 (50.9%)	IIIIIIIIII	(0.0%)
2	age [numeric]	Mean (sd) : 32.4 (23.9)	80 distinct values	:	0
		min < med < max:		::	(0.0%)
		1 < 28 < 80		:::	
		IQR (CV) : 41 (0.7)		:::::::::::	
3	race [factor]	1. 1	1616 (16.9%)	III	0
		2. 2	893 (9.3%)	I	(0.0%)
		3. 3	3449 (36.0%)	IIIIII	
		4. 4	2148 (22.4%)	IIII	
		5. 6	1033 (10.8%)	II	
		6. 7	439 (4.6%)		
4	bmi [numeric]	Mean (sd) : 25.6 (7.9)	436 distinct	:	706
		min < med < max:	values	:::	(7.4%)
		12.1 < 24.6 < 82.9		:::	
		IQR (CV) : 10.4 (0.3)		:::.	
				:::.	

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Missing
5	hdl [numeric]	Mean (sd) : 53.2 (15.2) min < med < max: 10 < 51 < 173 IQR (CV) : 19 (0.3)	116 distinct values	: : .: : : :.	2128 (22.2%)
6	bp_systolic [numeric]	Mean (sd) : 117.9 (18) min < med < max: 66 < 116 < 228 IQR (CV) : 20 (0.2)	71 distinct values	: : : .: :.	2571 (26.8%)
7	bp_diastolic [numeric]	Mean (sd) : 65.7 (15) min < med < max: 0 < 66 < 122 IQR (CV) : 16 (0.2)	59 distinct values	: : : : :.	2571 (26.8%)
8	waist [numeric]	Mean (sd) : 86.9 (22.5) min < med < max: 40.2 < 87.4 < 177.9 IQR (CV) : 31.6 (0.3)	1030 distinct values	: : : : :.	1091 (11.4%)
9	lifestyle [numeric]	Mean (sd) : 478.5 (642.1) min < med < max: 0 < 480 < 9999 IQR (CV) : 300 (1.3)	36 distinct values	: : : : : :.	2625 (27.4%)
10	education [factor]	1. 1 2. 2 3. 3 4. 4 5. 5 6. 7 7. 9	442 (7.9%) 761 (13.6%) 1261 (22.6%) 1715 (30.7%) 1406 (25.1%) 2 (0.0%) 5 (0.1%)	I II III IIII IIII : :	3986 (41.6%)
11	married [factor]	1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 77 8. 99	2866 (51.3%) 419 (7.5%) 637 (11.4%) 170 (3.0%) 1096 (19.6%) 401 (7.2%) 2 (0.0%) 1 (0.0%)	IIIIIIII I II : III I : :	3986 (41.6%)
12	depression [numeric]	Mean (sd) : 0.4 (0.8) min < med < max: 0 < 0 < 9 IQR (CV) : 0 (2.1)	0 : 3955 (75.5%) 1 : 876 (16.7%) 2 : 205 (3.9%) 3 : 194 (3.7%) 7 : 2 (0.0%) 9 : 3 (0.1%)	IIIIIIIIIIII III : : : : :	4343 (45.3%)
13	sleep [numeric]	Mean (sd) : 7 (3.2) min < med < max: 2 < 7 < 99 IQR (CV) : 2 (0.5)	12 distinct values	: : : : :	3300 (34.5%)

Density plots (numerical covariates)

Bar plots (categorical covariates)

```

diabetes_gender = ggplot(raw_data,
  aes(x = diabetes,
      fill = factor(gender,
                    levels = c("1", "2"),
                    labels = c("male", "female")))) +
  geom_bar(position = position_dodge(preserve = "single")) +
  scale_fill_brewer(palette = "Set2") +
  labs(fill = "gender")

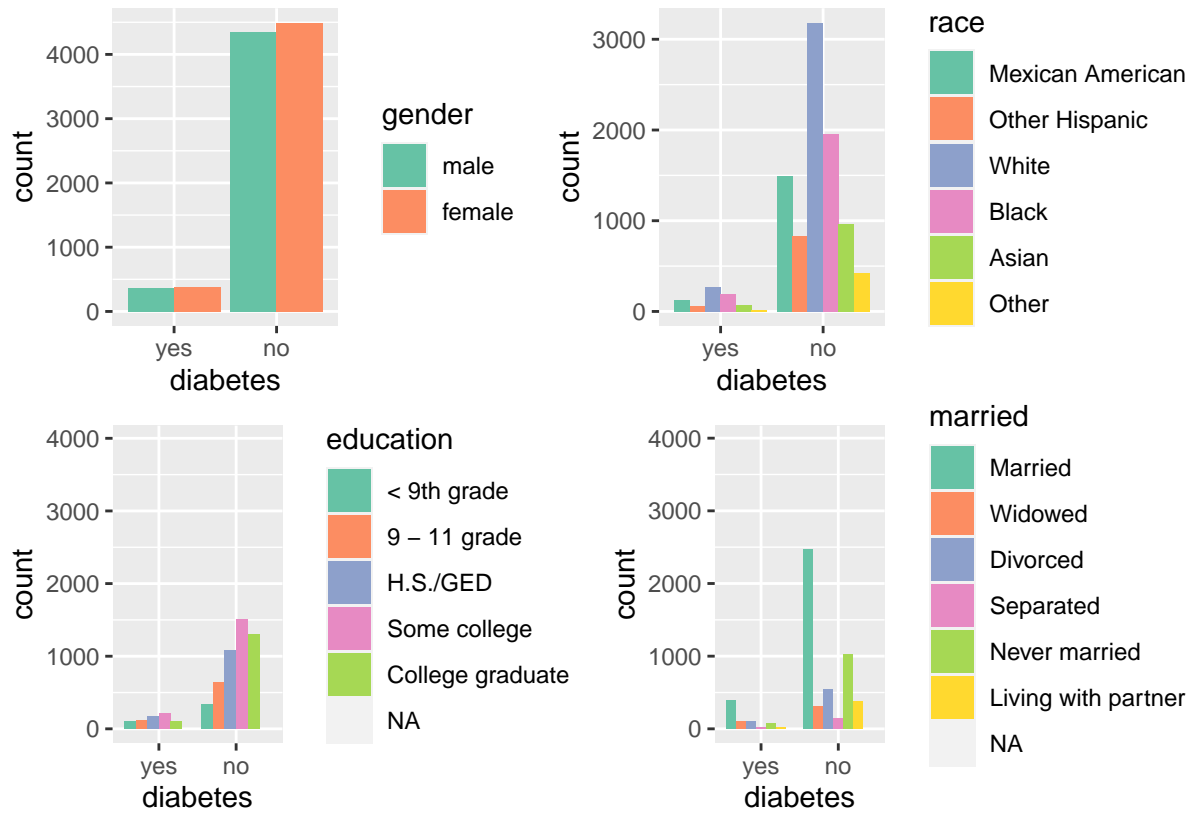
diabetes_race = ggplot(raw_data,
  aes(x = diabetes,
      fill = factor(race,
                    levels = c("1", "2", "3", "4", "6", "7"),
                    labels = c("Mexican American", "Other Hispanic", "White", "Black", "Asian", "O
  geom_bar(position = position_dodge(preserve = "single")) +
  scale_fill_brewer(palette = "Set2") +
  labs(fill = "race")

diabetes_education = ggplot(raw_data,
  aes(x = diabetes,
      fill = factor(education,
                    levels = c("1", "2", "3", "4", "5"),
                    labels = c("< 9th grade", "9 - 11 grade", "H.S./GED", "Some college", "College
  geom_bar(position = position_dodge(preserve = "single")) +
  scale_fill_brewer(palette = "Set2") +
  labs(fill = "education")

diabetes_married = ggplot(raw_data,
  aes(x = diabetes,
      fill = factor(married,
                    levels = c("1", "2", "3", "4", "5", "6"),
                    labels = c("Married", "Widowed", "Divorced", "Separated", "Never married", "Li
  geom_bar(position = position_dodge(preserve = "single")) +
  scale_fill_brewer(palette = "Set2") +
  labs(fill = "married")

(diabetes_gender + diabetes_race) / (diabetes_education + diabetes_married)

```

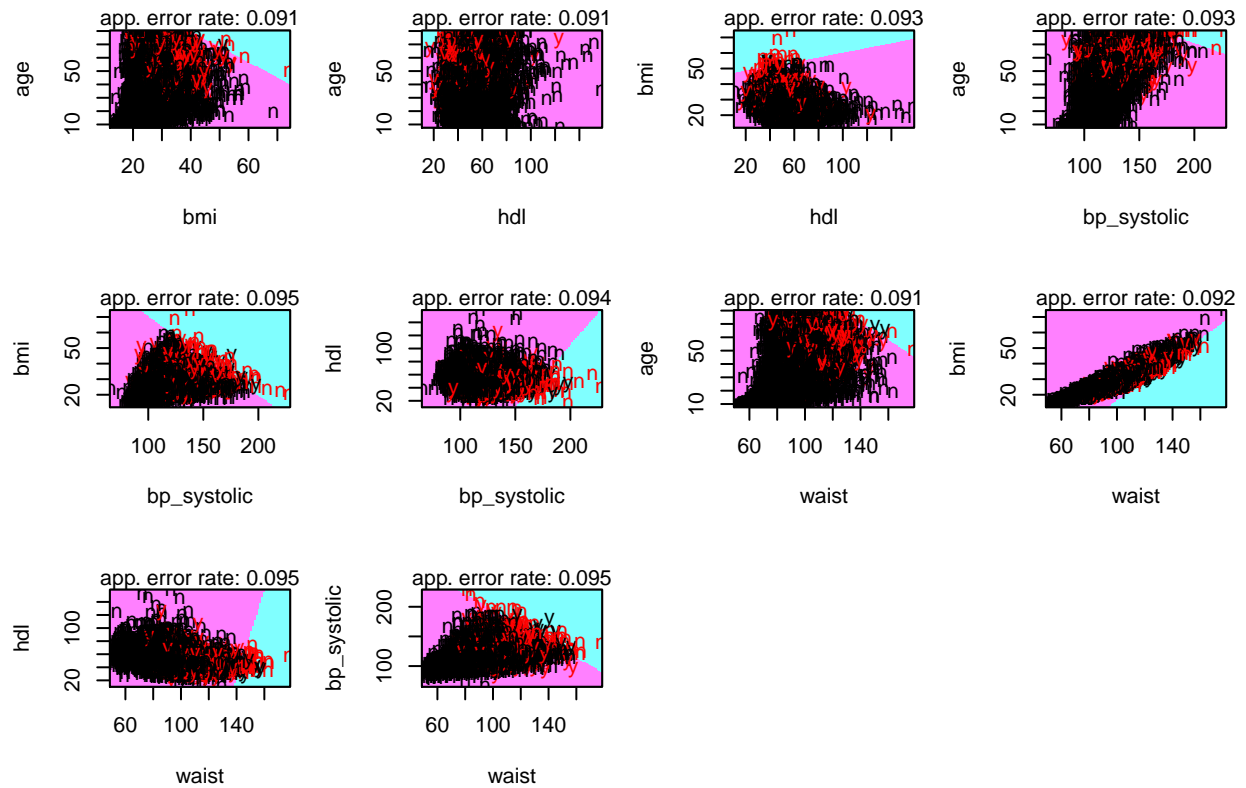
Partition plots

```
set.seed(1)
rowTrain <- createDataPartition(y = raw_data$diabetes,
                                p = 0.7,
                                list = FALSE)

# Exploratory analysis: LDA/QDA/NB based on every combination of two variables

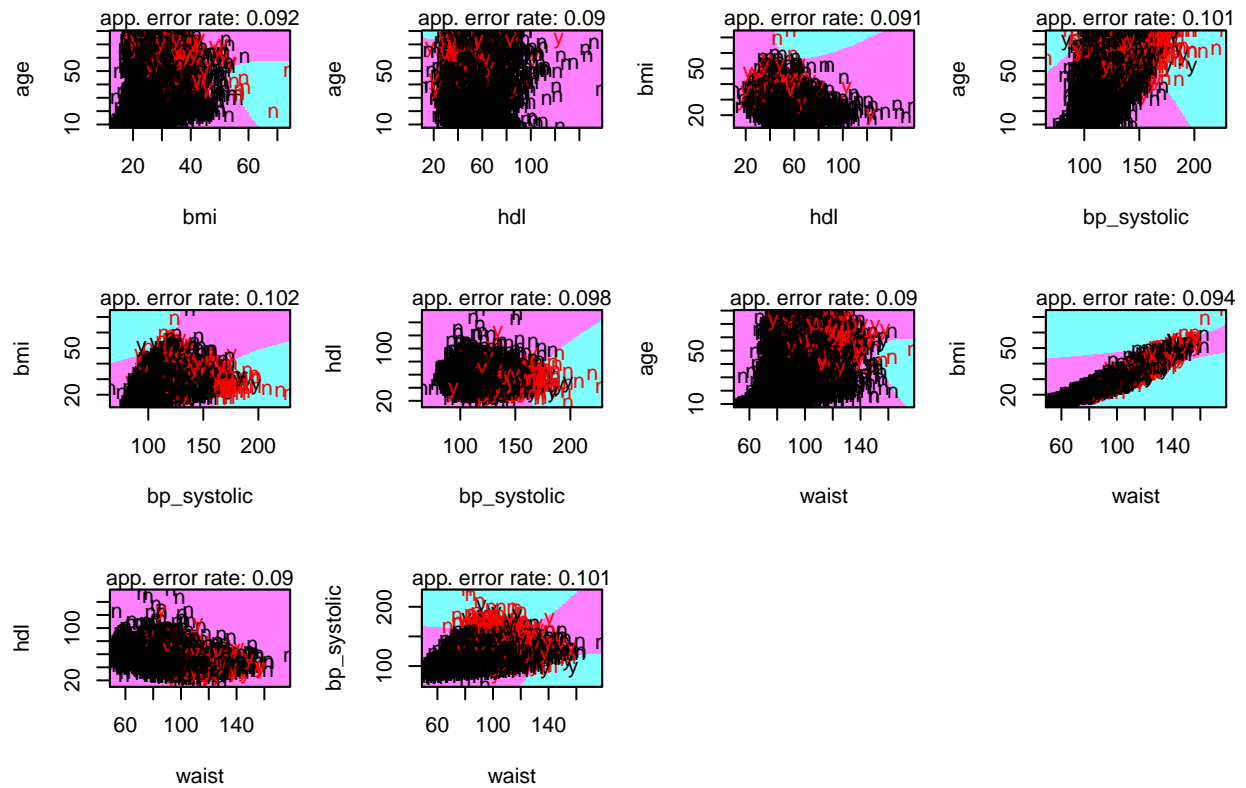
klaR::partimat(diabetes ~ age + bmi + hdl + bp_systolic + waist,
               data = raw_data, subset = rowTrain, method = "lda")
```

Partition Plot



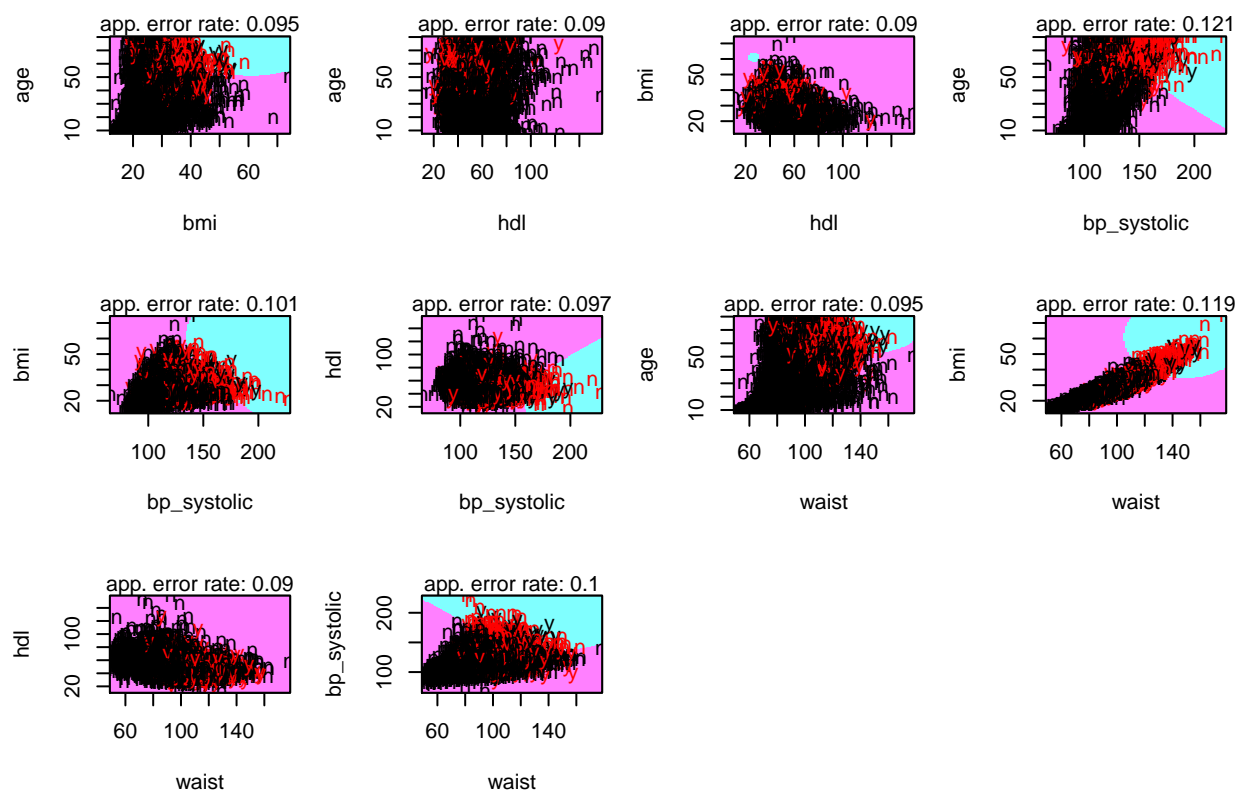
```
klaR::partimat(diabetes ~ age + bmi + hdl + bp_systolic + waist,
  data = raw_data, subset = rowTrain, method = "qda")
```

Partition Plot



```
klaR::partimat(diabetes ~ age + bmi + hdl + bp_systolic + waist,
  data = raw_data, subset = rowTrain, method = "naiveBayes")
```

Partition Plot



Models

training data

```
# Missing data omitted
diabetes_data <- na.omit(raw_data)
summary(diabetes_data)
```

```

ID      gender  race      education  married      age
Min. :73557 1:2077 1: 562 1: 286 1 :2227 Min. :20.00
1st Qu.:76165 2:2169 2: 400 2: 563 5 : 825 1st Qu.:34.00
Median :78695 3:1880 3: 940 3 : 495 Median :48.00
Mean :78660 4: 824 4:1334 6 : 304 Mean :48.37
3rd Qu.:81161 6: 451 5:1122 2 : 268 3rd Qu.:62.00
Max. :83727 7: 129 7: 0 4 : 126 Max. :80.00
9: 1 (Other): 1
bmi hdl bp_systolic bp_diastolic
Min. :14.10 Min. : 10.00 Min. : 66.0 Min. : 0.00
1st Qu.:24.10 1st Qu.: 42.00 1st Qu.:110.0 1st Qu.: 62.00
Median :27.70 Median : 50.00 Median :120.0 Median : 70.00
Mean :28.83 Mean : 52.84 Mean :122.7 Mean : 69.75
```

3rd Qu.:32.20 3rd Qu.: 62.00 3rd Qu.:132.0 3rd Qu.: 78.00
 Max. :74.10 Max. :173.00 Max. :228.0 Max. :122.00

waist	lifestyle	depression	sleep	diabetes
Min. : 55.50	Min. : 0.0	Min. :0.0000	Min. : 2.000	yes: 529
1st Qu.: 87.20	1st Qu.: 240.0	1st Qu.:0.0000	1st Qu.: 6.000	no :3717
Median : 97.30	Median : 420.0	Median :0.0000	Median : 7.000	
Mean : 98.83	Mean : 449.6	Mean :0.3479	Mean : 6.939	
3rd Qu.:108.20	3rd Qu.: 540.0	3rd Qu.:0.0000	3rd Qu.: 8.000	
Max. :177.90	Max. :9999.0	Max. :9.0000	Max. :99.000	

```
set.seed(1)
trainRows <- createDataPartition(diabetes_data$diabetes, p = 0.8, list = FALSE)

# training data
x <- diabetes_data[trainRows, -c(1, 15)]
y <- diabetes_data$diabetes[trainRows]

# test data
x2 <- diabetes_data[-trainRows, -c(1, 15)]
y2 <- diabetes_data$diabetes[-trainRows]
```

Nonlinear models

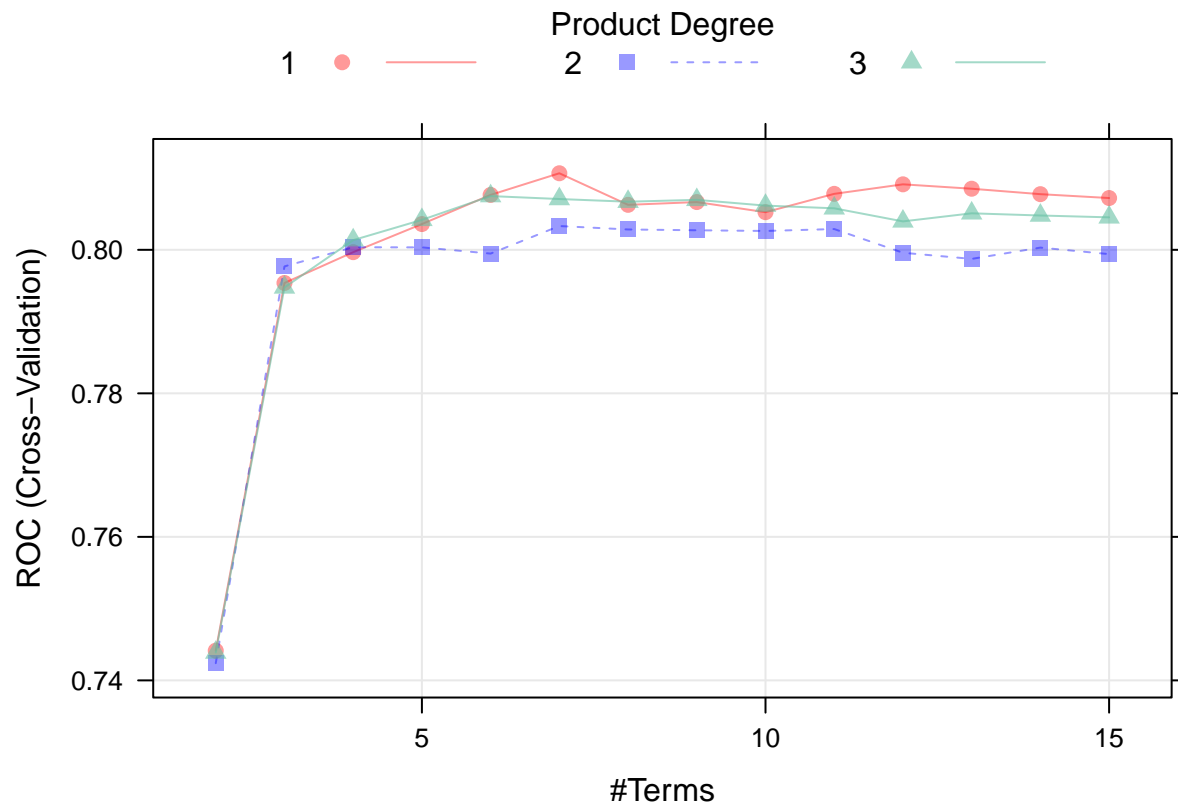
```
ctrl <- trainControl(method = "cv",
                     summaryFunction = twoClassSummary,
                     classProbs = TRUE)

## Non-linear Logistic regression: GAM, MARS
# GAM
#set.seed(1)
#model.gam <- train(x = x,
#                  y = y,
#                  method = "gam",
#                  metric = "ROC",
#                  trControl = ctrl)

#model.gam$finalModel

# MARS
set.seed(1)
model.mars <- train(x = x,
                   y = y,
                   method = "earth",
                   tuneGrid = expand.grid(degree = 1:3,
                                         nprune = 2:15),
                   metric = "ROC",
                   trControl = ctrl)

plot(model.mars)
```



```
coef(model.mars$finalModel)
```

(Intercept)	h(waist-86.9)	h(74-hdl)	h(86-bp_diastolic)
4.92044992	-0.05025456	-0.02809120	-0.01171198
race3	h(bmi-41.4)	h(age-28)	
0.72597496	0.11873361	-0.05851595	

```
## Non-linear Discriminant analysis: QDA, Naive Bayes (NB)
```

```
# QDA = for continuous features
```

```
#set.seed(1)
```

```
#model.qda <- train(x = x,
```

```
#
#           y = y,
#           method = "qda",
#           metric = "ROC",
#           trControl = ctrl)
```

```
# NB
```

```
set.seed(1)
```

```
nbGrid <- expand.grid(usekernel = c(FALSE,TRUE),
                      fL = 1,
                      adjust = seq(.2, 2.5, by = .2))
```

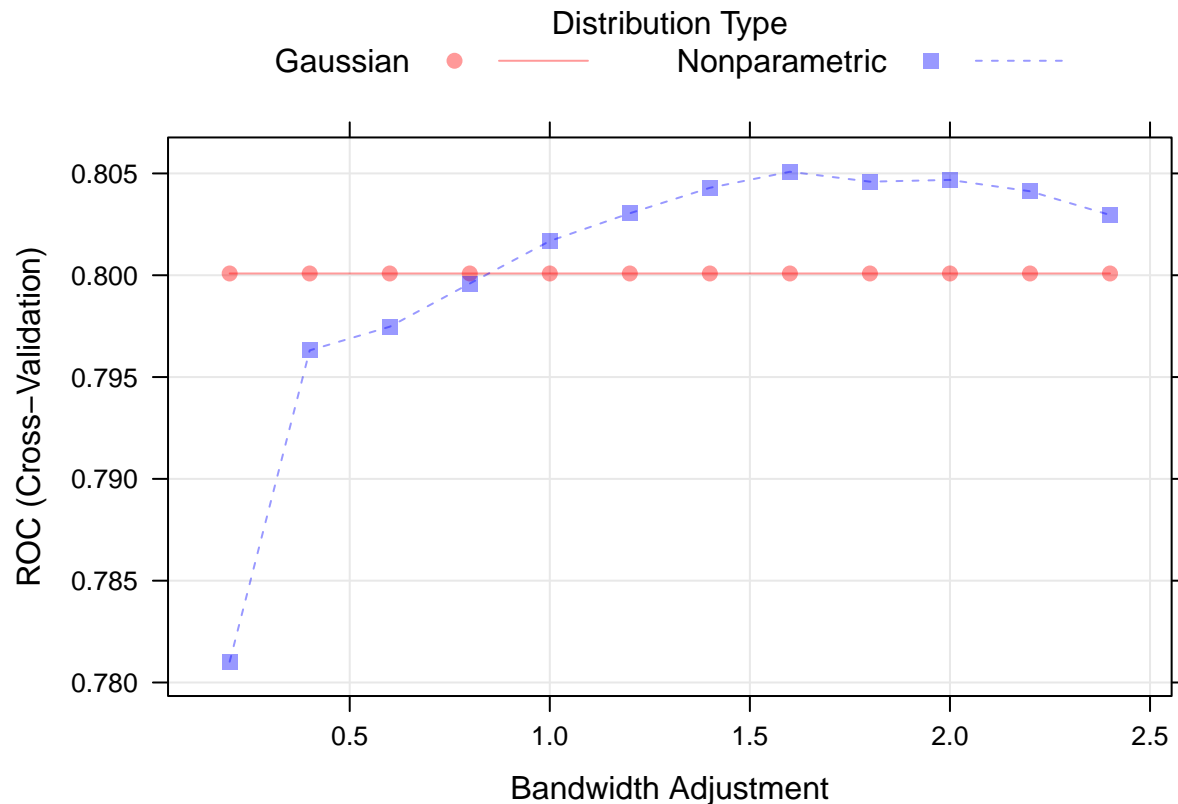
```
model.nb <- train(x = x,
                  y = y,
```

```

method = "nb",
tuneGrid = nbGrid,
metric = "ROC",
trControl = ctrl)

plot(model.nb)

```



```

res <- resamples(list(MARS = model.mars, NB = model.nb))
summary(res)

```

Call: summary.resamples(object = res)

Models: MARS, NB Number of resamples: 10

ROC Min. 1st Qu. Median Mean 3rd Qu. Max. NA's MARS 0.7731281 0.7945418 0.8005599 0.8106704
0.8210825 0.8801973 0 NB 0.7457468 0.7865422 0.8022083 0.8050834 0.8180215 0.8622661 0

Sens Min. 1st Qu. Median Mean 3rd Qu. Max. NA's MARS 0.02380952 0.0952381 0.1162791 0.1178848
0.1578073 0.1666667 0 NB 0.18604651 0.2558140 0.2738095 0.2689922 0.3035714 0.3333333 0

Spec Min. 1st Qu. Median Mean 3rd Qu. Max. NA's MARS 0.966443 0.9798658 0.9831932 0.9842056
0.9915825 0.996633 0 NB 0.909396 0.9403882 0.9495910 0.9479018 0.9621212 0.962963 0