

QMSS Final Project

The Physician Compare website was created by the Centers for Medicare & Medicaid Services (CMS) in December 2010 as required by the Affordable Care Act (ACA) of 2010 to help patients assess and find doctors and hospitals. This dataset contains the information supplied to patients via that website, including patient satisfaction surveys and performance scores across over 100 metrics.

Looking at individual physician scores:

- MIPS
- Performance by measure category
- Organization MIPS

Question/Problem: How can we better help patients assess and find doctors, where the scoring and rating come in a format not easily accessible or understandable by the average individual?

Approach/Methods: Supervised learning for binary classification utilizing the MIPS as a target with other physician scoring methods as predictors (which we know some of the metrics are direct factors of the individual MIPS scoring, such as the IA, ACI, and Quality category scorings). Potential methods outlined below, including generalized linear models and tree methods.

```
suppressPackageStartupMessages(library(dplyr))
suppressPackageStartupMessages(library(stringr))
suppressPackageStartupMessages(library(tidyr))
suppressPackageStartupMessages(library(dmm))
suppressPackageStartupMessages(library(pcaPP))
suppressPackageStartupMessages(library(caret))
suppressPackageStartupMessages(library(splines))
```

```
set.seed(70856775)
```

Potential Methods for Binary Classification:

Using overall MIPS for individuals where MIPS \geq 75, the positive payment adjustment threshold.

- could apply spline to other MIPS, ACI scorings since they're somewhat discrete in nature.
 - ACI \geq 0 : clinician reported ACI category
 - ACI \geq 50: clinician achieved base score for ACI
 - MIPS $<$ 30: Negative Payment Adjustment
- Predictive MIPS \geq 75, essentially.
- Methods to try:
 - glmnet for binary classification (elastic model/penalized logit)
 - glm logit model with polynomials?
 - tree model if we can make it work? (Single Tree, Random Forest, Boosting, Dbarts???)
 - PLSDA or LDA
 - nnet or MARS

```
# Professional.Enrollment.ID is non-predictive, same with PAC_id
dr_scores <- read.csv("full_doctor_scoring.csv",
                     sep = ",", na = c("NA", "N/A"))
# remove majority of missingness by removing each task type
```

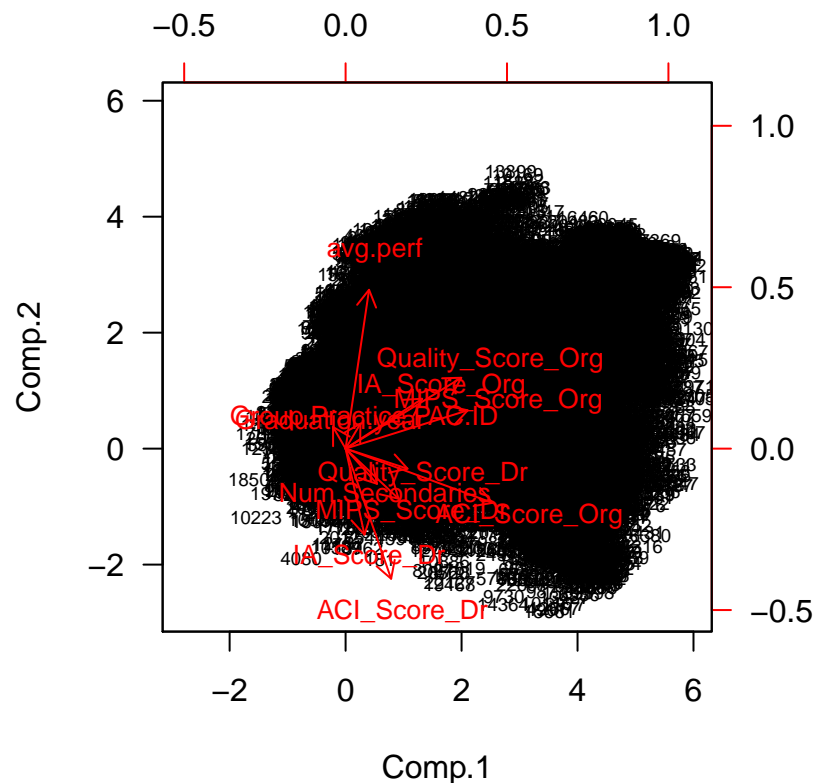
```
no_tasks <- dr_scores[c(1:23, 57)]
# add a "UKN" for unknown category
no_tasks$Credential <- factor(no_tasks$Credential,
                             levels = c(levels(dr_scores$Credential),
                                         "UKN"))
no_tasks$Credential[is.na(no_tasks$Credential)] <- "UKN"
# -1 for not reporting individual ACI Score (NA value)
no_tasks$ACI_Score_Dr[is.na(no_tasks$ACI_Score_Dr)] <- -1
```

PCA to Explore Correlation of Variables

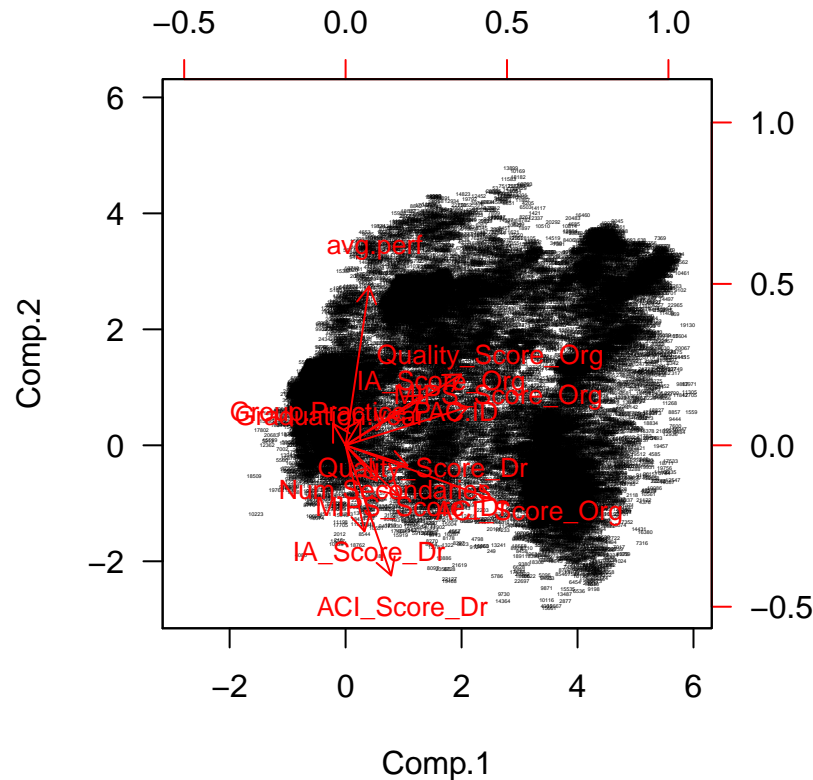
An initial PCA analysis to look understand how rows cluster based on the column variables. This is without considering the binary classifying aspect of MIPS ≥ 75 , and is just to understand some of the structure of the scores in the data.

```
# still minor missingness in Graduation Year and Quality Score
numerics <- no_tasks[, c(8, 10:11, 15:18, 20:24)] %>% na.omit(.)
pr_out <- PCAproj(numerics, scale = sd)
par(mar = c(5, 4, 3, 3) + 0.1, las = 1)

# visualization of any immediate outliers and the variables
biplot(pr_out, scale = 0, cex = c(0.6, 0.8))
```



```
# clusters within the observations visible
biplot(pr_out, scale = 0, cex = c(0.13, 0.8))
```



As mostly expected, the MIPS scores for the hospital organizations that each doctor works at are more correlated to each other than they are to the MIPS scores for each individual doctor. While the number of secondary specialties each doctor has is more correlated to the individual doctor MIPS scores, the group practice ID is more correlated to the organization scores. This is also fairly in line with our expectations that measures for the practice organization would cluster separately from the measures for the individual doctors. Interestingly, the average performance for an individual doctor across task categories seems to be more correlated to the organization scores though. Since performance is measured through individual patient reporting, their experience with the organization itself may be taken into consideration and bias their score, even if the doctor-patient interaction itself was positive.

In the second PCA plot, with the observation labels less cluttered from size, we see that there seem to be distinct clusters of observations. This would indicate groupings within the observations with distinguishing measurement characteristics. Many points do deviate from the groupings themselves. Still, overall, it doesn't appear that any observations seems like a major outlier, as seen in the first plot.

Prepare dataset for training/testing

```
# we do not expect PAC_id or Professional.Enrollment.ID to be predictive
# create factor for classification prediction (Individual Doctor MIPS Score >= 75)
# Score_source_org has no variance -- they're all group

no_tasks <- na.omit(no_tasks)[, c(3:13, 15:ncol(no_tasks))]
```

```

no_tasks$MIPS75_Dr <- factor(no_tasks$MIPS_Score_Dr >= 75.0,
                             levels = c(TRUE, FALSE),
                             labels = c("yes", "no"))

# too many levels for partitioning separately
no_tasks$Prim.Schl <- interaction(no_tasks$Primary.specialty,
                                 no_tasks$Medical.school.name,
                                 sep = ":", drop = TRUE)

# variables we expect to be predictive
preds <- c("Gender", "Prim.Schl", "MIPS_Score_Org",
           "Num.Secondaries", "Quality_Score_Org",
           "IA_Score_Org", "IA_Score_Dr", "Graduation.year",
           "ACI_Score_Dr", "Quality_Score_Dr",
           "avg.perf", "MIPS75_Dr")
no_tasks <- no_tasks[, preds]

remain <- group_by(no_tasks, Prim.Schl) %>%
  summarise(., count = n()) %>%
  filter(., count > 50) %>%
  .$Prim.Schl

no_tasks <- subset(no_tasks, Prim.Schl %in% remain)

```

Maybe explore PCA again for the things we are directly looking at and using the data that we've trimmed down to

```

in_train <- createDataPartition(no_tasks$Prim.Schl,
                                p = 0.8, list = F)

training <- no_tasks[in_train, ]
testing <- no_tasks[-in_train, ]

```

Let's test with a basic linear model:

```

# testing without stratification, med school name, primary speciality
ols <- lm(MIPS75_Dr == "yes" ~ ., data = training)
yhat <- predict(ols, newdata = testing)
z_ols <- factor(yhat > 0.5, levels = c(TRUE, FALSE),
                labels = c("yes", "no"))
confusionMatrix(z_ols, reference = testing$MIPS75_Dr)

```

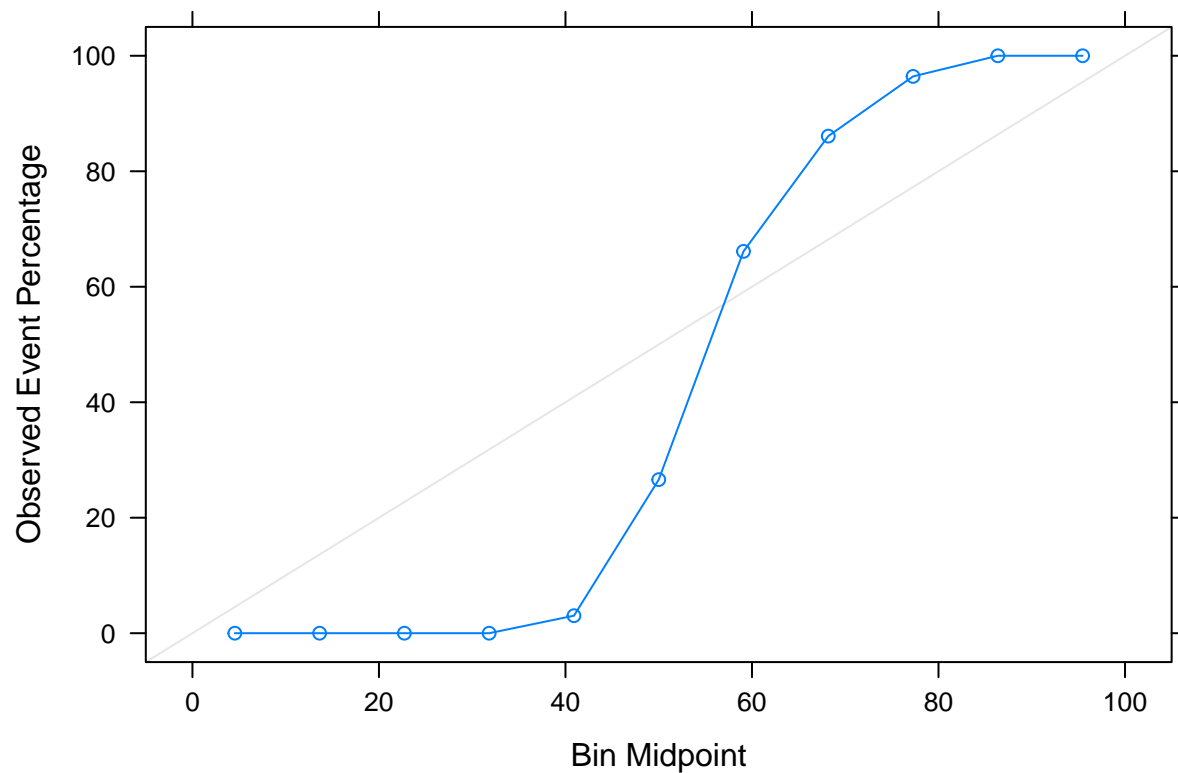
```

## Confusion Matrix and Statistics
##
##           Reference
## Prediction  yes   no
##           yes 1186  99
##           no   13 475

```

```
##
##           Accuracy : 0.9368
##           95% CI : (0.9245, 0.9477)
##    No Information Rate : 0.6763
##    P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.8499
##
##    McNemar's Test P-Value : 9.61e-16
##
##           Sensitivity : 0.9892
##           Specificity : 0.8275
##           Pos Pred Value : 0.9230
##           Neg Pred Value : 0.9734
##           Prevalence : 0.6763
##           Detection Rate : 0.6689
##    Detection Prevalence : 0.7248
##           Balanced Accuracy : 0.9083
##
##           'Positive' Class : yes
##
```

```
calibration(MIPS75_Dr ~ yhat, data = testing) %>%
  plot(.)
```



ADJUST ANALYSIS FOR CHANGES IN FACTORIZATION

The stratified model seems to fit a bit better compared to the model trained and tested without medical school name or primary specialty. This is especially for bins of lower probabilities, such as 0.5, which is fairly expected, as a larger number of observations from stratification and including medical school name and primary specialty would likely provide much more information. These factors, medical school name and primary specialty, probably strongly influence scores. Both the doctor's medical school name and the types of diseases and treatments could bias how scorers evaluate. For example, a scorer may be predisposed to giving a higher score if the doctor being evaluated attended a more prestigious school. Certain fields of medicine, such as surgery or oncology, could incur much higher costs or have higher mortality rates compared to fields such as family medicine, skewing factors of the final MIPS score, such as cost and quality scoring. Still, even with fairly good accuracies, we can see that the model isn't calibrated in an ideal manner.

GLM Models – Logit and GLMnet Penalized

We should include scaling for numericals where it doesn't need to be splined or kept as is?

Num.secondaries, Quality score, IA score is 0-40 so maybe this too, avg.perf

```
# continous measurements can be scaled
# ACI is really a mixture of discrete and continuous

pp_names <- c("Num.Secondaries", "Quality_Score_Org",
              "Quality_Score_Dr", "IA_Score_Org",
              "IA_Score_Dr", "avg.perf")
pp <- list(center = pp_names, scale = pp_names)
```

Resplit to only put half in each to make functions run a bit faster

```
in_train <- createDataPartition(no_tasks$MIPS75_Dr,
                                p = 0.6, list = F)
training <- no_tasks[in_train, ]
testing <- no_tasks[-in_train, ]
```

```
logit <- glm(MIPS75_Dr ~ ., data = training,
             family = binomial(link = "logit"))
z <- predict(logit, newdata = testing,
             type = "response") > 0.5
z <- factor(z, levels = c(TRUE, FALSE),
           labels = c("no", "yes"), order = T)
confusionMatrix(z, testing$MIPS75_Dr)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  yes   no
##           yes 2313   86
##           no   78 1108
##
```

```
##               Accuracy : 0.9543
##               95% CI : (0.9469, 0.9609)
##      No Information Rate : 0.6669
##      P-Value [Acc > NIR] : <2e-16
##
##               Kappa : 0.8969
##
##  Mcnemar's Test P-Value : 0.5846
##
##               Sensitivity : 0.9674
##               Specificity : 0.9280
##      Pos Pred Value : 0.9642
##      Neg Pred Value : 0.9342
##      Prevalence : 0.6669
##      Detection Rate : 0.6452
##      Detection Prevalence : 0.6692
##      Balanced Accuracy : 0.9477
##
##      'Positive' Class : yes
##
```

Wow glm fit way better than I thought so quickly

```
# glmnet PENALIZATION
ctrl <- trainControl(method = "cv", number = 3)
enet <- train(formula(logit), data = training,
              method = "glmnet", trControl = ctrl,
              tuneLength = 10, preProcess = pp)
enet_hat <- predict(enet, newdata = testing)
confusionMatrix(enet_hat, reference = testing$MIPS75_Dr)
```

```
## Confusion Matrix and Statistics
##
##               Reference
## Prediction  yes   no
##      yes 2311   95
##      no   80 1099
##
##               Accuracy : 0.9512
##               95% CI : (0.9436, 0.958)
##      No Information Rate : 0.6669
##      P-Value [Acc > NIR] : <2e-16
##
##               Kappa : 0.8898
##
##  Mcnemar's Test P-Value : 0.2899
##
##               Sensitivity : 0.9665
##               Specificity : 0.9204
##      Pos Pred Value : 0.9605
##      Neg Pred Value : 0.9321
##      Prevalence : 0.6669
```

```
##          Detection Rate : 0.6446
##    Detection Prevalence : 0.6711
##      Balanced Accuracy : 0.9435
##
##      'Positive' Class : yes
##
```

Linear Models with Polynomials and Splines – how linear is our data? Do splines make it better? WHY IS IT LEARNING SO WELL ALREADY

Potential polynomials:

- Num.Secondaries

Potential splines:

- ACI: spline at -1:0, 0:50, 50:
- raw MIPS: spline at 0:30, 30:75, 75:

```
# with polynomials and splines
poly <- glm(MIPS75_Dr ~ . +
            bs(MIPS_Score_Org, knots = c(30, 75)) +
            bs(ACI_Score_Dr, knots = c(0, 50)) +
            poly(Num.Secondaries, degree = 2),
            data = training,
            family = binomial(link = "logit"))
poly_z <- predict(poly, newdata = testing,
                  type = "response") > 0.5
poly_z <- factor(poly_z, levels = c(TRUE, FALSE),
                  labels = c("no", "yes"), order = T)
confusionMatrix(poly_z, testing$MIPS75_Dr)
```

```
## Confusion Matrix and Statistics
##
##          Reference
## Prediction  yes   no
##          yes 2384   10
##          no    7 1184
##
##          Accuracy : 0.9953
##          95% CI : (0.9924, 0.9972)
##    No Information Rate : 0.6669
##    P-Value [Acc > NIR] : <2e-16
##
##          Kappa : 0.9893
##
##    McNemar's Test P-Value : 0.6276
##
##          Sensitivity : 0.9971
##          Specificity : 0.9916
##    Pos Pred Value : 0.9958
##    Neg Pred Value : 0.9941
```



```
##           Prevalence : 0.6669
##           Detection Rate : 0.6650
##           Detection Prevalence : 0.6678
##           Balanced Accuracy : 0.9943
##
##           'Positive' Class : yes
##
```

```
# glmnet PENALIZATION
# with polynomials and splines
poly_el <- train(formula(poly), data = training,
                  method = "glmnet", trControl = ctrl,
                  tuneLength = 10, preProcess = pp)
poly_yh <- predict(poly_el, newdata = testing)
confusionMatrix(poly_yh, reference = testing$MIPS75_Dr)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  yes   no
##           yes 2388   14
##           no    3 1180
##
##           Accuracy : 0.9953
##           95% CI : (0.9924, 0.9972)
##           No Information Rate : 0.6669
##           P-Value [Acc > NIR] : < 2e-16
##
##           Kappa : 0.9893
##
##           McNemar's Test P-Value : 0.01529
##
##           Sensitivity : 0.9987
##           Specificity : 0.9883
##           Pos Pred Value : 0.9942
##           Neg Pred Value : 0.9975
##           Prevalence : 0.6669
##           Detection Rate : 0.6661
##           Detection Prevalence : 0.6700
##           Balanced Accuracy : 0.9935
##
##           'Positive' Class : yes
##
```

DBARTS/trees and then NNET?