

Fragile Families Challenge

Group 5

Introduction



Figure 1:

• BACKGROUND

The Fragile Families Challenge is a mass collaboration that will combine predictive modeling, causal inference, and in-depth interviews to yield insights that can improve the lives of disadvantaged children in the United States. By working together we can discover things that none of us can discover individually.

The Fragile Families Challenge is based on the Fragile Families and Child Wellbeing Study, which has followed thousands of American families for more than 15 years. During this time, the Fragile Families study collected information about the children, their parents, their schools, and their larger environments.

Social Scientists \longleftrightarrow Data Scientists

Figure 2:

• HELP THE WORLD:

Designed to produce scientific knowledge that can be used to improve the lives of disadvantaged children in the United States. Even more than that, we hope the Fragile Families Challenge can serve as a model for how social scientists and data scientists can collaborate on problems of societal importance.

- WIN PRIZES:

We will award prizes to participants who make important contributions to the project. All prize winners will be given an all-expenses paid trip to Princeton University for the scientific workshop at the end of the project.

Research Data

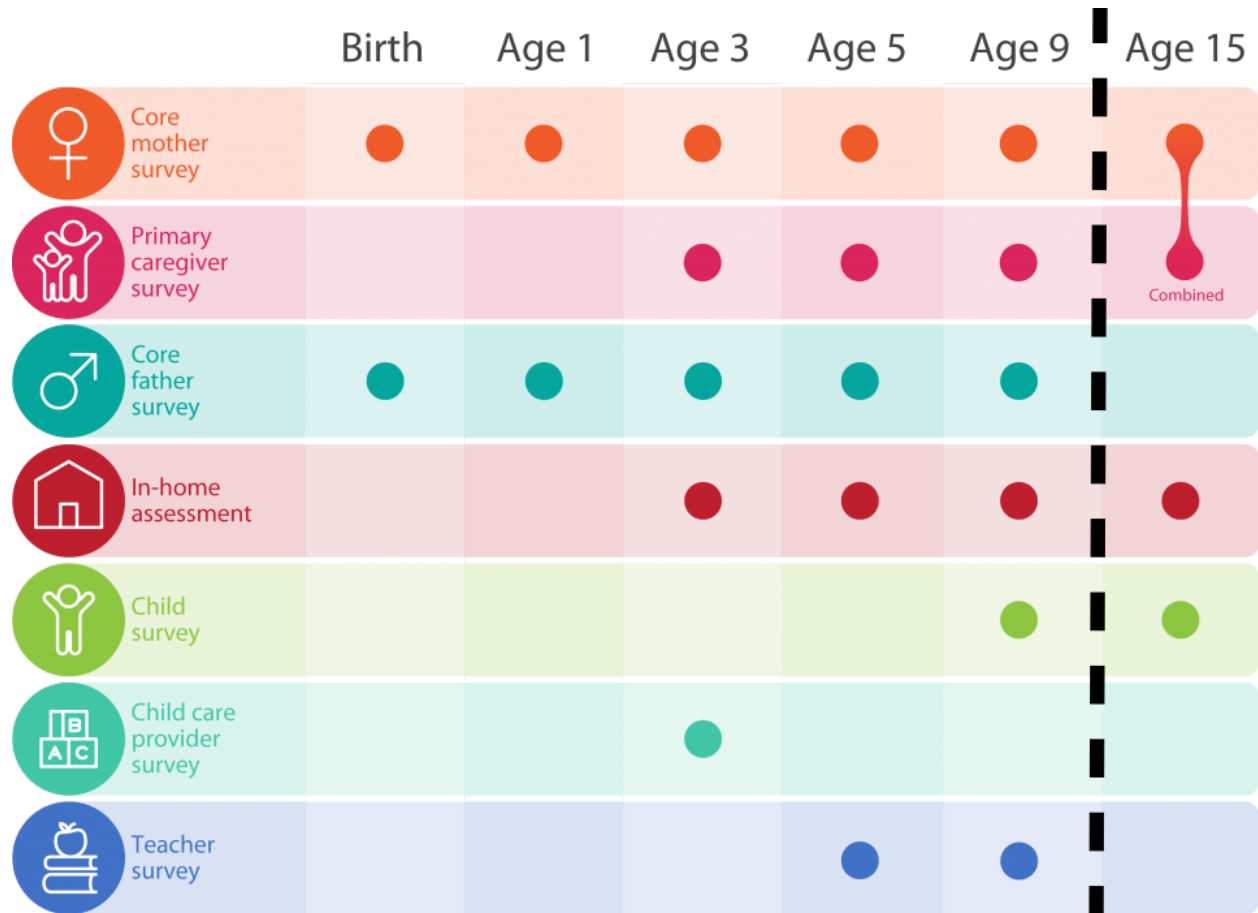


Figure 3:

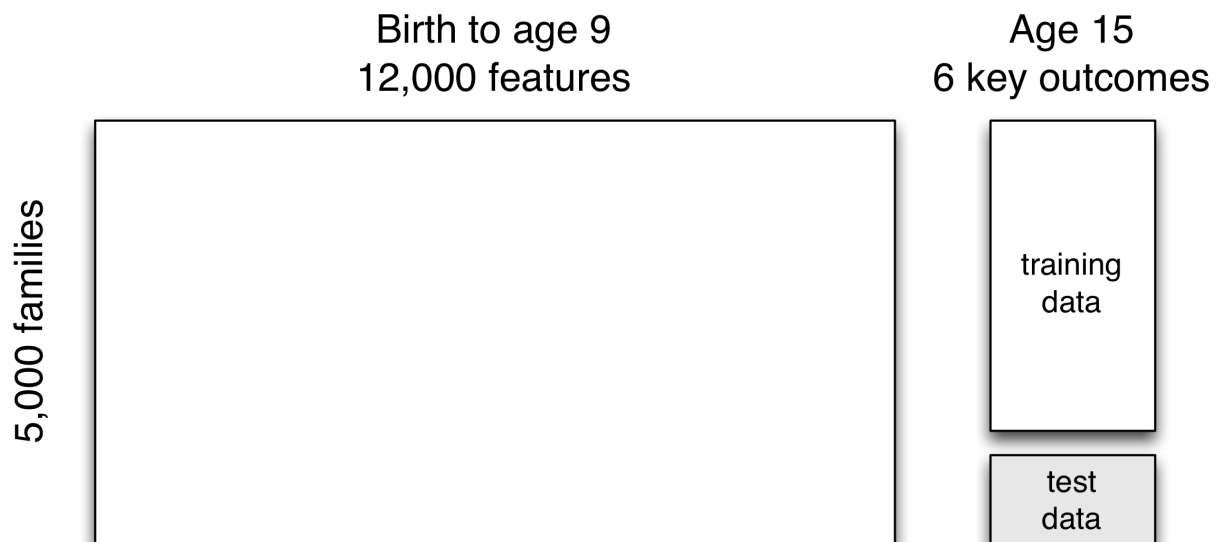


Figure 4:

Continuous outcomes:

- ▶ GPA
- ▶ Grit
- ▶ Material hardship

Binary outcomes:

- ▶ Housing eviction
- ▶ Layoff of a caregiver
- ▶ Job training for a caregiver

Figure 5:

Data Matrix

Outcomes

Step 0: Clean Data

- **-9 Not in wave** – Did not participate in survey/data collection component
- **-8 Out of range** – Response not possible; rarely used
- **-7 Not applicable** (also -10/-14) – Rarely used for survey questions
- **-6 Valid skip** – Intentionally not asked question; question does not apply to respondent or response known based on prior information.
- **-5 Not asked “Invalid skip”** – Respondent not asked question in the version of the survey they received.
- **-3 Missing** – Data is missing due to some other reason; rarely used
- **-2 Don’t know** – Respondent asked question; Responded “Don’t Know”.
- **-1 Refuse** – Respondent asked question; Refused to answer question

Data frame: 4242x12943



Delete the variables that have NAs $\geq 80\%$ or have the same value

4242x9575



Delete the observations that the responses are NAs

1165x9575



numerical 1020

categorical 8546

remove 9 variables cannot be used

```
library(readr)
library(Matrix)
library(mlr)
library(xgboost)
library(plyr)
```

```
bg=read_csv('background.csv')
train=read_csv('train.csv')
```

```
# Delete the records that corresponding to NA in the train
```

```

naid<- which(is.na(train$gpa))
bgtrain<- bg[setdiff(train$challengeID,train$challengeID[naid]),]

# Delete the variables that are 80% NAs
n <- ncol(bgtrain)
na_count<- rep(NA, n)
for(i in 1:n){
    na_count[i]<- sum(is.na(bgtrain[,i]))
}
na_index<- c(1:n)[na_count>=0.8*nrow(bgtrain)]
bgtrain<- bgtrain[,-na_index]

# Delete the variables that have the same value
namelist=names(bgtrain)
for (f in namelist) {
    if(nrow(unique(bgtrain[f]))==1) {
        bgtrain[f]=NULL
    }
}

# Delete the variables that can neither be used as character
n<- ncol(bgtrain)
t<- rep(NA, n)
for(i in 1:n){
    t[i]<- typeof(bgtrain[[1,i]])
}
# unique(t)
char_index<- c(1:n)[t=="character"]
inter_index<- c(1:n)[t=="integer"]
dou_index<- c(1:n)[t=="double"]

# separate the variable by different type, character are those that cannot
# be convert to any type of data, doubles are continuous data. And intergers may
# contain part of continuous data. Thus, we should also handle these variables.
bgtrain_char<- bgtrain[, char_index]
bgtrain_inter<- bgtrain[, inter_index]
bgtrain_dou<- bgtrain[, dou_index]

# For interger-type variables, check the factor number it contains. And we consider those have
# more than 20 factors as continuous number.
fac_num<- as.numeric(apply(bgtrain_inter, 2, function(vec)
{return(length(table(vec)))} ))
fac_index<- which(fac_num<=20)
bgtrain_dou<- cbind(bgtrain_inter[,-fac_index], bgtrain_dou)
bgtrain_inter<- bgtrain_inter[,fac_index]

# Fianlly, we get two seperate dataset to deal with in next steps
bgtrain_dou<- data.frame(bgtrain_dou)
bgtrain_inter<- data.frame(bgtrain_inter)

```

For Numerical Variables

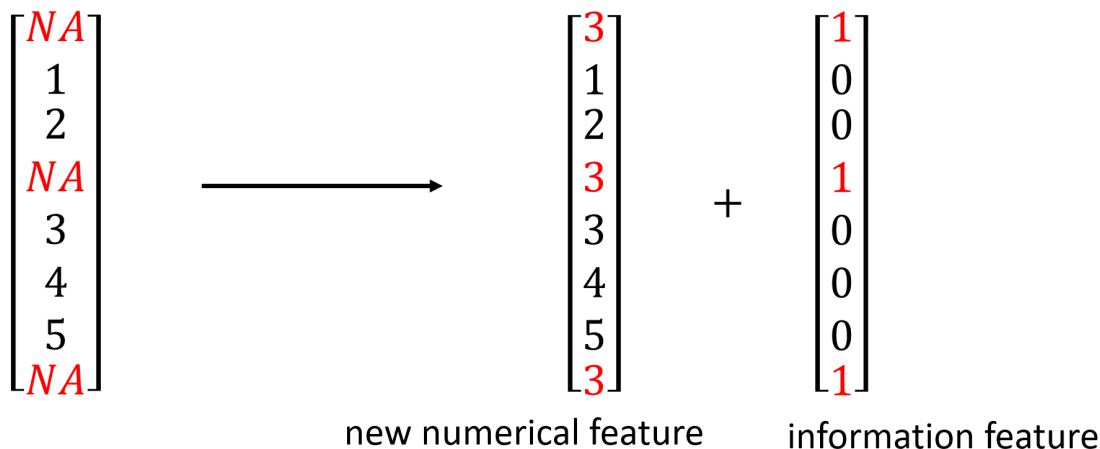


Figure 6: Numerical

Deal with Numerical Variables

```
### fill NA and create new data
fill_each_column <- function(each_col){
  na_label <- is.na(each_col)
  if(sum(na_label) > 0){
    cate_col <- ifelse(na_label == T, 1, 0) ## T = 1, is NA
    fill <- median(each_col, na.rm = T)
    each_col[na_label] <- fill
    return(list(NEW_COLUMN = each_col, NEW_CATE = cate_col))
  }else{
    ## no NA in a column
    return(list(NEW_COLUMN = rep(0,length(each_col)), NEW_CATE = rep(0,length(each_col))))
  }
}
```

Deal with Categorical Variables

Step 1: Regression and Prediction on GPA

Model 1: 9-year-old kid numerical variables + numerical information variables

```
my.text <- "^[a-z]{1,3}5)"
indices <- grepl(my.text,colnames(bgtrain_dou))
mat <- bgtrain_dou[,indices]

## INPUT data= continuous dataframe
gp5data_dou <- mat
```

For Categorical Variables

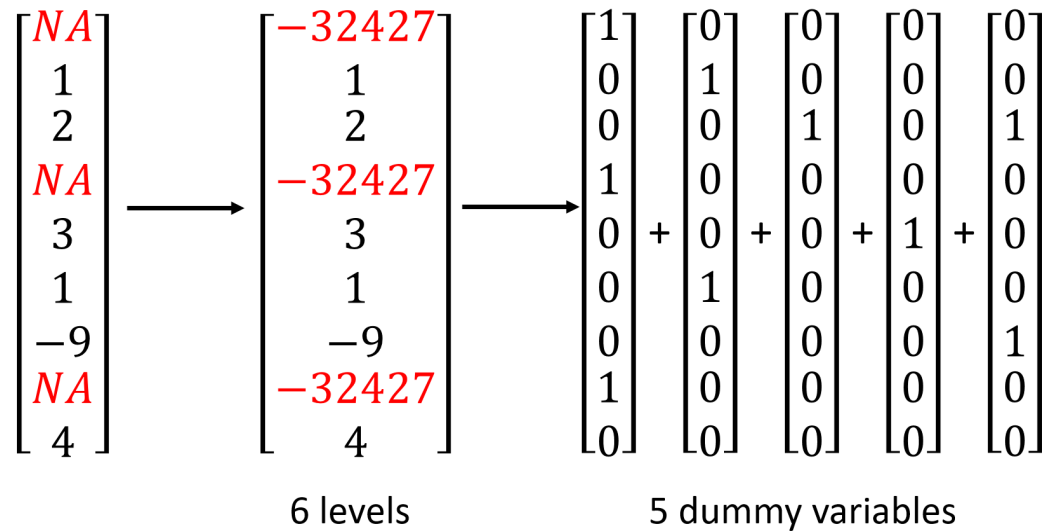


Figure 7: Categorical

```
gp5data_dou <- matrix(unlist(gp5data_dou), nrow(gp5data_dou))
## 1 for new column, 2 for new categorical features
gp5data_dou_RMNA <- apply(gp5data_dou, 2, function(col){fill_each_column(col)[[1]]})
cate_dou_na <- apply(gp5data_dou, 2, function(col){fill_each_column(col)[[2]]})

# remove NON-NA columns
non_na_dou <- colSums(cate_dou_na) == 0
cate_dou_na <- cate_dou_na[,!non_na_dou]
bg_dou <- bg[,colnames(mat)]

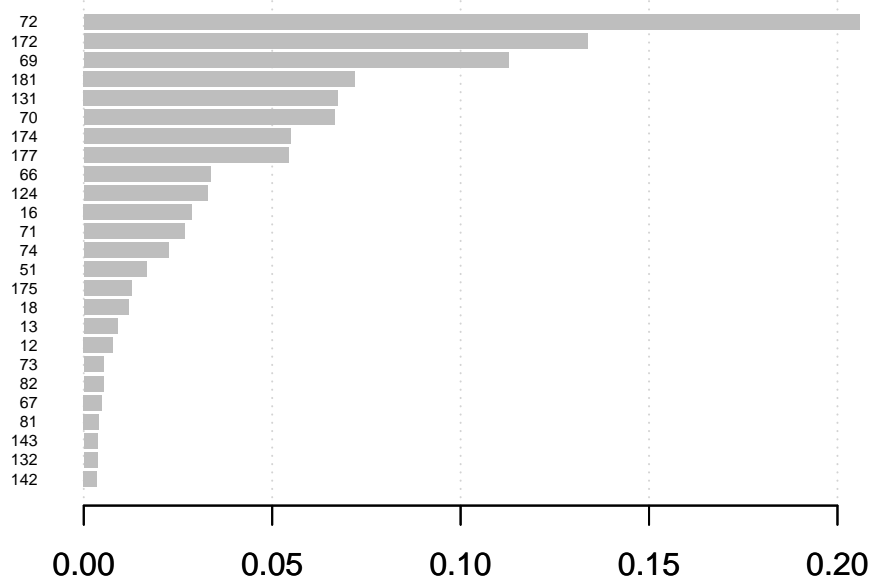
gp5pred_dou <- bg_dou
gp5pred_dou <- matrix(unlist(gp5pred_dou), nrow(gp5pred_dou))
## 1 for new column, 2 for new categorical features
gp5pred_dou_RMNA <- apply(gp5pred_dou, 2, function(col){fill_each_column(col)[[1]]})
cate_pred_dou_na <- apply(gp5pred_dou, 2, function(col){fill_each_column(col)[[2]]})

dtrain=xgb.DMatrix(cbind(data.matrix(mat),cate_dou_na),label=train$gpa[which(!is.na(train$gpa))])
dtest=xgb.DMatrix(cbind(data.matrix(bg_dou),cate_pred_dou_na))
#dtest=xgb.DMatrix(new_features_test)
params=list(
  objective='reg:linear',
  subsample=0.9,
  colsample_bytree=0.8,
  eta=0.05,
  max_depth=1
)

xgb.cv(nfold=10,data=dtrain,params = params,nround=300, print_every_n = 50)
```

```
## [1] train-rmse:2.344313+0.006535 test-rmse:2.343521+0.061291
## [51] train-rmse:0.634044+0.003701 test-rmse:0.646974+0.031971
## [101] train-rmse:0.594031+0.003845 test-rmse:0.618214+0.035590
## [151] train-rmse:0.583825+0.003640 test-rmse:0.619410+0.034911
## [201] train-rmse:0.575686+0.003542 test-rmse:0.620738+0.034979
## [251] train-rmse:0.568553+0.003505 test-rmse:0.622452+0.034688
## [300] train-rmse:0.562188+0.003508 test-rmse:0.623996+0.034575
```

```
model=xgb.train(data=dtrain,params=params,nrounds=100)
imp=xgb.importance(model=model)
xgb.plot.importance(importance_matrix = imp)
```



```
sub=predict(model,dtest)
```

Model 2: 9-year-old kid categorical and KMeans

```
##### K-MEANS func
new_features_kmeans <- function(data, K){
  ## Return cluster.id
  kmeans_results <- kmeans(data, centers = K, iter.max = 500)
  return(kmeans_results$cluster)
}
#####
generate_new_f_kmeans <- function(data, cluster.id){
  new_data_lm <- data.frame(CCLUS = cluster.id,
```



```

                                Records = data)
  new_data_lm_done <- ddply((new_data_lm), .(CLUS), colMeans)
  return(DATA = t(new_data_lm_done[, -1]))
}

bg_dou <- bg[, colnames(bgtrain_dou)]

gp5pred_dou <- bg_dou
gp5pred_dou <- matrix(unlist(gp5pred_dou), nrow(gp5pred_dou))
## 1 for new column, 2 for new categorical features
gp5pred_dou_RMNA <- apply(gp5pred_dou, 2, function(col){fill_each_column(col)[[1]]})
cate_pred_dou_na <- apply(gp5pred_dou, 2, function(col){fill_each_column(col)[[2]]})

my.text <- "[a-z]{1,3}5"
indices <- grepl(my.text, colnames(bgtrain_inter))
mat <- bgtrain_inter[, indices]

#grp5_cate_422 <- bgtrain_inter
grp5_cate_422 <- mat
grp5_cate_422[is.na(grp5_cate_422)] <- -32767
for (iter in 1:ncol(grp5_cate_422 )){
  grp5_cate_422[, iter] <- as.factor(grp5_cate_422[, iter])
}
aaa <- apply(grp5_cate_422, 2, createDummyFeatures)

aaa1 <- matrix(ncol = 0, nrow = nrow(grp5_cate_422))
for (iter in 1:length(aaa)){
  aaa1 <- cbind(aaa1, aaa[[iter]][, 2:ncol(aaa[[iter]])])
}

cs.id <- new_features_kmeans(t(aaa1), K = 250)
new_features_train <- generate_new_f_kmeans(t(aaa1), cs.id)

## bg_dou_cate begin
bg_inter <- bg[, colnames(mat)]
grp5_pre_cate_422 <- bg_inter
grp5_pre_cate_422[is.na(grp5_pre_cate_422)] <- -32767
grp5_pre_cate_422 <- data.frame(grp5_pre_cate_422)

for (iter in 1:ncol(grp5_pre_cate_422 )){
  set1 <- colnames(aaa[[iter]])
  grp5_pre_cate_422[!(grp5_pre_cate_422[, iter] %in% as.numeric(set1)), iter] <- -32767
  grp5_pre_cate_422[, iter] <- as.factor(grp5_pre_cate_422[, iter])
}

aaa_pre <- apply(grp5_pre_cate_422, 2, createDummyFeatures)
## This is a list
aaa_pre1 <- matrix(ncol = 0, nrow = nrow(grp5_pre_cate_422))
for (iter in 1:length(aaa_pre)){
  aaa_temp <- aaa_pre[[iter]][, colnames(aaa[[iter]])]
  aaa_temp <- aaa_temp[, 2:ncol(aaa[[iter]])]
  aaa_pre1 <- cbind(aaa_pre1, aaa_temp)
}

```

```

## k-means features
new_features_test <- generate_new_f_kmeans(t(aaa_pre1), cs.id)

## XGBOOST FEATURES
dtrain=xgb.DMatrix(data.matrix(new_features_train),label=train$gpa[which(!is.na(train$gpa))])
dtest=xgb.DMatrix(data.matrix(new_features_test))
#dim(dtrain)
#dim(dtest)

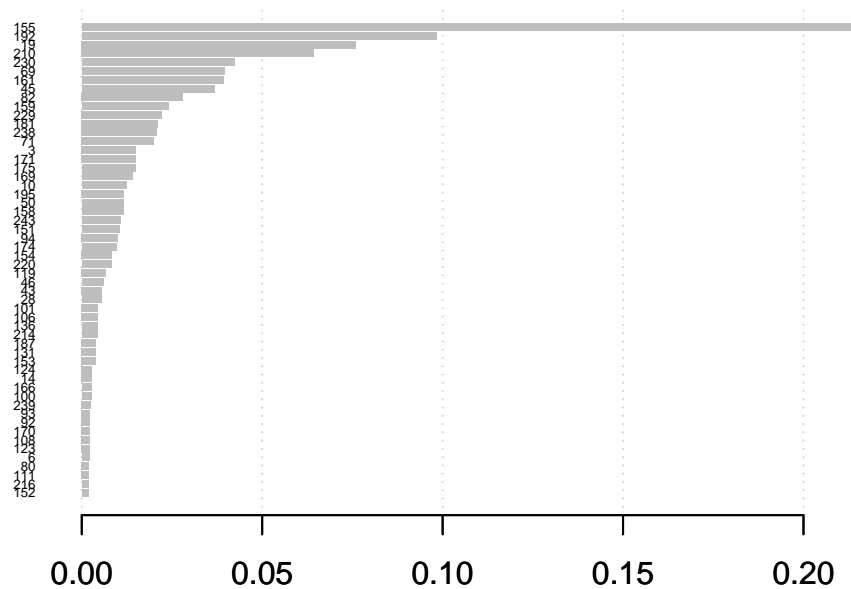
params=list(
  objective='reg:linear',
  subsample=0.9,
  colsample_bytree=1, ## This can be modified
  eta=0.05,
  max_depth=1
)

xgb.cv(nfold=10,data = dtrain, params = params, nround = 300, print_every_n = 50)

## [1] train-rmse:2.344215+0.004731 test-rmse:2.343908+0.045937
## [51] train-rmse:0.633219+0.005467 test-rmse:0.644482+0.039936
## [101] train-rmse:0.591517+0.005685 test-rmse:0.614752+0.051562
## [151] train-rmse:0.580458+0.005557 test-rmse:0.613497+0.052556
## [201] train-rmse:0.571719+0.005549 test-rmse:0.613119+0.053026
## [251] train-rmse:0.564176+0.005600 test-rmse:0.614441+0.053494
## [300] train-rmse:0.557504+0.005666 test-rmse:0.615141+0.053533

model8=xgb.train(data=dtrain,params=params,nrounds=200)
imp=xgb.importance(model=model8)
xgb.plot.importance(importance_matrix = imp)

```



```
cate_result <- predict(model8,dtest)
```

Bagging the two results, that is, given equal weights to the two results and get our final prediction

```
cate_result <- predict(model8,dtest)
nie <- read_csv('prediction1.csv')
yyy <- nie
yyy$gpa <- 0.5 * nie$gpa + 0.5 * cate_result
hist(yyy$gpa)
write_csv(yyy,'prediction.csv')
```

Step 3: Regression and Prediction on Grit & Material Hardship

Grit: Categorical variables related to 9-year-old kids

```
naid<- which(is.na(train$grit))
bgtrain<- bg[setdiff(train$challengeID,train$challengeID[naid]),]
n <- ncol(bgtrain)
na_count<- rep(NA, n)
for(i in 1:n){
  na_count[i]<- sum(is.na(bgtrain[,i]))
}
```

```

}
na_index<- c(1:n)[na_count>=0.8*nrow(bgtrain)]
bgtrain<- bgtrain[,-na_index]

namelist=names(bgtrain)
for (f in namelist) {
  if(nrow(unique(bgtrain[f]))==1) {
    bgtrain[f]=NULL
  }
}
n<- ncol(bgtrain)
t<- rep(NA, n)
for(i in 1:n){
  t[i]<- typeof(bgtrain[[1,i]])
}
char_index<- c(1:n)[t=="character"]
inter_index<- c(1:n)[t=="integer"]
dou_index<- c(1:n)[t=="double"]

bgtrain_char<- bgtrain[, char_index]
bgtrain_inter<- bgtrain[, inter_index]
bgtrain_dou<- bgtrain[, dou_index]

fac_num<- as.numeric(apply(bgtrain_inter, 2, function(vec)
{return(length(table(vec)))} ))
fac_index<- which(fac_num<=15)
bgtrain_dou<- cbind(bgtrain_inter[,-fac_index], bgtrain_dou)
bgtrain_inter<- bgtrain_inter[,fac_index]

bgtrain_dou<- data.frame(bgtrain_dou)
bgtrain_inter<- data.frame(bgtrain_inter)

### fill NA and create new data
fill_each_column <- function(each_col){
  na_label <- is.na(each_col)
  if(sum(na_label) > 0){
    cate_col <- ifelse(na_label == T, 1, 0) ## T = 1, is NA
    fill <- median(each_col, na.rm = T)
    each_col[na_label] <- fill
    return(list(NEW_COLUMN = each_col, NEW_CATE = cate_col))
  }else{
    ## no NA in a column
    return(list(NEW_COLUMN = rep(0,length(each_col)), NEW_CATE = rep(0,length(each_col))))
  }
}

### INPUT data= continuous dataframe
gp5data_dou <- bgtrain_dou
gp5data_dou <- matrix(unlist(gp5data_dou), nrow(gp5data_dou))
## 1 for new column, 2 for new categorical features
gp5data_dou_RMNA <- apply(gp5data_dou, 2, function(col){fill_each_column(col)[[1]]})
cate_dou_na <- apply(gp5data_dou, 2, function(col){fill_each_column(col)[[2]]})

## remove NON-NA columns

```

```

non_na_dou <- colSums(cate_dou_na) == 0
cate_dou_na <- cate_dou_na[,!non_na_dou]

bg_dou <- bg[,colnames(bgtrain_dou)]
gp5pred_dou <- bg_dou
gp5pred_dou <- matrix(unlist(gp5pred_dou), nrow(gp5pred_dou))
## 1 for new column, 2 for new categorical features
gp5pred_dou_RMNA <- apply(gp5pred_dou, 2, function(col){fill_each_column(col)[[1]]})
cate_pred_dou_na <- apply(gp5pred_dou, 2, function(col){fill_each_column(col)[[2]]})

my.text <- "(k5|p5|o5|n5|hv5)"
indices <- grepl(my.text, colnames(bgtrain_inter))
mat <- bgtrain_inter[,indices]

grp5_cate_422 <- mat
grp5_cate_422[is.na(grp5_cate_422)] <- -32767
for (iter in 1:ncol(grp5_cate_422 )){
  grp5_cate_422[,iter] <- as.factor(grp5_cate_422[,iter])
}
aaa <- apply(grp5_cate_422, 2, createDummyFeatures)

aaa1 <- matrix(ncol = 0, nrow = nrow(grp5_cate_422))
for (iter in 1:length(aaa)){
  aaa1 <- cbind(aaa1, aaa[[iter]][,2:ncol(aaa[[iter]])])
}

bg_inter <- bg[, colnames(mat)]
grp5_pre_cate_422 <- bg_inter
grp5_pre_cate_422[is.na(grp5_pre_cate_422)] <- -32767
grp5_pre_cate_422 <- data.frame(grp5_pre_cate_422)

for (iter in 1:ncol(grp5_pre_cate_422 )){
  set1 <- colnames(aaa[[iter]])
  grp5_pre_cate_422[!(grp5_pre_cate_422[,iter] %in% as.numeric(set1)),iter] <- -32767
  grp5_pre_cate_422[,iter] <- as.factor(grp5_pre_cate_422[,iter])
}

aaa_pre <- apply(grp5_pre_cate_422, 2, createDummyFeatures) ## This is a list
aaa_pre1 <- matrix(ncol = 0, nrow = nrow(grp5_pre_cate_422))
for (iter in 1:length(aaa_pre)){
  aaa_temp <- aaa_pre[[iter]][,colnames(aaa[[iter]])]
  aaa_temp <- aaa_temp[,2:ncol(aaa[[iter]])]
  aaa_pre1 <- cbind(aaa_pre1, aaa_temp)
}
## bg_dou_cate end

dtrain=xgb.DMatrix(cbind(data.matrix(aaa1),cate_dou_na),label=train$grit[which(!is.na(train$grit))])
dtest=xgb.DMatrix(cbind(data.matrix(aaa_pre1),cate_pred_dou_na))
params=list(
  objective='reg:linear',
  subsample=0.9,
  colsample_bytree=0.6,
  eta=0.05,

```

```

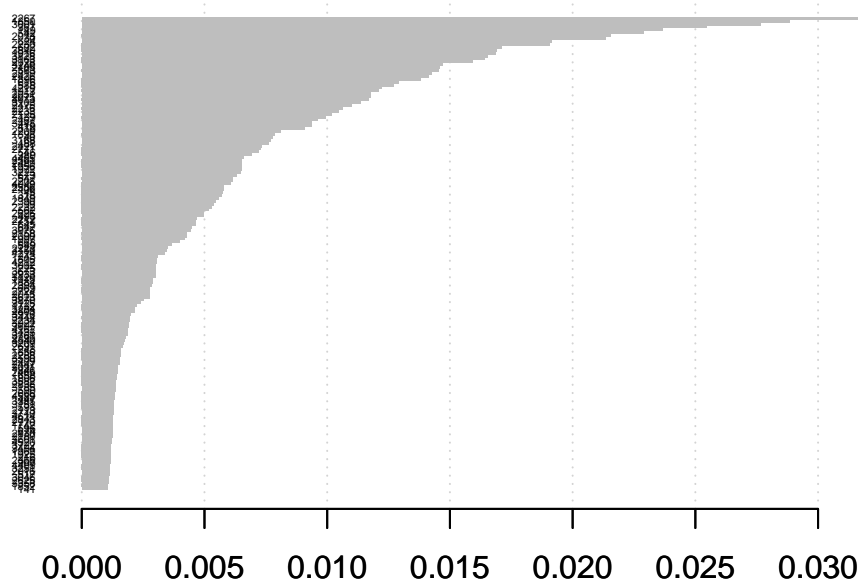
    max_depth=1
)

xgb.cv(nfold=10,data=dtrain,params = params,nround=600,print_every_n = 100)

## [1]  train-rmse:2.823209+0.002900    test-rmse:2.823005+0.026304
## [101]  train-rmse:0.459582+0.002625    test-rmse:0.478059+0.022517
## [201]  train-rmse:0.443973+0.002596    test-rmse:0.473397+0.023916
## [301]  train-rmse:0.432242+0.002610    test-rmse:0.471921+0.023654
## [401]  train-rmse:0.422539+0.002640    test-rmse:0.471670+0.024206
## [501]  train-rmse:0.414120+0.002611    test-rmse:0.471992+0.023796
## [600]  train-rmse:0.406709+0.002576    test-rmse:0.473161+0.022890

model8=xgb.train(data=dtrain,params=params,nrounds=500)
imp=xgb.importance(model=model8)
xgb.plot.importance(importance_matrix = imp)

```



```
sub=predict(model8,dtest)
```

Material Hardship: Mother's categorical variables, and use Kmeans on categorical dummy variables

```

naid<- which(is.na(train$materialHardship))
bgtrain<- bg[setdiff(train$challengeID,train$challengeID[naid]),]
n <- ncol(bgtrain)

```

```

na_count<- rep(NA, n)
for(i in 1:n){
  na_count[i]<- sum(is.na(bgtrain[,i]))
}
na_index<- c(1:n)[na_count>=0.8*nrow(bgtrain)]
bgtrain<- bgtrain[,-na_index]

# Delete the variables that have the same value
namelist=names(bgtrain)
for (f in namelist) {
  if(nrow(unique(bgtrain[f]))==1) {
    bgtrain[f]=NULL
  }
}

n<- ncol(bgtrain)
t<- rep(NA, n)
for(i in 1:n){
  t[i]<- typeof(bgtrain[[1,i]])
}

char_index<- c(1:n)[t=="character"]
inter_index<- c(1:n)[t=="integer"]
dou_index<- c(1:n)[t=="double"]
bgtrain_char<- bgtrain[, char_index]
bgtrain_inter<- bgtrain[, inter_index]
bgtrain_dou<- bgtrain[, dou_index]
fac_index<- which(fac_num<=15)
bgtrain_dou<- cbind(bgtrain_inter[,-fac_index], bgtrain_dou)
bgtrain_inter<- bgtrain_inter[,fac_index]
bgtrain_dou<- data.frame(bgtrain_dou)
bgtrain_inter<- data.frame(bgtrain_inter)
### fill NA and create new data
fill_each_column <- function(each_col){
  na_label <- is.na(each_col)
  if(sum(na_label) > 0){
    cate_col <- ifelse(na_label == T, 1, 0) ## T = 1, is NA
    fill <- median(each_col, na.rm = T)
    each_col[na_label] <- fill
    return(list(NEW_COLUMN = each_col, NEW_CATE = cate_col))
  }else{
    ## no NA in a column
    return(list(NEW_COLUMN = rep(0,length(each_col)), NEW_CATE = rep(0,length(each_col))))
  }
}

}
my.text <- "^{(m)}"
indices <- grepl(my.text, colnames(bgtrain_dou))
mat <- bgtrain_dou[,indices]

## INPUT data= continuous dataframe
gp5data_dou <- mat

```

```

gp5data_dou <- matrix(unlist(gp5data_dou), nrow(gp5data_dou))
## 1 for new column, 2 for new categorical features
gp5data_dou_RMNA <- apply(gp5data_dou, 2, function(col){fill_each_column(col)[[1]]})
cate_dou_na <- apply(gp5data_dou, 2, function(col){fill_each_column(col)[[2]]})
## remove NON-NA columns
non_na_dou <- colSums(cate_dou_na) == 0
cate_dou_na <- cate_dou_na[,!non_na_dou]

bg_dou <- bg[,colnames(mat)]
gp5pred_dou <- bg_dou
gp5pred_dou <- matrix(unlist(gp5pred_dou), nrow(gp5pred_dou))
## 1 for new column, 2 for new categorical features
gp5pred_dou_RMNA <- apply(gp5pred_dou, 2, function(col){fill_each_column(col)[[1]]})
cate_pred_dou_na <- apply(gp5pred_dou, 2, function(col){fill_each_column(col)[[2]]})
dtrain=xgb.DMatrix(data.matrix(mat),label=train$materialHardship[which(!is.na(train$materialHardship))])
dtest=xgb.DMatrix(data.matrix(bg_dou))
params=list(
  objective='reg:linear',
  subsample=0.9,
  colsample_bytree=0.8,
  eta=0.05,
  max_depth=1
)

model=xgb.train(data=dtrain,params=params,nrounds=200)
imp=xgb.importance(model=model)
my.text <- "^{(m)}"
indices <- grepl(my.text, colnames(bgtrain_inter))
mat <- bgtrain_inter[,indices]

grp5_cate_422 <- mat
grp5_cate_422[is.na(grp5_cate_422)] <- -32767
for (iter in 1:ncol(grp5_cate_422 )){
  grp5_cate_422[,iter] <- as.factor(grp5_cate_422[,iter])
}
aaa <- apply(grp5_cate_422, 2, createDummyFeatures)

aaa1 <- matrix(ncol = 0, nrow = nrow(grp5_cate_422))
for (iter in 1:length(aaa)){
  aaa1 <- cbind(aaa1, aaa[[iter]][,2:ncol(aaa[[iter]])])
}
bg_inter <- bg[, colnames(mat)]
grp5_pre_cate_422 <- bg_inter
grp5_pre_cate_422[is.na(grp5_pre_cate_422)] <- -32767
grp5_pre_cate_422 <- data.frame(grp5_pre_cate_422)
for (iter in 1:ncol(grp5_pre_cate_422 )){
  set1 <- colnames(aaa[[iter]])
  grp5_pre_cate_422[!(grp5_pre_cate_422[,iter] %in% as.numeric(set1)),iter] <- -32767
  grp5_pre_cate_422[,iter] <- as.factor(grp5_pre_cate_422[,iter])
}
####

aaa_pre <- apply(grp5_pre_cate_422, 2, createDummyFeatures) ## This is a list

```



```

aaa_pre1 <- matrix(ncol = 0, nrow = nrow(grp5_pre_cate_422))
for (iter in 1:length(aaa_pre)){
  aaa_temp <- aaa_pre[[iter]][,colnames(aaa[[iter])]]
  aaa_temp <- aaa_temp[,2:ncol(aaa[[iter])]]
  aaa_pre1 <- cbind(aaa_pre1, aaa_temp)
}
## bg_dou_cate end

new_features_kmeans <- function(data, K){
  ## Return cluster.id
  kmeans_results <- kmeans(t(data), centers = K, iter.max = 100)
  return(kmeans_results$cluster)
}

generate_new_f_kmeans <- function(data, cluster.id){
  new_data_lm <- data.frame(CLUS = cluster.id,
                           Records = t(data))
  new_data_lm_done <- ddply((new_data_lm), .(CLUS), colMeans)
  return(DATA = t(new_data_lm_done[, -1]))
}

cs.id <- new_features_kmeans(data = aaa1, K = 150)
new_features_train <- generate_new_f_kmeans(aaa1, cs.id)
new_features_test <- generate_new_f_kmeans(aaa_pre1, cs.id)

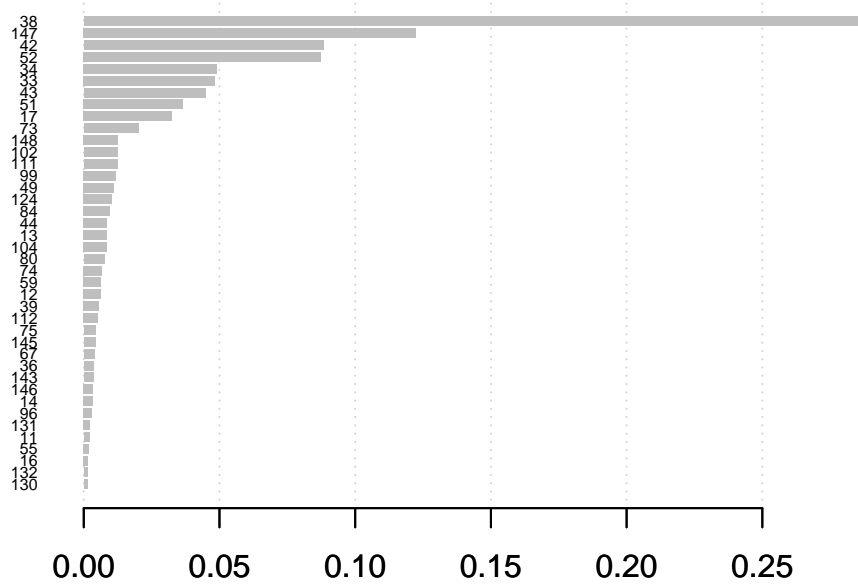
dtrain=xgb.DMatrix(cbind(data.matrix(new_features_train),cate_dou_na),label=train$materialHardship[which(
dtest=xgb.DMatrix(cbind(data.matrix(new_features_test),cate_pred_dou_na))
params=list(
  objective='reg:linear',
  subsample=0.9,
  colsample_bytree=0.8,
  eta=0.05,
  max_depth=1
)

xgb.cv(nfold=10,data=dtrain,params = params,nround=300, print_every_n = 50)

## [1] train-rmse:0.407060+0.000754 test-rmse:0.407118+0.007179
## [51] train-rmse:0.144453+0.001188 test-rmse:0.147241+0.009714
## [101] train-rmse:0.137975+0.001284 test-rmse:0.143170+0.011590
## [151] train-rmse:0.135852+0.001310 test-rmse:0.143221+0.012101
## [201] train-rmse:0.134185+0.001326 test-rmse:0.143203+0.012431
## [251] train-rmse:0.132777+0.001358 test-rmse:0.143395+0.012545
## [300] train-rmse:0.131542+0.001375 test-rmse:0.143801+0.012811

model8=xgb.train(data=dtrain,params=params,nrounds=200)
imp=xgb.importance(model=model8)
xgb.plot.importance(importance_matrix = imp)

```



Final Results

Results							
#	User	GPA ▲	Grit ▲	Material hardship ▲	Eviction ▲	Layoff ▲	Job training ▲
1	kai_niubi	0.36440 (1)	0.21997 (33)	0.02880 (44)	0.05341 (27)	0.17435 (24)	0.20224 (22)
1	ADSgrp5	0.36440 (1)	0.21292 (8)	0.02453 (2)	0.05341 (27)	0.17435 (24)	0.20224 (22)
2	zn	0.36456 (2)	0.21997 (33)	0.02880 (44)	0.05341 (27)	0.17435 (24)	0.20224 (22)
3	ovarol	0.36571 (3)	0.21812 (21)	0.02481 (6)	0.05660 (35)	0.17422 (22)	0.20225 (23)
4	OldDriver.ffc	0.36630 (4)	0.21251 (7)	0.02431 (1)	0.05273 (19)	0.17187 (8)	0.20060 (11)
5	kai_666	0.36733 (5)	0.21997 (33)	0.02880 (44)	0.05341 (27)	0.17435 (24)	0.20224 (22)
6	wjlei1990	0.36742 (6)	0.21176 (4)	0.02431 (1)	0.05233 (16)	0.17149 (7)	0.20251 (26)
7	kgenova	0.36837 (7)	0.21502 (16)	0.02461 (3)	0.05112 (7)	0.17336 (14)	0.19382 (1)
8	malte	0.36912 (8)	0.21569 (17)	0.02484 (7)	0.05208 (11)	0.17461 (25)	0.20079 (13)
9	yjpeng	0.36974 (9)	0.21450 (13)	0.02491 (9)	0.05218 (12)	0.17120 (5)	0.20499 (29)

Figure 8: Leader Board