# R Notebook

## Parametros:

**Measure =** Accuracy
**Columns =** sampling, weight_space, underbagging, learner
**Performance =** tuning_measure
**Filter keys =** imba.rate
**Filter values =** 0.05

```
library("scmamp")
library(dplyr)
```

## Tratamento dos dados

Carregando data set compilado

```
ds = read.csv("/home/rodrigo/Dropbox/UNICAMP/IC/estudo_cost_learning/SummaryResults/summary_compilation
```

```
ds = filter(ds, learner != "classif.rusboost")
summary(ds)
```

```
##                    learner        weight_space
##  classif.ksvm         :17100    Mode :logical
##  classif.randomForest:17100    FALSE:41040
##  classif.rusboost    :    0    TRUE :10260
##  classif.xgboost      :17100    NA's :0
##
##
##
##                                    measure         sampling      underbagging
##  Accuracy                     :10260    ADASYN:10260    Mode :logical
##  Area under the curve         :10260    FALSE :30780    FALSE:41040
##  F1 measure                   :10260    SMOTE :10260    TRUE :10260
##  G-mean                       :10260                    NA's :0
##  Matthews correlation coefficient:10260
##
##
##  tuning_measure    holdout_measure    holdout_measure_residual
##  Min.   :-0.1277   Min.   :-0.2120    Min.   :-0.4658
##  1st Qu.: 0.6911   1st Qu.: 0.4001    1st Qu.: 0.1994
##  Median : 0.9700   Median : 0.8571    Median : 0.5581
##  Mean   : 0.7903   Mean   : 0.6718    Mean   : 0.5298
##  3rd Qu.: 0.9975   3rd Qu.: 0.9900    3rd Qu.: 0.8755
##  Max.   : 1.0000   Max.   : 1.0000    Max.   : 1.0000
##  NA's   :1077      NA's   :1077       NA's   :1077
##  iteration_count                dataset        imba.rate
##  Min.   :1         abalone          : 900    Min.   :0.0010
##  1st Qu.:1         adult            : 900    1st Qu.:0.0100
##  Median :2         bank             : 900    Median :0.0300
##  Mean   :2         car              : 900    Mean   :0.0286
```

```
## 3rd Qu.:3         cardiotocography-10clases:  900   3rd Qu.:0.0500
## Max.   :3         cardiotocography-3clases :  900   Max.   :0.0500
## NA's   :1077   (Other)                    :45900
```

Filtrando pela metrica

```
ds = filter(ds, measure == params$measure)
```

Filtrando o data set

```
if(params$filter_keys != 'NULL' && !is.null(params$filter_keys)){
  dots = paste0(params$filter_keys," == '",params$filter_values,"'")
  ds = filter_(ds, .dots = dots)
}

summary(ds)
```

```
##                   learner      weight_space
##   classif.ksvm       :1230   Mode :logical
##   classif.randomForest:1230   FALSE:2952
##   classif.rusboost    :   0   TRUE :738
##   classif.xgboost     :1230   NA's :0
##
##
##
##                                  measure       sampling    underbagging
##   Accuracy                     :3690   ADASYN: 738   Mode :logical
##   Area under the curve         :   0   FALSE :2214   FALSE:2952
##   F1 measure                   :   0   SMOTE : 738   TRUE :738
##   G-mean                       :   0                 NA's :0
##   Matthews correlation coefficient:  0
##
##
##   tuning_measure   holdout_measure   holdout_measure_residual
##   Min.   :0.2470   Min.   :0.04739   Min.   :0.0367
##   1st Qu.:0.9494   1st Qu.:0.94505   1st Qu.:0.3902
##   Median :0.9688   Median :0.96078   Median :0.7223
##   Mean   :0.9425   Mean   :0.93413   Mean   :0.6602
##   3rd Qu.:0.9908   3rd Qu.:0.98413   3rd Qu.:0.9315
##   Max.   :1.0000   Max.   :1.00000   Max.   :1.0000
##   NA's   :42       NA's   :42        NA's   :42
##   iteration_count         dataset        imba.rate
##   Min.   :1       abalone      : 45   Min.   :0.05
##   1st Qu.:1       adult        : 45   1st Qu.:0.05
##   Median :2       annealing    : 45   Median :0.05
##   Mean   :2       arrhythmia   : 45   Mean   :0.05
##   3rd Qu.:3       balance-scale: 45   3rd Qu.:0.05
##   Max.   :3       bank         : 45   Max.   :0.05
##   NA's   :42      (Other)      :3420
```

Computando as médias das iteracoes

```
ds = group_by(ds, learner , weight_space , measure , sampling , underbagging , dataset , imba.rate)
ds = summarise(ds, tuning_measure = mean(tuning_measure), holdout_measure = mean(holdout_measure),
            holdout_measure_residual = mean(holdout_measure_residual))

ds = as.data.frame(ds)
```

Criando dataframe

```r
# Dividindo o ds em n, um para cada técnica
splited_df = ds %>% group_by_at(.vars = params$columns) %>% do(vals = as.data.frame(.)) %>% select(vals)

# Juntando cada uma das partes horizontalmente em um data set
df_tec_wide = do.call("cbind", splited_df)

# Renomeando duplicacao de nomes
colnames(df_tec_wide) = make.unique(colnames(df_tec_wide))

# Selecionando apenas as medidas da performance escolhida
df_tec_wide_residual = select(df_tec_wide, matches(paste("^", params$performance, "$|", params$performan

# Renomeando colunas
new_names = NULL
for(i in (1:length(splited_df))){
  id = toString(sapply(splited_df[[i]][1, params$columns], as.character))
  new_names = c(new_names, id)
}
colnames(df_tec_wide_residual) = new_names

# Verificando a dimensao do df
dim(df_tec_wide_residual)
```

```
## [1] 82 15
```

```r
# Renomeando a variavel
df = df_tec_wide_residual

head(df)
```

```
##   ADASYN, FALSE, FALSE, classif.ksvm
## 1                           0.9453798
## 2                           0.9607452
## 3                           0.9529558
## 4                           0.8433805
## 5                           1.0000000
## 6                           0.9774219
##   ADASYN, FALSE, FALSE, classif.randomForest
## 1                                   0.9386682
## 2                                          NA
## 3                                   0.9932040
## 4                                   0.9932530
## 5                                   1.0000000
## 6                                   0.9765231
##   ADASYN, FALSE, FALSE, classif.xgboost FALSE, FALSE, FALSE, classif.ksvm
## 1                             0.9573086                        0.9266667
## 2                             0.9734333                        0.9447727
## 3                             0.9881291                        0.9494180
## 4                             0.9823549                        0.9474120
## 5                             1.0000000                        1.0000000
## 6                             0.9733316                        0.9500000
##   FALSE, FALSE, FALSE, classif.randomForest
## 1                                 0.9500000
## 2                                 0.9613617
```

```
## 3                                 0.9747104
## 4                                 0.9792501
## 5                                 1.0000000
## 6                                 0.9497024
##   FALSE, FALSE, FALSE, classif.xgboost FALSE, FALSE, TRUE, classif.ksvm
## 1                            0.9500000                        0.6054167
## 2                            0.9622584                        0.8949096
## 3                            0.9718904                        0.9331186
## 4                            0.9728548                        0.8485392
## 5                            1.0000000                        0.9902311
## 6                            0.9488095                        0.5355159
##   FALSE, FALSE, TRUE, classif.randomForest
## 1                                0.6205556
## 2                                       NA
## 3                                0.8516012
## 4                                0.9155280
## 5                                0.9939215
## 6                                0.8056548
##   FALSE, FALSE, TRUE, classif.xgboost FALSE, TRUE, FALSE, classif.ksvm
## 1                           0.6205556                       0.9345833
## 2                           0.8093866                       0.9481514
## 3                           0.8521510                       0.9522324
## 4                           0.9282494                       0.9474120
## 5                           0.9426498                       1.0000000
## 6                           0.7994048                       0.9500000
##   FALSE, TRUE, FALSE, classif.randomForest
## 1                                0.9500000
## 2                                0.9614898
## 3                                0.9763968
## 4                                0.9744651
## 5                                1.0000000
## 6                                0.9497024
##   FALSE, TRUE, FALSE, classif.xgboost SMOTE, FALSE, FALSE, classif.ksvm
## 1                           0.9500000                       0.9448830
## 2                           0.9628509                       0.9602482
## 3                           0.9701983                       0.9613247
## 4                           0.9744882                       0.8173401
## 5                           1.0000000                       1.0000000
## 6                           0.9498016                       0.9852235
##   SMOTE, FALSE, FALSE, classif.randomForest
## 1                                 0.9364035
## 2                                 0.9678472
## 3                                 0.9879112
## 4                                 0.9932660
## 5                                 1.0000000
## 6                                 0.9764515
##   SMOTE, FALSE, FALSE, classif.xgboost
## 1                            0.9538012
## 2                            0.9747191
## 3                            0.9873090
## 4                            0.9772727
## 5                            1.0000000
## 6                            0.9720134
```

```r
summary(df)
```

```
##   ADASYN, FALSE, FALSE, classif.ksvm
##   Min.   :0.7408
##   1st Qu.:0.9578
##   Median :0.9792
##   Mean   :0.9640
##   3rd Qu.:0.9949
##   Max.   :1.0000
##   NA's   :1
##   ADASYN, FALSE, FALSE, classif.randomForest
##   Min.   :0.7035
##   1st Qu.:0.9717
##   Median :0.9915
##   Mean   :0.9755
##   3rd Qu.:0.9983
##   Max.   :1.0000
##   NA's   :4
##   ADASYN, FALSE, FALSE, classif.xgboost FALSE, FALSE, FALSE, classif.ksvm
##   Min.   :0.7120                         Min.   :0.9264
##   1st Qu.:0.9674                         1st Qu.:0.9496
##   Median :0.9875                         Median :0.9521
##   Mean   :0.9733                         Mean   :0.9585
##   3rd Qu.:0.9962                         3rd Qu.:0.9649
##   Max.   :1.0000                         Max.   :1.0000
##
##   FALSE, FALSE, FALSE, classif.randomForest
##   Min.   :0.9442
##   1st Qu.:0.9543
##   Median :0.9700
##   Mean   :0.9708
##   3rd Qu.:0.9846
##   Max.   :1.0000
##   NA's   :1
##   FALSE, FALSE, FALSE, classif.xgboost FALSE, FALSE, TRUE, classif.ksvm
##   Min.   :0.9440                        Min.   :0.4146
##   1st Qu.:0.9543                        1st Qu.:0.7063
##   Median :0.9708                        Median :0.9151
##   Mean   :0.9712                        Mean   :0.8332
##   3rd Qu.:0.9856                        3rd Qu.:0.9682
##   Max.   :1.0000                        Max.   :0.9982
##
##   FALSE, FALSE, TRUE, classif.randomForest
##   Min.   :0.4222
##   1st Qu.:0.7694
##   Median :0.8938
##   Mean   :0.8424
##   3rd Qu.:0.9601
##   Max.   :1.0000
##   NA's   :3
##   FALSE, FALSE, TRUE, classif.xgboost FALSE, TRUE, FALSE, classif.ksvm
##   Min.   :0.3600                       Min.   :0.9346
##   1st Qu.:0.7693                       1st Qu.:0.9497
##   Median :0.8715                       Median :0.9521
```

```
## Mean   :0.8369                    Mean   :0.9587
## 3rd Qu.:0.9487                    3rd Qu.:0.9639
## Max.   :1.0000                    Max.   :1.0000
##
## FALSE, TRUE, FALSE, classif.randomForest
## Min.   :0.9442
## 1st Qu.:0.9529
## Median :0.9689
## Mean   :0.9708
## 3rd Qu.:0.9870
## Max.   :1.0000
## NA's   :1
## FALSE, TRUE, FALSE, classif.xgboost SMOTE, FALSE, FALSE, classif.ksvm
## Min.   :0.9440                    Min.   :0.7394
## 1st Qu.:0.9543                    1st Qu.:0.9577
## Median :0.9708                    Median :0.9792
## Mean   :0.9712                    Mean   :0.9622
## 3rd Qu.:0.9851                    3rd Qu.:0.9931
## Max.   :1.0000                    Max.   :1.0000
##
## SMOTE, FALSE, FALSE, classif.randomForest
## Min.   :0.7259
## 1st Qu.:0.9718
## Median :0.9911
## Mean   :0.9754
## 3rd Qu.:0.9978
## Max.   :1.0000
## NA's   :4
## SMOTE, FALSE, FALSE, classif.xgboost
## Min.   :0.6952
## 1st Qu.:0.9679
## Median :0.9874
## Mean   :0.9732
## 3rd Qu.:0.9964
## Max.   :1.0000
##
```

## Verificando a média de cada coluna selecionada

```
for(i in (1:dim(df)[2])){
  print(paste("Media da coluna ", colnames(df)[i], " = ", mean(df[,i], na.rm = TRUE), sep=""))
}
```

```
## [1] "Media da coluna ADASYN, FALSE, FALSE, classif.ksvm = 0.964025236569767"
## [1] "Media da coluna ADASYN, FALSE, FALSE, classif.randomForest = 0.975473207484394"
## [1] "Media da coluna ADASYN, FALSE, FALSE, classif.xgboost = 0.973346351429288"
## [1] "Media da coluna FALSE, FALSE, FALSE, classif.ksvm = 0.958493157114079"
## [1] "Media da coluna FALSE, FALSE, FALSE, classif.randomForest = 0.970841181840409"
## [1] "Media da coluna FALSE, FALSE, FALSE, classif.xgboost = 0.971249595483984"
## [1] "Media da coluna FALSE, FALSE, TRUE, classif.ksvm = 0.833173562108117"
## [1] "Media da coluna FALSE, FALSE, TRUE, classif.randomForest = 0.842420653375876"
## [1] "Media da coluna FALSE, FALSE, TRUE, classif.xgboost = 0.836888855277703"
```

```
## [1] "Media da coluna FALSE, TRUE, FALSE, classif.ksvm = 0.958658104029159"
## [1] "Media da coluna FALSE, TRUE, FALSE, classif.randomForest = 0.97083851619639"
## [1] "Media da coluna FALSE, TRUE, FALSE, classif.xgboost = 0.971153937701804"
## [1] "Media da coluna SMOTE, FALSE, FALSE, classif.ksvm = 0.962176401536434"
## [1] "Media da coluna SMOTE, FALSE, FALSE, classif.randomForest = 0.975446308141522"
## [1] "Media da coluna SMOTE, FALSE, FALSE, classif.xgboost = 0.973221273184089"
```

# Fazendo teste de normalidade

```
plotDensities(data = na.omit(df))
```



# Testando as diferencas

```
friedmanTest(df)
```

```
##
##  Friedman's rank sum test
##
## data:  df
## Friedman's chi-squared = 636.99, df = 14, p-value < 2.2e-16
```

# Testando as diferencas par a par

```
test <- nemenyiTest (df, alpha=0.05)
abs(test$diff.matrix) > test$statistic
```

```
##         ADASYN, FALSE, FALSE, classif.ksvm
##   [1,]                           FALSE
##   [2,]                           FALSE
##   [3,]                           FALSE
##   [4,]                            TRUE
##   [5,]                           FALSE
##   [6,]                           FALSE
##   [7,]                            TRUE
##   [8,]                            TRUE
##   [9,]                            TRUE
##  [10,]                            TRUE
##  [11,]                           FALSE
##  [12,]                           FALSE
##  [13,]                           FALSE
##  [14,]                           FALSE
##  [15,]                           FALSE
##         ADASYN, FALSE, FALSE, classif.randomForest
##   [1,]                                 FALSE
##   [2,]                                 FALSE
##   [3,]                                 FALSE
##   [4,]                                  TRUE
##   [5,]                                  TRUE
##   [6,]                                  TRUE
##   [7,]                                  TRUE
##   [8,]                                  TRUE
##   [9,]                                  TRUE
##  [10,]                                  TRUE
##  [11,]                                  TRUE
##  [12,]                                  TRUE
##  [13,]                                 FALSE
##  [14,]                                 FALSE
##  [15,]                                 FALSE
##         ADASYN, FALSE, FALSE, classif.xgboost
##   [1,]                             FALSE
##   [2,]                             FALSE
##   [3,]                             FALSE
##   [4,]                              TRUE
##   [5,]                              TRUE
##   [6,]                              TRUE
##   [7,]                              TRUE
##   [8,]                              TRUE
##   [9,]                              TRUE
##  [10,]                              TRUE
##  [11,]                              TRUE
##  [12,]                              TRUE
##  [13,]                             FALSE
##  [14,]                             FALSE
##  [15,]                             FALSE
```

```
##         FALSE, FALSE, FALSE, classif.ksvm
##  [1,]                              TRUE
##  [2,]                              TRUE
##  [3,]                              TRUE
##  [4,]                             FALSE
##  [5,]                              TRUE
##  [6,]                              TRUE
##  [7,]                             FALSE
##  [8,]                              TRUE
##  [9,]                              TRUE
## [10,]                             FALSE
## [11,]                              TRUE
## [12,]                              TRUE
## [13,]                              TRUE
## [14,]                              TRUE
## [15,]                              TRUE
##         FALSE, FALSE, FALSE, classif.randomForest
##  [1,]                                     FALSE
##  [2,]                                      TRUE
##  [3,]                                      TRUE
##  [4,]                                      TRUE
##  [5,]                                     FALSE
##  [6,]                                     FALSE
##  [7,]                                      TRUE
##  [8,]                                      TRUE
##  [9,]                                      TRUE
## [10,]                                      TRUE
## [11,]                                     FALSE
## [12,]                                     FALSE
## [13,]                                     FALSE
## [14,]                                      TRUE
## [15,]                                      TRUE
##         FALSE, FALSE, FALSE, classif.xgboost
##  [1,]                                FALSE
##  [2,]                                 TRUE
##  [3,]                                 TRUE
##  [4,]                                 TRUE
##  [5,]                                FALSE
##  [6,]                                FALSE
##  [7,]                                 TRUE
##  [8,]                                 TRUE
##  [9,]                                 TRUE
## [10,]                                 TRUE
## [11,]                                FALSE
## [12,]                                FALSE
## [13,]                                FALSE
## [14,]                                 TRUE
## [15,]                                 TRUE
##         FALSE, FALSE, TRUE, classif.ksvm
##  [1,]                             TRUE
##  [2,]                             TRUE
##  [3,]                             TRUE
##  [4,]                            FALSE
##  [5,]                             TRUE
```

```
##  [6,]                                  TRUE
##  [7,]                                  FALSE
##  [8,]                                  FALSE
##  [9,]                                  FALSE
## [10,]                                  FALSE
## [11,]                                  TRUE
## [12,]                                  TRUE
## [13,]                                  TRUE
## [14,]                                  TRUE
## [15,]                                  TRUE
##         FALSE, FALSE, TRUE, classif.randomForest
##  [1,]                                       TRUE
##  [2,]                                       TRUE
##  [3,]                                       TRUE
##  [4,]                                       TRUE
##  [5,]                                       TRUE
##  [6,]                                       TRUE
##  [7,]                                       FALSE
##  [8,]                                       FALSE
##  [9,]                                       FALSE
## [10,]                                       TRUE
## [11,]                                       TRUE
## [12,]                                       TRUE
## [13,]                                       TRUE
## [14,]                                       TRUE
## [15,]                                       TRUE
##         FALSE, FALSE, TRUE, classif.xgboost FALSE, TRUE, FALSE, classif.ksvm
##  [1,]                                  TRUE                             TRUE
##  [2,]                                  TRUE                             TRUE
##  [3,]                                  TRUE                             TRUE
##  [4,]                                  TRUE                            FALSE
##  [5,]                                  TRUE                             TRUE
##  [6,]                                  TRUE                             TRUE
##  [7,]                                 FALSE                            FALSE
##  [8,]                                 FALSE                             TRUE
##  [9,]                                 FALSE                             TRUE
## [10,]                                  TRUE                            FALSE
## [11,]                                  TRUE                             TRUE
## [12,]                                  TRUE                             TRUE
## [13,]                                  TRUE                             TRUE
## [14,]                                  TRUE                             TRUE
## [15,]                                  TRUE                             TRUE
##         FALSE, TRUE, FALSE, classif.randomForest
##  [1,]                                     FALSE
##  [2,]                                      TRUE
##  [3,]                                      TRUE
##  [4,]                                      TRUE
##  [5,]                                     FALSE
##  [6,]                                     FALSE
##  [7,]                                      TRUE
##  [8,]                                      TRUE
##  [9,]                                      TRUE
## [10,]                                      TRUE
## [11,]                                     FALSE
```

```
## [12,]                                        FALSE
## [13,]                                        FALSE
## [14,]                                         TRUE
## [15,]                                         TRUE
##       FALSE, TRUE, FALSE, classif.xgboost
##  [1,]                                        FALSE
##  [2,]                                         TRUE
##  [3,]                                         TRUE
##  [4,]                                         TRUE
##  [5,]                                        FALSE
##  [6,]                                        FALSE
##  [7,]                                         TRUE
##  [8,]                                         TRUE
##  [9,]                                         TRUE
## [10,]                                         TRUE
## [11,]                                        FALSE
## [12,]                                        FALSE
## [13,]                                        FALSE
## [14,]                                         TRUE
## [15,]                                         TRUE
##       SMOTE, FALSE, FALSE, classif.ksvm
##  [1,]                                        FALSE
##  [2,]                                        FALSE
##  [3,]                                        FALSE
##  [4,]                                         TRUE
##  [5,]                                        FALSE
##  [6,]                                        FALSE
##  [7,]                                         TRUE
##  [8,]                                         TRUE
##  [9,]                                         TRUE
## [10,]                                         TRUE
## [11,]                                        FALSE
## [12,]                                        FALSE
## [13,]                                        FALSE
## [14,]                                        FALSE
## [15,]                                        FALSE
##       SMOTE, FALSE, FALSE, classif.randomForest
##  [1,]                                        FALSE
##  [2,]                                        FALSE
##  [3,]                                        FALSE
##  [4,]                                         TRUE
##  [5,]                                         TRUE
##  [6,]                                         TRUE
##  [7,]                                         TRUE
##  [8,]                                         TRUE
##  [9,]                                         TRUE
## [10,]                                         TRUE
## [11,]                                         TRUE
## [12,]                                         TRUE
## [13,]                                        FALSE
## [14,]                                        FALSE
## [15,]                                        FALSE
##       SMOTE, FALSE, FALSE, classif.xgboost
##  [1,]                                        FALSE
```

```
## [2,]                                    FALSE
## [3,]                                    FALSE
## [4,]                                     TRUE
## [5,]                                     TRUE
## [6,]                                     TRUE
## [7,]                                     TRUE
## [8,]                                     TRUE
## [9,]                                     TRUE
## [10,]                                    TRUE
## [11,]                                    TRUE
## [12,]                                    TRUE
## [13,]                                   FALSE
## [14,]                                   FALSE
## [15,]                                   FALSE
```

# Plotando os ranks

```r
print(colMeans(rankMatrix(df)))
```

```
##          ADASYN, FALSE, FALSE, classif.ksvm
##                                    6.347561
## ADASYN, FALSE, FALSE, classif.randomForest
##                                    4.048780
##       ADASYN, FALSE, FALSE, classif.xgboost
##                                    4.384146
##           FALSE, FALSE, FALSE, classif.ksvm
##                                   10.365854
##  FALSE, FALSE, FALSE, classif.randomForest
##                                    7.676829
##       FALSE, FALSE, FALSE, classif.xgboost
##                                    7.378049
##            FALSE, FALSE, TRUE, classif.ksvm
##                                   12.646341
##   FALSE, FALSE, TRUE, classif.randomForest
##                                   13.225610
##        FALSE, FALSE, TRUE, classif.xgboost
##                                   13.573171
##            FALSE, TRUE, FALSE, classif.ksvm
##                                   10.396341
##   FALSE, TRUE, FALSE, classif.randomForest
##                                    7.829268
##        FALSE, TRUE, FALSE, classif.xgboost
##                                    7.548780
##           SMOTE, FALSE, FALSE, classif.ksvm
##                                    6.079268
##  SMOTE, FALSE, FALSE, classif.randomForest
##                                    4.012195
##       SMOTE, FALSE, FALSE, classif.xgboost
##                                    4.487805
```

# Plotando grafico de Critical Diference

```
result = tryCatch({
    plotCD(df, alpha=0.05, cex = 0.35)
}, error = function(e) {})
```