# R Notebook

## Parametros:

**Measure =**  F1 measure
**Columns =**  sampling, weight_space, underbagging, learner
**Performance =**  tuning_measure
**Filter keys =**  imba.rate
**Filter values =**  0.05

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

## Tratamento dos dados

Carregando data set compilado

```r
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                      :   0    ADASYN: 738   Mode :logical
##  Area under the curve          :   0    FALSE :2214   FALSE:2952
##  F1 measure                    :3690    SMOTE : 738   TRUE :738
##  G-mean                        :   0                  NA's :0
##  Matthews correlation coefficient:  0
##
##
##  tuning_measure   holdout_measure   holdout_measure_residual
##  Min.   :0.0000   Min.   :0.0000    Min.   :0.00000
##  1st Qu.:0.3333   1st Qu.:0.1000    1st Qu.:0.07022
##  Median :0.8198   Median :0.5000    Median :0.32530
##  Mean   :0.6671   Mean   :0.4905    Mean   :0.39891
##  3rd Qu.:0.9848   3rd Qu.:0.8333    3rd Qu.:0.73016
##  Max.   :1.0000   Max.   :1.0000    Max.   :1.00000
##  NA's   :51       NA's   :51        NA's   :51
##  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   :51      (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

```
# 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
```

```
# Renomeando a variavel
df = df_tec_wide_residual

head(df)
```

```
##   ADASYN, FALSE, FALSE, classif.ksvm
## 1                         0.9472249
## 2                         0.9611172
## 3                         0.9511782
## 4                         0.8095020
## 5                         1.0000000
## 6                         0.9769593
##   ADASYN, FALSE, FALSE, classif.randomForest
## 1                                  0.9401067
## 2                                         NA
## 3                                  0.9932657
## 4                                  0.9932778
## 5                                  1.0000000
## 6                                  0.9750455
##   ADASYN, FALSE, FALSE, classif.xgboost FALSE, FALSE, FALSE, classif.ksvm
## 1                             0.9576632                        0.1290537
## 2                             0.9733028                        0.2772359
## 3                             0.9882364                        0.2927554
## 4                             0.9842866                        0.0000000
## 5                             1.0000000                        1.0000000
## 6                             0.9733838                        0.1247397
##   FALSE, FALSE, FALSE, classif.randomForest
## 1                               0.005291005
## 2                               0.454977821
```

```
## 3                              0.727298407
## 4                              0.781040564
## 5                              1.000000000
## 6                              0.177957494
##   FALSE, FALSE, FALSE, classif.xgboost FALSE, FALSE, TRUE, classif.ksvm
## 1                         0.02289746                         0.1485313
## 2                         0.46619116                         0.2542338
## 3                         0.67464110                         0.3086281
## 4                         0.70432099                         0.1271359
## 5                         1.00000000                         0.9025814
## 6                         0.22381693                         0.1623156
##   FALSE, FALSE, TRUE, classif.randomForest
## 1                              0.1649924
## 2                                     NA
## 3                              0.4162512
## 4                              0.5370950
## 5                              0.9489338
## 6                              0.2908245
##   FALSE, FALSE, TRUE, classif.xgboost FALSE, TRUE, FALSE, classif.ksvm
## 1                         0.1616499                         0.12382963
## 2                         0.3103932                         0.28566897
## 3                         0.3827636                         0.32873851
## 4                         0.5873374                         0.00000000
## 5                         0.7294844                         1.00000000
## 6                         0.2858927                         0.09682558
##   FALSE, TRUE, FALSE, classif.randomForest
## 1                              0.005291005
## 2                                      NA
## 3                              0.722524478
## 4                              0.651940035
## 5                              1.000000000
## 6                              0.177957494
##   FALSE, TRUE, FALSE, classif.xgboost SMOTE, FALSE, FALSE, classif.ksvm
## 1                         0.02420406                         0.9464710
## 2                         0.45160605                         0.9596461
## 3                         0.64215249                         0.9595356
## 4                         0.69400353                         0.7646842
## 5                         1.00000000                         1.0000000
## 6                         0.21351370                         0.9851804
##   SMOTE, FALSE, FALSE, classif.randomForest
## 1                              0.9412533
## 2                                     NA
## 3                              0.9877630
## 4                              0.9933152
## 5                              1.0000000
## 6                              0.9766771
##   SMOTE, FALSE, FALSE, classif.xgboost
## 1                              0.9540043
## 2                              0.9738483
## 3                              0.9896043
## 4                              0.9778476
## 5                              1.0000000
## 6                              0.9720548
```

```r
summary(df)
```

```
##  ADASYN, FALSE, FALSE, classif.ksvm
##  Min.   :0.7913
##  1st Qu.:0.9544
##  Median :0.9791
##  Mean   :0.9630
##  3rd Qu.:0.9947
##  Max.   :1.0000
##  NA's   :1
##  ADASYN, FALSE, FALSE, classif.randomForest
##  Min.   :0.7179
##  1st Qu.:0.9739
##  Median :0.9920
##  Mean   :0.9770
##  3rd Qu.:0.9984
##  Max.   :1.0000
##  NA's   :5
##  ADASYN, FALSE, FALSE, classif.xgboost FALSE, FALSE, FALSE, classif.ksvm
##  Min.   :0.7310                         Min.   :0.0000
##  1st Qu.:0.9666                         1st Qu.:0.0000
##  Median :0.9876                         Median :0.2281
##  Mean   :0.9720                         Mean   :0.2947
##  3rd Qu.:0.9962                         3rd Qu.:0.4719
##  Max.   :1.0000                         Max.   :1.0000
##
##  FALSE, FALSE, FALSE, classif.randomForest
##  Min.   :0.0000
##  1st Qu.:0.1984
##  Median :0.6329
##  Mean   :0.5394
##  3rd Qu.:0.8303
##  Max.   :1.0000
##  NA's   :1
##  FALSE, FALSE, FALSE, classif.xgboost FALSE, FALSE, TRUE, classif.ksvm
##  Min.   :0.0000                        Min.   :0.09246
##  1st Qu.:0.2468                        1st Qu.:0.19564
##  Median :0.6565                        Median :0.39808
##  Mean   :0.5765                        Mean   :0.43307
##  3rd Qu.:0.8417                        3rd Qu.:0.63027
##  Max.   :1.0000                        Max.   :0.98214
##
##  FALSE, FALSE, TRUE, classif.randomForest
##  Min.   :0.09861
##  1st Qu.:0.26788
##  Median :0.47097
##  Mean   :0.49105
##  3rd Qu.:0.71859
##  Max.   :1.00000
##  NA's   :3
##  FALSE, FALSE, TRUE, classif.xgboost FALSE, TRUE, FALSE, classif.ksvm
##  Min.   :0.09351                      Min.   :0.0000
##  1st Qu.:0.23896                      1st Qu.:0.0000
##  Median :0.43019                      Median :0.2065
```

```
##   Mean   :0.46400                      Mean   :0.2870
##   3rd Qu.:0.68242                      3rd Qu.:0.4463
##   Max.   :1.00000                      Max.   :1.0000
##
##   FALSE, TRUE, FALSE, classif.randomForest
##   Min.   :0.0000
##   1st Qu.:0.2506
##   Median :0.5752
##   Mean   :0.5486
##   3rd Qu.:0.8651
##   Max.   :1.0000
##   NA's   :3
##   FALSE, TRUE, FALSE, classif.xgboost SMOTE, FALSE, FALSE, classif.ksvm
##   Min.   :0.0000                      Min.   :0.7607
##   1st Qu.:0.2960                      1st Qu.:0.9546
##   Median :0.6406                      Median :0.9784
##   Mean   :0.5741                      Mean   :0.9599
##   3rd Qu.:0.8384                      3rd Qu.:0.9931
##   Max.   :1.0000                      Max.   :1.0000
##
##   SMOTE, FALSE, FALSE, classif.randomForest
##   Min.   :0.7311
##   1st Qu.:0.9730
##   Median :0.9915
##   Mean   :0.9772
##   3rd Qu.:0.9978
##   Max.   :1.0000
##   NA's   :4
##   SMOTE, FALSE, FALSE, classif.xgboost
##   Min.   :0.6908
##   1st Qu.:0.9684
##   Median :0.9877
##   Mean   :0.9747
##   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 = 0.963045110849921"
## [1] "Media da coluna ADASYN, FALSE, FALSE, classif.randomForest = 0.976993535327279"
## [1] "Media da coluna ADASYN, FALSE, FALSE, classif.xgboost = 0.971972609222412"
## [1] "Media da coluna FALSE, FALSE, FALSE, classif.ksvm = 0.294747069679735"
## [1] "Media da coluna FALSE, FALSE, FALSE, classif.randomForest = 0.539379103154103"
## [1] "Media da coluna FALSE, FALSE, FALSE, classif.xgboost = 0.576466997099341"
## [1] "Media da coluna FALSE, FALSE, TRUE, classif.ksvm = 0.43307502070709"
## [1] "Media da coluna FALSE, FALSE, TRUE, classif.randomForest = 0.491051630464319"
## [1] "Media da coluna FALSE, FALSE, TRUE, classif.xgboost = 0.464002478013471"
```

```
## [1] "Media da coluna FALSE, TRUE, FALSE, classif.ksvm = 0.287034479003018"
## [1] "Media da coluna FALSE, TRUE, FALSE, classif.randomForest = 0.548585456558965"
## [1] "Media da coluna FALSE, TRUE, FALSE, classif.xgboost = 0.574089095970635"
## [1] "Media da coluna SMOTE, FALSE, FALSE, classif.ksvm = 0.959902879940769"
## [1] "Media da coluna SMOTE, FALSE, FALSE, classif.randomForest = 0.977181342299523"
## [1] "Media da coluna SMOTE, FALSE, FALSE, classif.xgboost = 0.974658612182877"
```

# 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 = 770.48, 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,]                           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.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,]                               FALSE
##  [9,]                               FALSE
## [10,]                               FALSE
## [11,]                                TRUE
## [12,]                                TRUE
## [13,]                                TRUE
## [14,]                                TRUE
## [15,]                                TRUE
##       FALSE, FALSE, FALSE, classif.randomForest
##  [1,]                                        TRUE
##  [2,]                                        TRUE
##  [3,]                                        TRUE
##  [4,]                                        TRUE
##  [5,]                                       FALSE
##  [6,]                                       FALSE
##  [7,]                                       FALSE
##  [8,]                                       FALSE
##  [9,]                                       FALSE
## [10,]                                        TRUE
## [11,]                                       FALSE
## [12,]                                       FALSE
## [13,]                                        TRUE
## [14,]                                        TRUE
## [15,]                                        TRUE
##       FALSE, FALSE, FALSE, classif.xgboost
##  [1,]                                   TRUE
##  [2,]                                   TRUE
##  [3,]                                   TRUE
##  [4,]                                   TRUE
##  [5,]                                  FALSE
##  [6,]                                  FALSE
##  [7,]                                  FALSE
##  [8,]                                  FALSE
##  [9,]                                  FALSE
## [10,]                                   TRUE
## [11,]                                  FALSE
## [12,]                                  FALSE
## [13,]                                   TRUE
## [14,]                                   TRUE
## [15,]                                   TRUE
##       FALSE, FALSE, TRUE, classif.ksvm
##  [1,]                               TRUE
##  [2,]                               TRUE
##  [3,]                               TRUE
##  [4,]                              FALSE
##  [5,]                              FALSE
```

9

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

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

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

```
##          ADASYN, FALSE, FALSE, classif.ksvm
##                                    4.225610
## ADASYN, FALSE, FALSE, classif.randomForest
##                                    3.560976
##       ADASYN, FALSE, FALSE, classif.xgboost
##                                    3.841463
##           FALSE, FALSE, FALSE, classif.ksvm
##                                   12.652439
##  FALSE, FALSE, FALSE, classif.randomForest
##                                   10.109756
##        FALSE, FALSE, FALSE, classif.xgboost
##                                    9.012195
##            FALSE, FALSE, TRUE, classif.ksvm
##                                   11.195122
##    FALSE, FALSE, TRUE, classif.randomForest
##                                   10.603659
##        FALSE, FALSE, TRUE, classif.xgboost
##                                   11.219512
##            FALSE, TRUE, FALSE, classif.ksvm
##                                   12.786585
##    FALSE, TRUE, FALSE, classif.randomForest
##                                   10.176829
##        FALSE, TRUE, FALSE, classif.xgboost
##                                    9.219512
##           SMOTE, FALSE, FALSE, classif.ksvm
##                                    4.042683
##   SMOTE, FALSE, FALSE, classif.randomForest
##                                    3.560976
##       SMOTE, FALSE, FALSE, classif.xgboost
##                                    3.792683
```

# Plotando grafico de Critical Diference

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