

R Notebook

Parametros:

Measure = Matthews correlation coefficient
Columns = sampling, weight_space, underbagging
Performance = holdout_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      cardiocography-10clases: 900 3rd Qu.:0.0500
## Max. :3      cardiocography-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           :  0  SMOTE : 738  TRUE :738
## G-mean              :  0              NA's :0
## Matthews correlation coefficient:3690
##
##
## tuning_measure  holdout_measure  holdout_measure_residual
## Min.   :-0.1277  Min.   :-0.21201  Min.   :-0.45710
## 1st Qu.: 0.3764  1st Qu.: 0.06131  1st Qu.: 0.05637
## Median : 0.8057  Median : 0.55190  Median : 0.23378
## Mean   : 0.6629  Mean   : 0.49274  Mean   : 0.32193
## 3rd Qu.: 0.9728  3rd Qu.: 0.82456  3rd Qu.: 0.56442
## Max.   : 1.0000  Max.   : 1.00000  Max.   : 1.00000
## NA's   :54      NA's   :54      NA's   :54
## 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    :54      (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$performance)))

# 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] 246 5
```

```
# Renomeando a variavel
df = df_tec_wide_residual

head(df)
```

```
## ADASYN, FALSE, FALSE FALSE, FALSE, FALSE FALSE, FALSE, TRUE
## 1 0.05919056 0.1037074 0.1222062
## 2 0.14927451 0.2609339 0.2801961
## 3 0.44081601 0.3065053 0.4408298
## 4 0.00000000 0.0000000 0.1702715
## 5 1.00000000 1.0000000 0.9530776
## 6 0.01727820 0.2013082 0.1850007
## FALSE, TRUE, FALSE SMOTE, FALSE, FALSE
## 1 0.07585303 0.09021363
## 2 0.29437987 0.18853292
## 3 0.32667105 0.36382909
## 4 0.00000000 0.00000000
## 5 1.00000000 0.79290891
## 6 0.14506241 0.08378119
```

```
summary(df)
```

```
## ADASYN, FALSE, FALSE FALSE, FALSE, FALSE FALSE, FALSE, TRUE
## Min. :-0.06657 Min. :-0.04044 Min. :-0.06927
## 1st Qu.: 0.19033 1st Qu.: 0.05938 1st Qu.: 0.24619
## Median : 0.55092 Median : 0.55293 Median : 0.47866
## Mean : 0.50783 Mean : 0.48350 Mean : 0.48099
## 3rd Qu.: 0.83178 3rd Qu.: 0.83361 3rd Qu.: 0.70069
## Max. : 1.00000 Max. : 1.00000 Max. : 1.00000
```

```
## NA's      :8                      NA's      :3
## FALSE, TRUE, FALSE SMOTE, FALSE, FALSE
## Min.      :-0.03836   Min.      :-0.0748
## 1st Qu.: 0.06407   1st Qu.: 0.1618
## Median : 0.54013   Median : 0.5248
## Mean      : 0.47801   Mean      : 0.5139
## 3rd Qu.: 0.81094   3rd Qu.: 0.8476
## Max.      : 1.00000   Max.      : 1.0000
## NA's      :3          NA's      :4
```

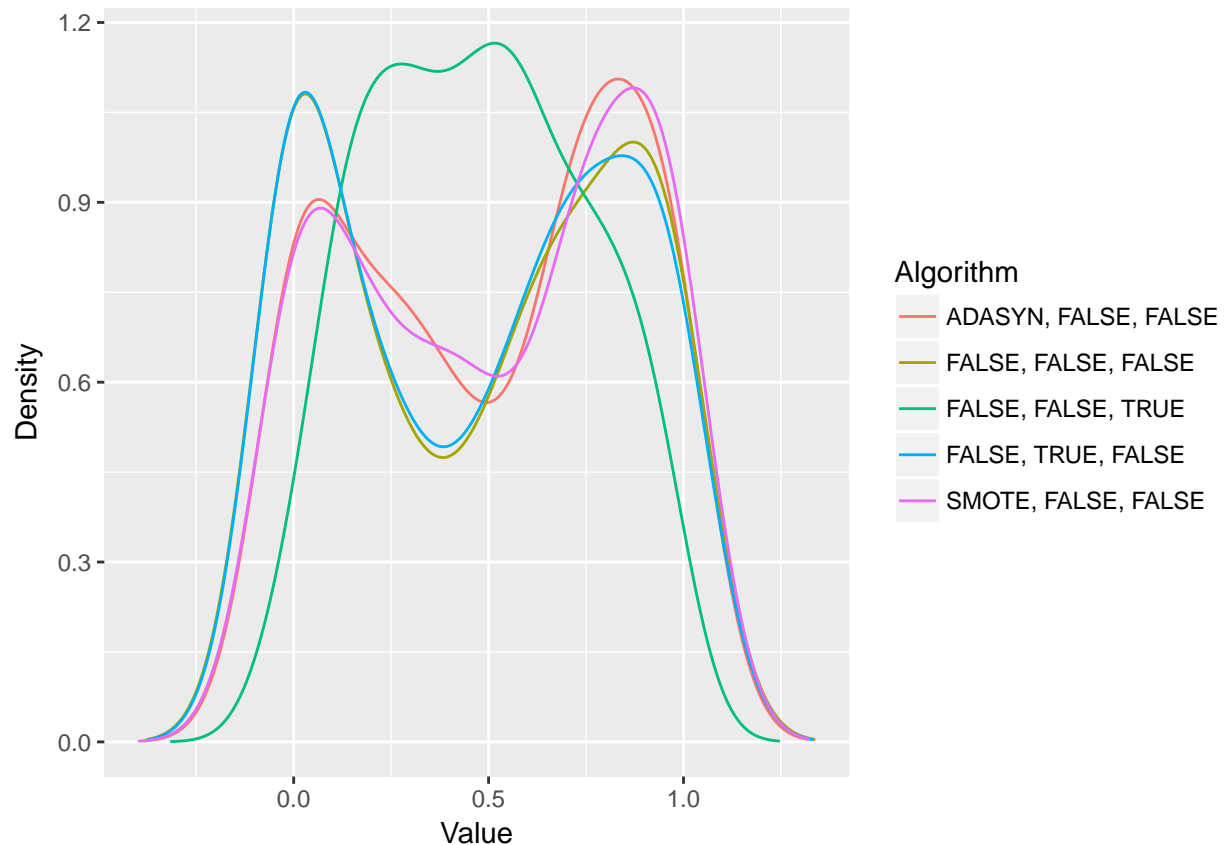
Verificando a média de cada coluna selecionada

```
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 = 0.507832450493462"
## [1] "Media da coluna FALSE, FALSE, FALSE = 0.483499945881345"
## [1] "Media da coluna FALSE, FALSE, TRUE = 0.480994553192514"
## [1] "Media da coluna FALSE, TRUE, FALSE = 0.478011885054923"
## [1] "Media da coluna SMOTE, FALSE, FALSE = 0.513885895989575"
```

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 = 21.718, df = 4, p-value = 0.0002281
```

Testando as diferencas par a par

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

```
##      ADASYN, FALSE, FALSE FALSE, FALSE, FALSE FALSE, FALSE, TRUE
## [1,]      FALSE      FALSE      TRUE
## [2,]      FALSE      FALSE      FALSE
## [3,]      TRUE      FALSE      FALSE
## [4,]      FALSE      FALSE      FALSE
## [5,]      FALSE      TRUE      TRUE
##      FALSE, TRUE, FALSE SMOTE, FALSE, FALSE
## [1,]      FALSE      FALSE
```

```
## [2,]          FALSE          TRUE
## [3,]          FALSE          TRUE
## [4,]          FALSE          TRUE
## [5,]          TRUE           FALSE
```

Plotando os ranks

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

```
## ADASYN, FALSE, FALSE  FALSE, FALSE, FALSE  FALSE, FALSE, TRUE
##          2.819106          3.109756          3.237805
##  FALSE, TRUE, FALSE  SMOTE, FALSE, FALSE
##          3.146341          2.686992
```

Plotando grafico de Critical Difference

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

