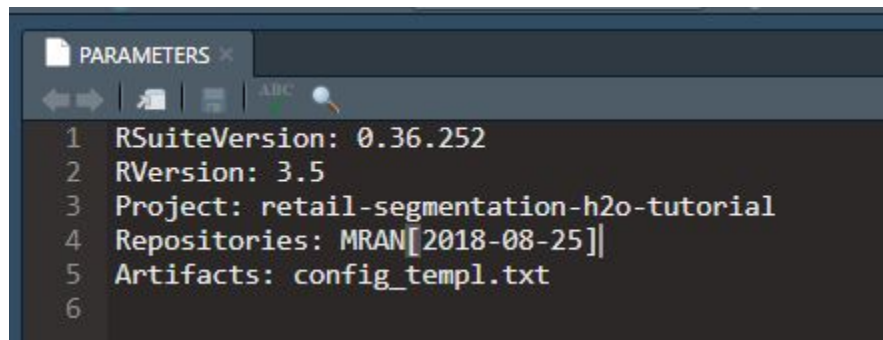# Retail Segmentation with *h2o*

# Initial steps

- Create the project
- Add the packages
- Set the repository reference date
- Set the debug level
- Add the data folder and copy the raw data
- Quick test running *master.R*
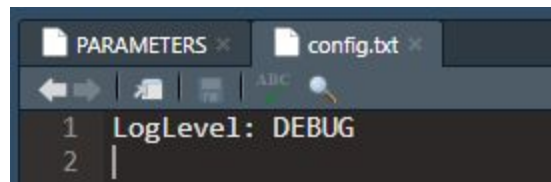
# PARAMETERS

Change the MRAN date

# config.txt

Set the log level to DEBUG

# Add data files

# Run master.R

- We run the scripts from the terminal

```
R:\rsuite-projects\retail-segmentation-h2o-tutorial>Rscript R/master.R
```

# Next steps

- Package
  - Add imports to `DESCRIPTION`
  - Add imports to `packages_import.R`
  - Add 1st script `find_best_model.R`
- Project
  - add 1st script `build_p2b_nosegmentation_model.R`
- Terminal
  - Install dependencies and build project
  - Run script

# terminal: Install and build

- Install dependencies
- Build project

```
R:\rsuite-projects\retail-segmentation-h2o-tutorial>rsuite proj depsinst
```

```
R:\rsuite-projects\retail-segmentation-h2o-tutorial>rsuite proj build
?[0m2019-05-04 13:00:15 INFO:rsuite:Installing segmentationmodels (for R 3.5) ...?
[0m?[0m?[0m
?[0m2019-05-04 13:00:21 INFO:rsuite:Successfuly installed 1 packages?[0m?[0m?[0m
```

# Add package **h2o**

# terminal: Install dependencies and build project



```
R:\rsuite-projects\retail-segmentation-h2o-tutorial>rsuite proj depsinst
?[0m2019-05-04 13:04:38 INFO:rsuite:Detecting repositories (for R 3.5)...?[0m?[0m?[0m
?[0m2019-05-04 13:04:38 INFO:rsuite:Will look for dependencies in ...?[0m?[0m?[0m
?[0m2019-05-04 13:04:38 INFO:rsuite:.          MRAN#1 = https://mran.microsoft.com/snapshot/2018-
08-25 (win.binary, source)?[0m?[0m?[0m
?[0m2019-05-04 13:04:38 INFO:rsuite:Collecting project dependencies (for R 3.5)...?[0m?[0m?[0m
?[0m2019-05-04 13:04:38 INFO:rsuite:Resolving dependencies (for R 3.5)...?[0m
?[0m2019-05-04 13:04:39 INFO:rsuite:Following installed packages will be upda
els?[0m?[0m?[0m
?[0m2019-05-04 13:04:40 INFO:rsuite:Detected 4 dependencies to install. Insta

?[0m2019-05-04 13:04:46 INFO:rsuite:All dependencies successfully installed.?
```

```
R:\rsuite-projects\retail-segmentation-h2o-tutorial>rsuite proj build
?[0m2019-05-04 13:05:27 INFO:rsuite:Installing segmentationmodels (for R 3.5)
?[0m2019-05-04 13:05:33 INFO:rsuite:Successfuly installed 1 packages?[0m?[0m
Warning message:
In if (nchar(msg) > 8192) { :
  the condition has length > 1 and only the first element will be used

R:\rsuite-projects\retail-segmentation-h2o-tutorial>
```

o-tutorial ▶ deployment ▶ libs ▶

| # | Name ^ |
|---|--------|
| 1 | bitops |
| 2 | h2o |
| 3 | jsonlite |
| 4 | logging |
| 5 | RCurl |
| 6 | segmentationmodels |
| 7 | .gitignore |

## project: Add the first script to the project

build_p2b_nosegmentation_model.R

# package: Add 1st script - find best model

find_best_model.R

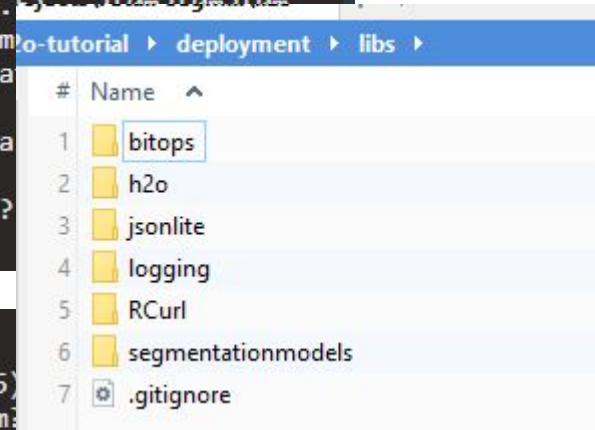# terminal: Install dependencies and build project



```
R:\rsuite-projects\retail-segmentation-h2o-tutorial>rsuite proj depsinst
?[0m2019-05-04 13:14:46 INFO:rsuite:Detecting repositories (for R 3.5)...?[0m?[0m?[0m
?[0m2019-05-04 13:14:46 INFO:rsuite:Will look for dependencies in ...?[0m?[0m?[0m
?[0m2019-05-04 13:14:46 INFO:rsuite:.        MRAN#1 = https://mran.microsoft.com/snapshot
/2018-08-25 (win.binary, source)?[0m?[0m?[0m
?[0m2019-05-04 13:14:46 INFO:rsuite:Collecting project dependencies (for R 3.5)...?[0m?[0m?
[0m
?[0m2019-05-04 13:14:46 INFO:rsuite:Resolving dependencies (for R 3.5)...?[0m?[0m?[0m
?[0m2019-05-04 13:14:47 INFO:rsuite:Following installed packages will be updated: segmentat
ionmodels?[0m?[0m?[0m
?[0m2019-05-04 13:14:47 INFO:rsuite:Detected 1 dependencies to install. Installing...?[0m?[
0m?[0m
?[0m2019-05-04 13:14:49 INFO:rsuite:All dependencies successfully installed.?[0m?[0m?[0m

R:\rsuite-projects\retail-segmentation-h2o-tutorial>
```

```
R:\rsuite-projects\retail-segmentation-h2o-tutorial>rsuite proj build
?[0m2019-05-04 13:15:47 INFO:rsuite:Installing segmentationmodels (for R 3.5) ...?[0m?[0m?[
0m
?[0m2019-05-04 13:15:54 INFO:rsuite:Successfuly installed 1 packages?[0m?[0m?[0m

R:\rsuite-projects\retail-segmentation-h2o-tutorial>
```

# new packages have been added

# terminal: Run the *no-segmentation model*

- It will run h2o but will stop on error because the package pROC is missing
- Install pROC outside the project (global) because pROC is called by the project not the package

```
> install.packages("pROC")
Installing package into 'R:/rsuite-projects/retail-segmentation-h2o-tutorial
x'
(as 'lib' is unspecified)
trying URL 'https://cran.rstudio.com/bin/windows/contrib/3.5/pROC_1.14.0.zip
Content type 'application/zip' length 1165753 bytes (1.1 MB)
downloaded 1.1 MB

package 'pROC' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
        C:\Users\msfz751\AppData\Local\Temp\Rtmpc9kZgw\downloaded_packages
>
```

# terminal: Run the same model again ...

- No it runs h2o and rPROC



```
R:\rsuite-projects\retail-segmentation-h2o-tutorial>Rscript R/build_p2b_nosegmentation_model.R
```

```
Warning message:
In h2o.clusterInfo() :
Your H2O cluster version is too old (10 months and 18 days)!
Please download and install the latest version from http://h2o.ai/download/
[1] 0
2019-05-04 13:23:05 INFO::--> H2O started
    |========================================================================| 100%
    |========================================================================| 100%
2019-05-04 13:23:09 INFO::--> Datasets imported into H2O cluster
    |========================================================================| 100%
[1] "R:\\rsuite-projects\\retail-segmentation-h2o-tutorial\\export\\glm_grid_model_0"
2019-05-04 13:23:11 INFO::--> Best model exported into export folder
2019-05-04 13:23:11 INFO::--> Best model with test AUC=0.646072656639073
    |========================================================================| 100%

R:\rsuite-projects\retail-segmentation-h2o-tutorial>
```

# Next steps

- Package
  - Add 2nd script `build_segmentation_models.R`
  - Add 3rd script `predict_segmentation_models.R`
- Project
  - add 2nd script `build_p2b_nosegmentation_local_models.R`
- Terminal
  - Install dependencies and build project
  - Run script `R/build_p2b_nosegmentation_local_models.R`

# package: Add the 2nd script - build models

```r
 build_segmentation_models.R ×
1    #' @export
2    build_segmentation_models <- function(training_frame, segmentation_vars, cluster_
3      segmentation_models <- list()
4
5      for (cluster_cnt in cluster_cnts) {
6        best_model <- NULL
7        for (round in 1:rounds) {
8          segmentation_model <- h2o.kmeans(training_frame = training_frame,
9                                           x = segmentation_vars,
10                                          k = cluster_cnt,
11                                          model_id = sprintf("segmentation_model_%s"
12                                          init = "PlusPlus",
13                                          standardize = TRUE)
14          model_withinss <- h2o.tot_withinss(segmentation_model)
15          model_betweenss <- h2o.betweenss(segmentation_model)
16
17          if (is.null(best_model)) {
18            best_model <- list(
19              segmentation_model = segmentation_model,
20              tot_withinss = model_withinss,
21              betweenss = model_betweenss)
22          } else if (best_model$tot_withinss/best_model$betweenss >
23                     model_withinss/model_betweenss) {
24            best_model <- list(
```

# package: Add the 3rd script - prediction

```r
predict_segmentation_models.R  ×

1   #'@export
2   predict_segmentation_models <- function(segmentation_models, train_df, test_df) {
3       lapply(X = segmentation_models,
4               FUN = function(segmentation_model) {
5                   list(
6                       k = segmentation_model@parameters$k,
7                       segment_train = h2o.assign(h2o.predict(segmentation_model, newdata = train_df),
8                                                   key = sprintf("retail_train_segment_assignment_k_%s",
9                                                       segmentation_model@parameters$k)),
10                      segment_test = h2o.assign(h2o.predict(segmentation_model, newdata = test_df),
11                                                  key = sprintf("retail_test_segment_assignment_k_%s",
12                                                      segmentation_model@parameters$k))
13                  )
14              })
15  }
16
```

# project: Add script for segmentation of local models



```
 R  build_p2b_segmentation_local_models.R   ×

14    # Setting .libPaths() to point to libs folder
15    source(file.path(script_path, "set_env.R"), chdir = T)
16
17    config <- load_config()
18    args <- args_parser()
19
20    ###############################################################################
21
22    library(data.table)
23    library(h2o)
24    library(logging)
25
26    h2o_local <- h2o.init(nthreads = 4,
27                          max_mem_size = "6g")
28    h2o.removeAll()
29
30    loginfo("--> H2O started")
31
32    set.seed(1234)
33
34    retail_train <- h2o.importFile(path = file.path(script_path, "../data/retail_train.csv"),
35                                   destination_frame = "retail_train",
```

# terminal: install dependencies and build project

- the new package scripts will be added

```
R:\rsuite-projects\retail-segmentation-h2o-tutorial>rsuite proj depsinst
?[0m2019-05-04 13:32:12 INFO:rsuite:Detecting repositories (for R 3.5)...?[0m?[0m?[0m
?[0m2019-05-04 13:32:12 INFO:rsuite:Will look for dependencies in ...?[0m?[0m?[0m
?[0m2019-05-04 13:32:12 INFO:rsuite:.          MRAN#1 = https://mran.microsoft.com/snapshot/2018-
08-25 (win.binary, source)?[0m?[0m?[0m
?[0m2019-05-04 13:32:12 INFO:rsuite:Collecting project dependencies (for R 3.5)...?[0m?[0m?[0m
?[0m2019-05-04 13:32:12 INFO:rsuite:Resolving dependencies (for R 3.5)...?[0m?[0m?[0m
?[0m2019-05-04 13:32:13 INFO:rsuite:Following installed packages will be updated: segmentationmod
els?[0m?[0m?[0m
?[0m2019-05-04 13:32:13 INFO:rsuite:No dependencies to install.?[0m?[0m?[0m

R:\rsuite-projects\retail-segmentation-h2o-tutorial>rsuite proj build
?[0m2019-05-04 13:32:22 INFO:rsuite:Installing segmentationmodels (for R 3.5) ...?[0m?[0m?[0m
?[0m2019-05-04 13:32:28 INFO:rsuite:Successfuly installed 1 packages?[0m?[0m?[0m

R:\rsuite-projects\retail-segmentation-h2o-tutorial>
```

# terminal: Run segmentation for local models

- Results from h2o and pROC

# project: Add script for segmentation of *standard* models



```r
build_p2b_segmentation_model.R ×
14    # Setting .libPaths() to point to libs folder
15    source(file.path(script_path, "set_env.R"), chdir = T)
16
17    config <- load_config()
18    args <- args_parser()
19
20    ################################################################################
21
22    library(data.table)
23    library(h2o)
24    library(logging)
25
26    h2o_local <- h2o.init(nthreads = 4,
27                          max_mem_size = "6g")
28    h2o.removeAll()
29
30    loginfo("--> H2O started")
31
32    set.seed(1234)
33
34    retail_train <- h2o.importFile(path = file.path(script_path, "..", "data/retail_train.csv"),
35                                   destination_frame = "retail_train",
36                                   header = TRUE,
37                                   sep = ";",
38                                   parse = TRUE)
39
40    retail_test <- h2o.importFile(path = file.path(script_path, "..", "data/retail_test.csv"),
41                                  destination_frame = "retail_test",
42                                  header = TRUE,
43                                  sep = ";",
44                                  parse = TRUE)
45
46    loginfo("--> Datasets imported into H2O cluster")
```

# terminal: Run the standard models

- In this case, we don't need to install dependencies or build the project because we didn't make any changes on the package, only on the project.
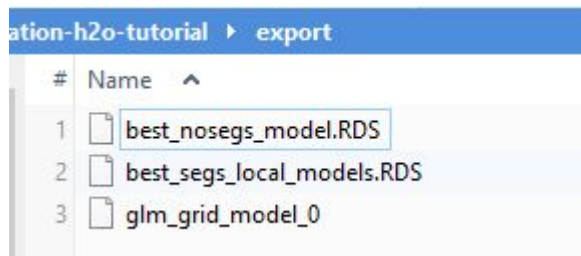
```
>Rscript R/build_p2b_segmentation_model.R
```

```
|============================================================| 100%
|============================================================| 100%
2019-05-04 13:56:36 INFO::--> Predicted segmentation models
|============================================================| 100%
|============================================================| 100%
|============================================================| 100%
|============================================================| 100%
2019-05-04 13:56:43 INFO::--> Built models with added segmentation assignment
2019-05-04 13:56:43 INFO::--> Best model [k = 2] with test AUC=0.647041773458421
|============================================================| 100%
```

# Observe an `export` folder has been created

- These are the outputs of the scripts

# Next steps

- Package
  - None
- Project
  - add 4th script `compare_models.R`
- Terminal
  - Run script `R/compare_models.R`

**project:** add the last script for model comparison



```
17    config <- load_config()
18    args <- args_parser()
19
20    #################################################################################
21
22    best_nosegs_model <- readRDS(file.path(script_path, "../export/best_nosegs_model.RDS"))
23    best_segsvar_model <- readRDS(file.path(script_path, "../export/best_segsvar_model.RDS"))
24    best_segs_local_models <- readRDS(file.path(script_path, "../export/best_segs_local_models.RDS"))
25
26
27    pROC::roc.test(roc1 = best_nosegs_model$roc,
28                   roc2 = best_segsvar_model$roc,
29                   alternative = "less")
30
31    pROC::roc.test(roc1 = best_nosegs_model$roc,
32                   roc2 = best_segs_local_models$roc,
33                   alternative = "less")
34
35    pROC::roc.test(roc1 = best_segsvar_model$roc,
36                   roc2 = best_segs_local_models$roc,
37                   alternative = "less")
38
```

# terminal: Run model comparison

```
R:\rsuite-projects\retail-segmentation-h2o-tutorial>Rscript R/compare_models.R

        DeLong's test for two correlated ROC curves

data:  best_nosegs_model$roc and best_segsvar_model$roc
Z = -3.195, p-value = 0.0006993
alternative hypothesis: true difference in AUC is less than 0
sample estimates:
AUC of roc1 AUC of roc2
  0.6460978   0.6470411


        DeLong's test for two ROC curves

data:  best_nosegs_model$roc and best_segs_local_models$roc
D = -1.8339, df = 179860, p-value = 0.03333
alternative hypothesis: true difference in AUC is less than 0
sample estimates:
AUC of roc1 AUC of roc2
  0.6460978   0.6512212


        DeLong's test for two ROC curves

data:  best_segsvar_model$roc and best_segs_local_models$roc
D = -1.4978, df = 179860, p-value = 0.06709
alternative hypothesis: true difference in AUC is less than 0
sample estimates:
AUC of roc1 AUC of roc2
  0.6470411   0.6512212
```