

# Deconvolution Analysis with CIBERSORT

Cristiane Esteves, Mariana Boroni - Bioinformatics and Computational Biology Lab (LBBC/INCA-RJ)

```
#Load libraries
library(c17r)/deconv_cibersort/deconv_cibersort/lib/"./how/manager/R/a86_64-pc-linux-gnu-library/4.2")
pkgs <- c("survival", "survminer", "data.table", "dplyr", "ggplot2", "r1071", "parallel", "preprocessCore", "core
plot", "RColorBrewer", "parallel", "gg dendro", "tidy" )
install.packages(pkgs)

#Load Packages
suppressPackageStartupMessages({
  library(tibble)
  library(dplyr)
  library(ggplot2)
  library(survival)
  library(survminer)
  library(r171)
  library(parallel)
  library(preprocessCore)
  library(data.table)
  library(corrplot)
  library(RColorBrewer)
  library(readr)
})

#Load script CIBERSORT and barplot function
#source("CIBERSORT.R")
#source("barplot_cibersort.R")

#Load signature matrix (LM22) and bulk RNA matrix (SKCM-Melanstasis)

LM22 is the signature genes file we used for Cibersort analyses (https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4730640/). The file contains
expression counts for 547 signature genes (547 rows) for 22 distinct human immune cells (22 columns).
```

```
lm22_signatures <- as.data.frame(fread("data/lm22.txt"))

## Warning in fread("data/lm22.txt"): Detected 22 column names but the data has
## 21 columns (i.e., invalid file). Added 1 extra default column name for the first
## column which is guessed to be row names or an index. Use setnames() afterwards
## If this guess is not correct, or fix the file write command that created the
## file to create a valid file.
```

```
lm22_signatures <- tibble::column_to_rownames(lm22_signatures, "VI")
print(head(lm22_signatures[,1:6]))
```

```
##           B cells Naive B cells memory Plasma cells T cells CD8
## ABCB4      555.71345      10.74423      7.225819      4.31128
## ABCB9      15.60394      22.08479      652.292228      24.22272
## ABCB1      215.30095      321.62102      38.618872      105.61388
## ABCB6      15.11795      16.64885      22.123737      13.42829
## ACPS       405.89738      195.20148      1120.184664      306.31252
## ADAM28     1943.74270      1148.12014      324.780800      22.68972
```

```
#Bulk TCGA-SKCM metastatic
skcm_bulk <- as.data.frame(fread("data/bulk.txt"))

## Warning in fread("data/bulk.txt"): Detected 366 column names but the data has
## 367 columns (i.e., invalid file). Added 1 extra default column name for the first
## column which is guessed to be row names or an index. Use setnames() afterwards
## If this guess is not correct, or fix the file write command that created the
## file to create a valid file.
```

```
skcm_bulk <- tibble::column_to_rownames(skcm_bulk, "VI")
print(head(skcm_bulk[,1:6]))
```

```
##           TCGA-BB-A5VV-06A-11B-A32P-07 TCGA-GH-A263-01A-11B-A187-07
## ABCB4      8.1243      1.8439
## ABCB9      1.1392      0.6741
## ABCB1      123.8629      0.6941
## ACB6      23.7555      0.0730
## ACPS       78.3980      76.1972
## ADAM28     43.2955      0.5015
## ADAM29     43.2955      0.5015
## ABCB4      1.8307      1.3882
## ABCB9      2.0190      2.5264
## ABCB1      6.1677      2.7528
## ACB6      1.1324      1.1902
## ACPS       46.2074      132.7469
## ADAM28     2.8939      0.7076
```

```
#Deconvolution Analysis - CIBERSORT

l.perm = No; permutations: set to >=100 to calculate p-values (default = 0)
i.QM = Quantile normalization of input mixture (default = TRUE) (downloading is recommended for RNA-Seq data)
iii.absolute = Run CIBERSORT in absolute mode (default = FALSE)
    note that cell subsets will be scaled by their absolute levels and will not be represented as fractions (to derive the default
    output: normalized absolute levels such that they sum to 1 for each mixture sample)
    the sum of all cell subsets in each mixture sample will be added to the output.
```

```
set.seed(42)
hl <- Sys.time()
results_cibersort <- CIBERSORT(lm22_signatures, skcm_bulk, perm = 1, absolute = F, QM = F)
hl2 <- Sys.time()
print(hl2 - hl)

## Time difference of 10.8416 mins
```

```
results_sign <- as.data.frame(results_cibersort)[which(as.data.frame(results_cibersort)$'P-value' <= 0.05),]
results_sign <- results_sign[1:102]
```

```
#Save CIBERSORT results to directory
saveRDS(results_sign, "data/cibersort/cancer_genome_analysis_africa/modules/RNA_deconvolution/Data_Deconvolut
ion/deconv_cibersort/results_cibersort.rds")
```

```
library(readr)

# Load Metastatic Melanoma (SKCM-TCGA) Clinical and Survival dataset
dedos_SKCM <- readRDS("data/Dedos_SKCM.rds")
subtypes <- read_csv("data/subtypes.csv")
```

```
## New names:
## Rows: 773 Columns: 11
## — Column specification
## ————— Delimiter: ",", chr
## (1): pan.sample.id, cancer.type, Subtype_MRNA, Subtype_MRNaseq, Subtype... db1
## (2): ..., Subtype_protein
## I use 'spec()' to retrieve the full column specification for this data. I
## Specify the column types or set 'show_col_types' = FALSE to quiet this message.
## • --> '...'

#Identify the quartile of each sample in each cell type
```

```
results_sign1 <- results_sign
for (i in 1:length(colnames(results_sign1))) {
  for (j in 1:10) {
    quant <- quantile(results_sign1[,i])
    results_sign1[which(results_sign1[,i] > quant[j]),i] <- j
  }
}
```

```
results_sign1Mixture <- rownames(results_sign1)
results_sign1Mixture <- paste("M", " ", results_sign1Mixture)
results_sign1Mixture <- substr(results_sign1Mixture,1,12)

#Aggregate CIBERSORT result with clinical and survival data according to patient ID
```

```
forest_data <- left_join(results_sign, dedos_SKCMsurvival_mat[,c(1,8,16,17,2,5)], by = c("Mixture" = "bcr_patient_
barcode"))
forest_data <- left_join(forest_data, subtypes[,c(1,10)], by = c("Mixture" = "pan.sample.id"))

#rename column
colnames(forest_data)[25] = "Subtype_Mutation"
```

## Barplot

Proportions of the expression predicted by Cibersort pf each celltype

```
# Check the columns (variable names) present in the clinical dataset
names(forest_data)
```

```
## [1] "B cells naive" "B cells memory"
## [2] "Plasma cells" "T cells CD8"
## [3] "T cells CD4 naive" "T cells CD4 memory resting"
## [4] "T cells CD4 memory activated" "T cells follicular helper"
## [5] "T cells regulatory (Tregs)" "T cells gamma delta"
## [6] "NK cells resting" "NK cells activated"
## [7] "Macrophages" "Macrophages_M0"
## [8] "Macrophages_M1" "Macrophages_M2"
## [9] "Dendritic cells resting" "Dendritic cells activated"
## [10] "Mast cells resting" "Mast cells activated"
## [11] "Eosinophils" "Neutrophils"
## [12] "Mixture" "Subtype_MRNaseq"
## [13] "OS" "OS_time"
## [14] "Gender" "age_at_initial_pathologic_diagnosis"
## [15] "Subtype_Mutation"
```

```
#Filter for columns sample and Subtype mutation
data_barplot <- forest_data[,c(3,25)]

#Make patient IDs unique values
data_barplotMixture <- make.names(data_barplotMixture, unique = T)
data_barplotMixture <- paste("M", " ", data_barplotMixture)
```

```
#One patient ID in columns
colnames(data_barplot) = data_barplotMixture
data_barplotMixture = NULL

# Add 'NA' in empty fields
data_barplotSubtype_Mutation[which(is.na(data_barplotSubtype_Mutation))] <- "nan"
data_barplotSubtype_Mutation[which(data_barplotSubtype_Mutation == "")] <- "nan"
```

```
# Make Mixture column (Patient ID) from cibersort result table as first column
res_cibersort <- forest_data[, c("Mixture", colnames(forest_data)[1:102])]
res_cibersortMixture <- make.names(res_cibersortMixture, unique = T)
res_cibersortMixture <- paste("M", " ", res_cibersortMixture)
```

```
#For the barplot in which each column is a patient with the clinical informations on the first row (stage, for a
sample) colored according to legend colors) and each bar
is divided by the proportion of immune cells types described also in the legend.
```

```
plot.ciber.beat(ciber.obj = res_cibersort, ann_info = data_barplot, sample.column = 1)
```

```
## Loading required package: gg dendro

## Loading required package: gridExtra
```

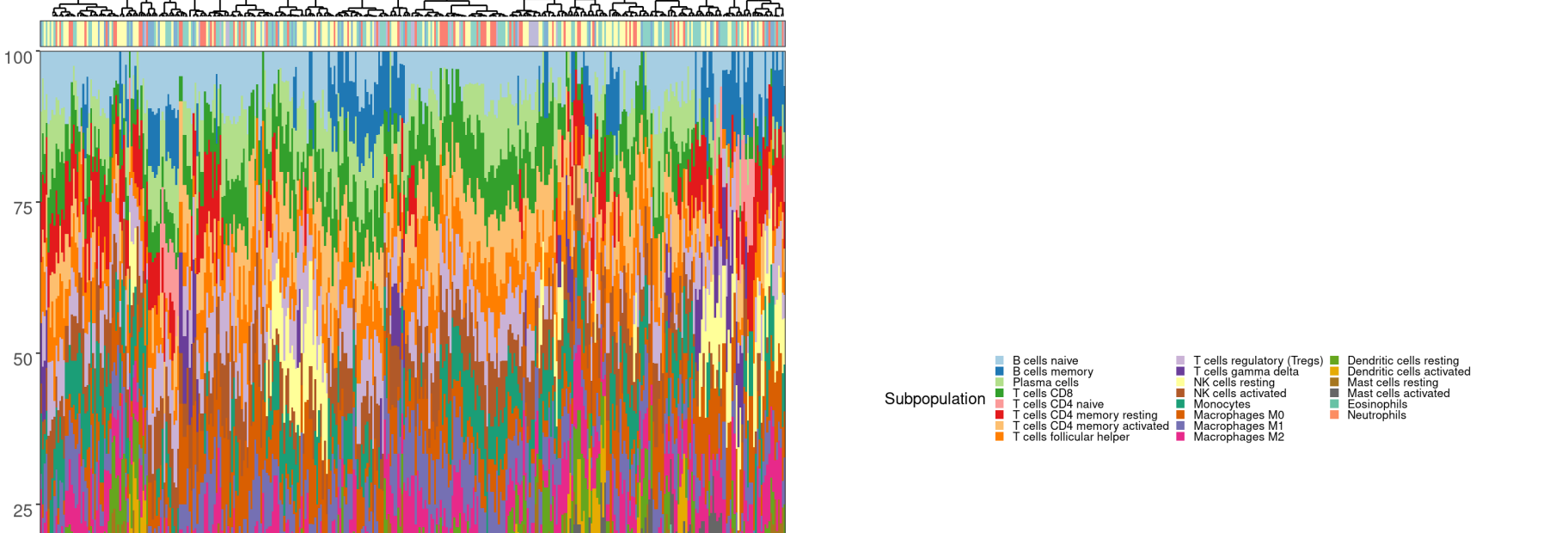
```
##
## Attaching package: 'gridExtra'
##
## The following object is masked from 'package:dplyr':
##
## combine
##
## Loading required package: grid
```

```
## Loading required package: cowplot

##
## Attaching package: 'cowplot'
```

```
## The following object is masked from 'package:ggpubr':
##
## get_legend

## Note: Using an external vector in selections is ambiguous.
## I use 'all_of(sample.column)' instead of 'sample.column' to silence this message.
## See <https://tidyselect.r-lib.org/reference/faq-external-vector.html>.
## This message is displayed once per session.
## Using Mixture as id variables
## Using Mixture as id variables
## Joining, by = "Mixture"
```



## Univariate and Multivariate (Cox Regression)/Survival analysis

The Cox proportional-hazards model (Cox, 1972) is essentially a regression model commonly used statistical in medical research for investigating the association between the survival time of patients and one or several quantitative variables (known as covariates), potentially affect patient prognosis.

In clinical investigations, there are many situations, where several known quantities (known as covariates), potentially affect patient prognosis. For instance, suppose two groups of patients are compared: those with and those without a specific genotype. If one of the groups also contains other individuals, any difference in survival may be attributable to genotype or age or indeed both. Hence, when investigating survival in relation to any one factor (i.e. the variable adjusted for), the effect of the other factors must be taken into account.

Statistical model is a frequently used tool that allows to analyze survival with respect to several factors simultaneously. Additionally, statistical model provides the effect size for each factor.

The Cox proportional-hazards model is one of the most important methods used for modelling survival analysis data. The next section introduces the basics of the Cox regression model.

References: <http://www.shdta.com/englist/wk/cox-proportional-hazards-model>  
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2394262/>

```
library(dplyr)
library(survival)
library(survminer)
```

```
# Univariate Cox
#Cox univariate analysis estimated the impact on survival of each cell type.
#For this, it is necessary to have the survival time of each patient and the event variable (in this case, death)
and provide this information as input to the survfit function
surv_object <- Surv(time = forest_data$OS.time, event = forest_data$OS)
```

```
colnames(forest_data)[1:122] <- grab(" ", " ", colnames(forest_data)[1:122])
colnames(forest_data)[9] <- "Tregs"
colnames(forest_data)
```

```
## [1] "B cells naive" "B cells memory"
## [2] "Plasma cells" "T cells CD8"
## [3] "T cells CD4 naive" "T cells CD4 memory resting"
## [4] "T cells CD4 memory activated" "T cells follicular helper"
## [5] "Tregs" "T cells gamma delta"
## [6] "NK cells resting" "NK cells activated"
## [7] "Macrophages" "Macrophages_M0"
## [8] "Macrophages_M1" "Macrophages_M2"
## [9] "Dendritic cells resting" "Dendritic cells activated"
## [10] "Mast cells resting" "Mast cells activated"
## [11] "Eosinophils" "Neutrophils"
## [12] "Mixture" "Subtype_MRNaseq"
## [13] "OS" "OS_time"
## [14] "Gender" "age_at_initial_pathologic_diagnosis"
## [15] "Subtype_Mutation"
```

```
# Get names of each column (immune cells)
covariables <- colnames(forest_data)[1:122/7:125]
univ_formulas <- apply(covariables, MARGIN=2, FUN=function(x) as.formula(paste("surv_object ~", x)))
univ_models <- lapply(univ_formulas, FUNCTION=function(x) coxph(x, data = forest_data))
univ_results <- apply(univ_models, MARGIN=2, FUN=function(x) {
  x <- summary(x)
  p.value <- signif(x$wald["p.value"], digits=2)
  wald.test <- signif(x$wald["test"], digits=2)
  beta <- signif(x$coef[1], digits=2)
  HR <- signif(x$coef[2], digits=2)
  HR.confint.lower <- signif(x$confint.int["lower .95%"], 2)
  HR.confint.upper <- signif(x$confint.int["upper .95%"], 2)
  HR <- paste(HR, " (",
    HR.confint.lower, " - ", HR.confint.upper, ")")
  res <- c(beta, HR, wald.test, p.value)
  names(res) <- c("beta", HR, "95% CI for HR", "wald.test", "p.value")
  return(res)
})
res_bisque <- as.data.frame(t(do.call(cbind, univ_results)))
```

```
## Warning in function (..., deparse.level = 1) : number of rows of result is not
## a multiple of vector length (arg 1)
```

```
res_bisque <- as.data.frame(res_bisque)
res_bisque$p.value <- as.character(res_bisque$p.value)
res_bisque$p.value <- as.numeric(res_bisque$p.value)
```

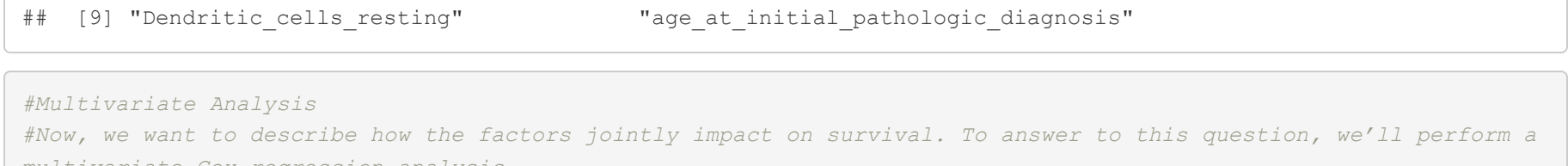
```
## Warning: NA's introduced by coercion

#Filter for pval <= 0.05
res_bisque_filt <- res_bisque[which(res_bisque$p.value <= 0.05),]
res_bisque_filt
```

```
#Check the immune cells that significantly impact each patient's survival (p val <= 0.05)
rownames(res_bisque_filt)
```

```
## [1] "Plasma cells" "T cells CD8"
## [2] "T cells CD4 memory resting" "T cells CD4 memory activated"
## [3] "NK cells activated" "Macrophages_M0"
## [4] "Macrophages_M1" "Macrophages_M2"
## [5] "Dendritic cells resting" "age_at_initial_pathologic_diagnosis"
```

```
#Multivariate Analysis
#Now, we want to describe how the factors jointly impact on survival. To answer to this question, we'll perform a
multivariate Cox regression analysis.
El <- as.formula(paste("Surv(forest_data$OS.time, event = forest_data$OS) ~",
  paste(c(rownames(res_bisque_filt)), collapse=" + ")
))
fit_coxph <- coxph(El, data = forest_data)
summary(fit_coxph)
ggforest(fit_coxph, data = forest_data, main = "Hazard Ratio Melanoma Metastasis")
```



The Cox regression results can be interpreted as follow:

Statistical significance: The column marked "Z" gives the Wald statistic value. It corresponds to the ratio of each regression coefficient to its standard error (Z = coefficient/coef). The Wald statistic evaluates, whether the beta (β) coefficient of a given variable is statistically significantly different from 0. From the output above, we can conclude that the variable we have highly statistically significant coefficients.

The regression coefficients: The second feature to note in the Cox model results is the age of the regression coefficients (coef). A positive sign means that the hazard ratio (hazard) is higher, and thus the prognosis worse, for subjects with higher values of that variable.

Hazard ratios: The exponentiated coefficients (exp(coef) = exp(0.53) = 0.59). Also known as hazard ratios, give the effect size of covariates. For example, being TCGA memory activated reduces the hazard by a factor of 0.78. Being TCGA memory activated is associated with good prognosis.

Confidence intervals of the hazard ratios: The summary output also gives upper and lower 95% confidence intervals for the hazard ratio (exp(coef)).

Global statistical significance of the model: Finally, the output gives p-values for three alternative tests for overall significance of the model: The likelihood-ratio test, Wald test, and score logrank statistics. These three methods are asymptotically equivalent. For large enough N, they will give similar results. For small N, they may differ somewhat. The Likelihood ratio test has better behavior for small sample sizes, so it is generally preferred.

```
##### Levels expressions NK macrophages and T cells CD4 memory activated in survival Analysis (Kaplan-Meier plot)
library(survminer)
library(survfit)
library(survminer)
forest_data$Macrophages_M1_group = ifelse(forest_data$Macrophages_M1 == mean(forest_data$Macrophages_M1), "High", "Low")

#Having fit a Cox model to the data, it's possible to visualize the predicted survival proportion at any given point in time for a particular risk group. The function survfit() estimates the survival proportion, by default at the mean values of covariates.
fit <- survfit(surv(OS.time, OS) ~ Macrophages_M1_group, data = forest_data)
ggsurvplot(fit, palette = c("087093", "#20B2AA"), xlab = "Survival time in years",
  surv.median.line = c("hv"), cunecolor = F, conf.int = F, risk.table = TRUE, pval = T,
  title = "Overall survival: TCGA-SKCM (Macrophages M1), risk.table.y.test.col = T, # colour risk table text annotations.
  risk.table.y.test = FALSE, font.main = c(10), font.legend = c(10), font.y = c(10), font.x = c(10), font.cunecolor = c(10),
  font.ticks.lab = c(10), legend.labs = c("Macrophages M1 High", "Macrophages M1 Low"), fontsize = 3, risk.table.height = 0.1, pval.size = 4, censor.size = 3,
  font.yticks.lab = c(10))

Overall survival: TCGA-SKCM (Macrophages M1)

Strata — Macrophages M1 High — Macrophages M1 Low

Survival time in years

Number at risk

Strata
— 155 115 75 50 34 24
— 145 87 54 34 22 16
0 1 2 3 4 5
```

```
forest_data$T_cells_CD4_memory_activated_group = ifelse(forest_data$T_cells_CD4_memory_activated == mean(forest_data$T_cells_CD4_memory_activated), "High", "Low")

#Having fit a Cox model to the data, it's possible to visualize the predicted survival proportion at any given point in time for a particular risk group. The function survfit() estimates the survival proportion, by default at the mean values of covariates.
fit <- survfit(surv(OS.time, OS) ~ T_cells_CD4_memory_activated_group, data = forest_data)
ggsurvplot(fit, palette = c("087093", "#20B2AA"), xlab = "Survival time in years",
  surv.median.line = c("hv"), cunecolor = F, conf.int = F, risk.table = TRUE, pval = T,
  title = "Overall survival: TCGA-SKCM (Macrophages M1), risk.table.y.test.col = T, # colour risk table text annotations.
  risk.table.y.test = FALSE, font.main = c(10), font.legend = c(10), font.y = c(10), font.x = c(10), font.cunecolor = c(10),
  font.ticks.lab = c(10), legend.labs = c("T cells CD4 memory activated High", "T cells CD4 memory activated Low"), fontsize = 3, risk.table.height = 0.1, pval.size = 4, censor.size = 3,
  font.yticks.lab = c(10))

Overall survival: TCGA-SKCM (T cells CD4 memory activated)

Strata — T cells CD4 memory activated High — T cells CD4 memory activated Low

Survival time in years

Number at risk

Strata
— 149 113 76 53 34 23
— 145 87 54 34 22 16
0 1 2 3 4 5
```