

Introduction to CytoDx

Zicheng Hu

2025-10-30

Contents

Introduction	2
Installation	2
Example: diagnosing AML using flow cytometry	3
Step 1: Prepare data	3
Step 2: Build CytoDx model	4
Step 3: Predict AML using testing data	4
Step 4: Find cell subsets associated with AML	5
Session Infomation	6

Introduction

CytoDx is a method that predicts clinical outcomes using single cell data without the need of cell gating. It first predicts the association between each cell and the outcome using a linear statistical model (Figure 1). The cell level predictions are then averaged within each sample to represent the sample level predictor. A second model is used to make prediction at the sample level (Figure 1). Compare to traditional gating based methods, CytoDX have multiple advantages.

1. Robustness. CytoDx is able to robustly predict clinical outcomes using data acquired by different cytometry platforms in different research institutes.
2. Interpretability. CytoDx identifies cell markers and cell subsets that are most associated with the clinical outcome, allowing researchers to interpret the result easily.
3. Simplicity. CytoDx associates cytometry data with clinical features without any cell gating steps, therefore simplifying the diagnostic process in clinical settings.

In section 2, we demonstrate how to install CytoDx.

In section 3, we show an example where CytoDx is used to diagnose acute myeloid leukemia (AML). In addition to perform diagnosis, we also show that CytoDx can be used to identify the cell subsets that are associated with AML.

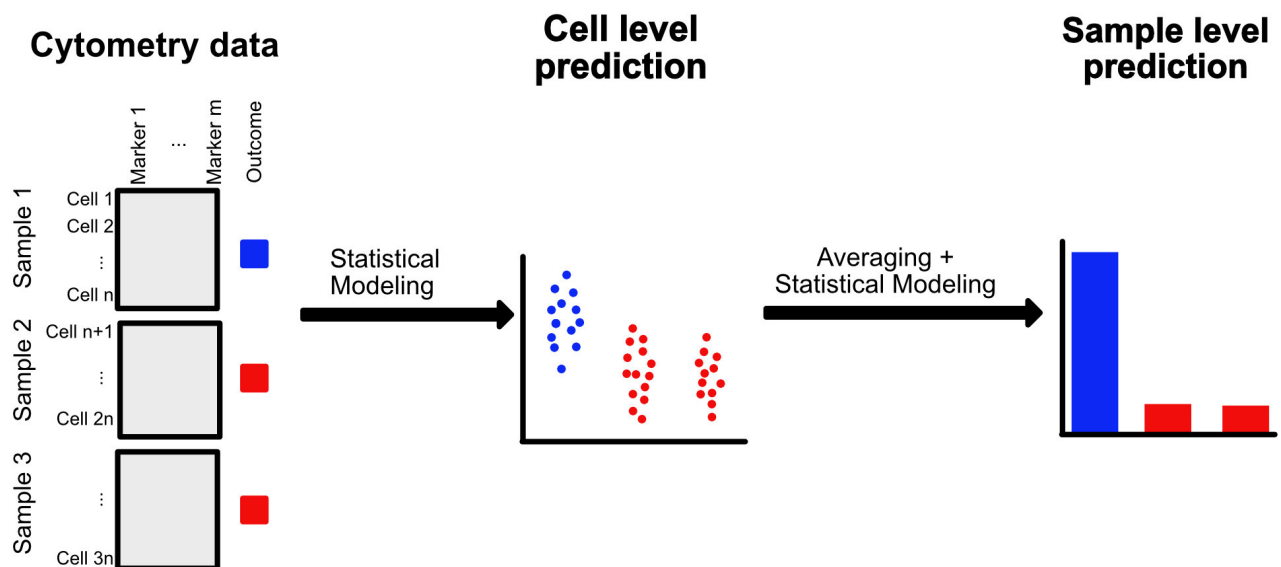


Figure 1

Installation

You can install the stable version of CytoDx from Bioconductor:

```
BiocManager::install("CytoDx")
```

You can also install the latest develop version of CytoDx from github:

```
devtools::install_github("hzc363/CytoDx")
```

Example: diagnosing AML using flow cytometry

In this example, we build a CytoDx model to diagnose acute myeloid leukemia (AML) using flow cytometry data. We train the model using data from 5 AML patients and 5 controls and test the performance in a test dataset.

Step 1: Prepare data

The CytoDx R package contains the fcs files and the ground truth (AML or normal) that are needed for our example. We first load the ground truth.

```
library(CytoDx)

# Find data in CytoDx package
path <- system.file("extdata",package="CytoDx")

# read the ground truth
fcs_info <- read.csv(file.path(path,"fcs_info.csv"))

# print out the ground truth
knitr::kable(fcs_info)
```

fcsName	Label	dataset
sample1.fcs	aml	test
sample2.fcs	aml	test
sample3.fcs	aml	test
sample4.fcs	aml	test
sample5.fcs	aml	test
sample6.fcs	normal	test
sample7.fcs	normal	test
sample8.fcs	normal	test
sample9.fcs	normal	test
sample10.fcs	normal	test
sample11.fcs	aml	train
sample12.fcs	aml	train
sample13.fcs	aml	train
sample14.fcs	aml	train
sample15.fcs	aml	train
sample16.fcs	normal	train
sample17.fcs	normal	train
sample18.fcs	normal	train
sample19.fcs	normal	train
sample20.fcs	normal	train

We then read the cytometry data for training samples using the fcs2DF function.

```
# Find the training data
train_info <- subset(fcs_info,fcs_info$dataset=="train")

# Specify the path to the cytometry files
fn <- file.path(path,train_info$fcsName)

# Read cytometry files using fcs2DF function
train_data <- fcs2DF(fcsFiles=fn,
```

```

y=train_info$Label,
assay="FCM",
b=1/150,
excludeTransformParameters=
  c("FSC-A","FSC-W","FSC-H","Time"))

```

The CytoDx is flexible to data transformations. It can be applied to rank transformed data to reduce batch effects. Here, we transform the original data to rank data.

```

# Performs rank transformation
x_train <- pRank(x=train_data[,1:7],xSample=train_data$xSample)

# Convert data frame into matrix. Here we included the 2-way interactions.
x_train <- model.matrix(~.*,x_train)

```

Step 2: Build CytoDx model

We use training data to build a predictive model.

```

# Build predictive model using the CytoDx.fit function
fit <- CytoDx.fit(x=x_train,
  y=(train_data$y=="aml"),
  xSample=train_data$xSample,
  family = "binomial",
  reg = FALSE)

```

Step 3: Predict AML using testing data

We first load and rank transform the test data.

```

# Find testing data
test_info <- subset(fcs_info,fcs_info$dataset=="test")

# Specify the path to cytometry files
fn <- file.path(path,test_info$fcsName)

# Read cytometry files using fcs2DF function
test_data <- fcs2DF(fcsFiles=fn,
  y=NULL,
  assay="FCM",
  b=1/150,
  excludeTransformParameters=
    c("FSC-A","FSC-W","FSC-H","Time"))

# Performs rank transformation
x_test <- pRank(x=test_data[,1:7],xSample=test_data$xSample)

# Convert data frame into matrix. Here we included the 2-way interactions.
x_test <- model.matrix(~.*,x_test)

```

We use the built CytoDx model to predict AML.

```

# Predict AML using CytoDx.ped function
pred <- CytoDx.pred(fit,xNew=x_test,xSampleNew=test_data$xSample)

```

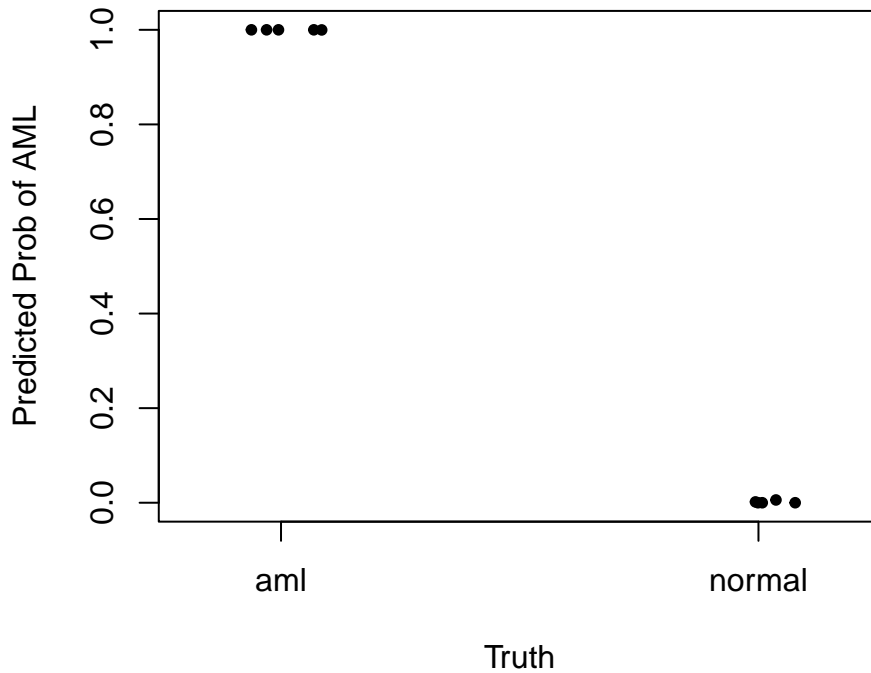
We plot the prediction. In this example, CytoDx classifies the sample into AML and normal perfectly.

```

# Combine prediction and truth
result <- data.frame("Truth"=test_info$Label,
                     "Prob"=pred$xNew.Pred.sample$y.Pred.s0)

# Plot the prediction
stripchart(result$Prob~result$Truth, jitter = 0.1,
           vertical = TRUE, method = "jitter", pch = 20,
           xlab="Truth",ylab="Predicted Prob of AML")

```



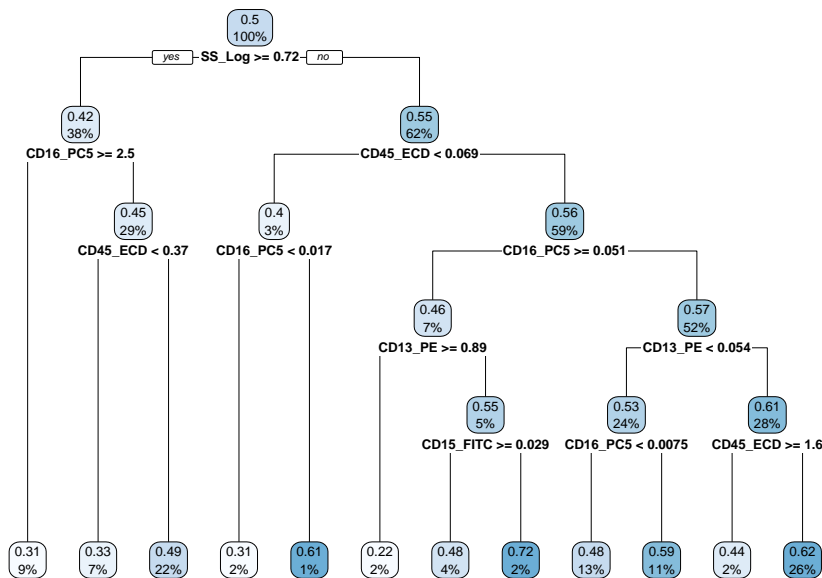
Step 4: Find cell subsets associated with AML

We use a decision tree to find cell subsets that are associated the AML. In this step, the original cytometry data should be used, rather than the ranked data.

```

# Use decision tree to find the cell subsets that are associated the AML.
TG <- treeGate(P = fit$train.Data.cell$y.Pred.s0,
               x= train_data[,1:7])

```



Session Information

`sessionInfo()`

R Under development (unstable) (2025-10-20 r88955) Platform: x86_64-pc-linux-gnu Running under: Ubuntu 24.04.3 LTS

Matrix products: default BLAS: /home/biocbuild/bbs-3.23-bioc/R/lib/libRblas.so LAPACK: /usr/lib/x86_64-linux-gnu/lapack/liblapack.so.3.12.0 LAPACK version 3.12.0

locale: [1] LC_CTYPE=en_US.UTF-8 LC_NUMERIC=C
[3] LC_TIME=en_GB LC_COLLATE=C
[5] LC_MONETARY=en_US.UTF-8 LC_MESSAGES=en_US.UTF-8
[7] LC_PAPER=en_US.UTF-8 LC_NAME=C
[9] LC_ADDRESS=C LC_TELEPHONE=C
[11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C

time zone: America/New_York tzcode source: system (glibc)

attached base packages: [1] stats graphics grDevices utils datasets methods base

other attached packages: [1] CytoDx_1.31.0

loaded via a namespace (and not attached): [1] Matrix_1.7-4 glmnet_4.1-10 dplyr_1.1.4

[4] compiler_4.6.0 tinytex_0.57 rpart_4.1.24
[7] tidyselect_1.2.1 Rcpp_1.1.0 Biobase_2.71.0
[10] cytolib_2.23.0 parallel_4.6.0 splines_4.6.0
[13] yaml_2.3.10 fastmap_1.2.0 lattice_0.22-7
[16] R6_2.6.1 RProtoBufLib_2.23.0 generics_0.1.4
[19] rpart.plot_3.1.3 shape_1.4.6.1 knitr_1.50
[22] BiocGenerics_0.57.0 iterators_1.0.14 tibble_3.3.0
[25] pillar_1.11.1 rlang_1.1.6 flowCore_2.23.0
[28] xfun_0.54 doParallel_1.0.17 cli_3.6.5
[31] withr_3.0.2 magrittr_2.0.4 digest_0.6.37
[34] foreach_1.5.2 grid_4.6.0 lifecycle_1.0.4
[37] S4Vectors_0.49.0 vctrs_0.6.5 evaluate_1.0.5
[40] glue_1.8.0 codetools_0.2-20 survival_3.8-3

[43] stats4_4.6.0 rmarkdown_2.30 matrixStats_1.5.0
[46] tools_4.6.0 pkgconfig_2.0.3 htmltools_0.5.8.1