

Alzheimer EDA

Paul K. Yu

2024-11-21

Exploratory data analysis

Do this once (set eval to TRUE only once)

```
# Install required libraries if not already installed
if (!requireNamespace("pheatmap", quietly = TRUE)) install.packages("pheatmap")
if (!requireNamespace("ggplot2", quietly = TRUE)) install.packages("ggplot2")
```

Get data

```
library(data.table)

# Get RNA data
expression <- as.data.frame(fread("alzheimer/GSE125583_raw_counts_GRCh38.p13_NCBI.tsv", header = TRUE, ))

# Remove NA
expression <- na.omit(expression)

# Remove duplicate genes
expression <- expression[!duplicated(expression$GeneID), ]

# Make gene ID as row name
rownames(expression) <- expression[[1]] # Set the first column as row names

# Remove gene ID column
expression <- expression[, -c(1)]

n <- ncol(expression) ## number of subjects
g <- nrow(expression) ## number of gene features

# Get outcome data
# Manually select only ID_REF and stage
outcome <- as.data.frame(t(fread("alzheimer/GSE125583_series_matrix.txt", header = FALSE, sep = "\t")))
outcome <- setNames(outcome, c("ID_REF", "stage"))
outcome <- outcome[-c(1), ]
outcome$stage <- as.factor(outcome$stage)
levels(outcome$stage)

## [1] "braak.score: I" "braak.score: II" "braak.score: III" "braak.score: IV"
## [5] "braak.score: NA" "braak.score: V" "braak.score: VI"
```

```

levels(outcome$stage) <- c(1,2,3,4,0,5,6)

# Make ID_REF as row name
rownames(outcome) <- outcome[[1]] # Set the first column as row names

# Remove NA
outcome <- na.omit(outcome)
matches <- colnames(eexpression) %in% rownames(outcome)
eexpression <- eexpression[, matches]
matches <- rownames(outcome) %in% colnames(eexpression)
outcome <- outcome[matches,]

```

Summary of the dataset

Get summary of outcome data

```
summary(outcome)
```

```

##      ID_REF      stage
## Length:289      1: 4
## Class :character 2:15
## Mode  :character 3:44
##                               4:99
##                               0: 2
##                               5:65
##                               6:60
##

```

Check for missing values

```
sum(is.na(eexpression))
```

```
## [1] 0
```

```
sum(is.na(outcome))
```

```
## [1] 0
```

Correlation analysis

```
library(pheatmap)
```

```
# Compute correlations
```

```
cor_matrix <- cor(eexpression, use = "pairwise.complete.obs")
```

```
# Visualize correlations as a heatmap
```

```

pheatmap(cor_matrix,
          show_rownames = FALSE,
          show_colnames = FALSE,
          clustering_distance_rows = "correlation",
          clustering_distance_cols = "correlation",
          fontsize = 4)

```

