

Untitled

Anishkumar Senthil

2025-11-05

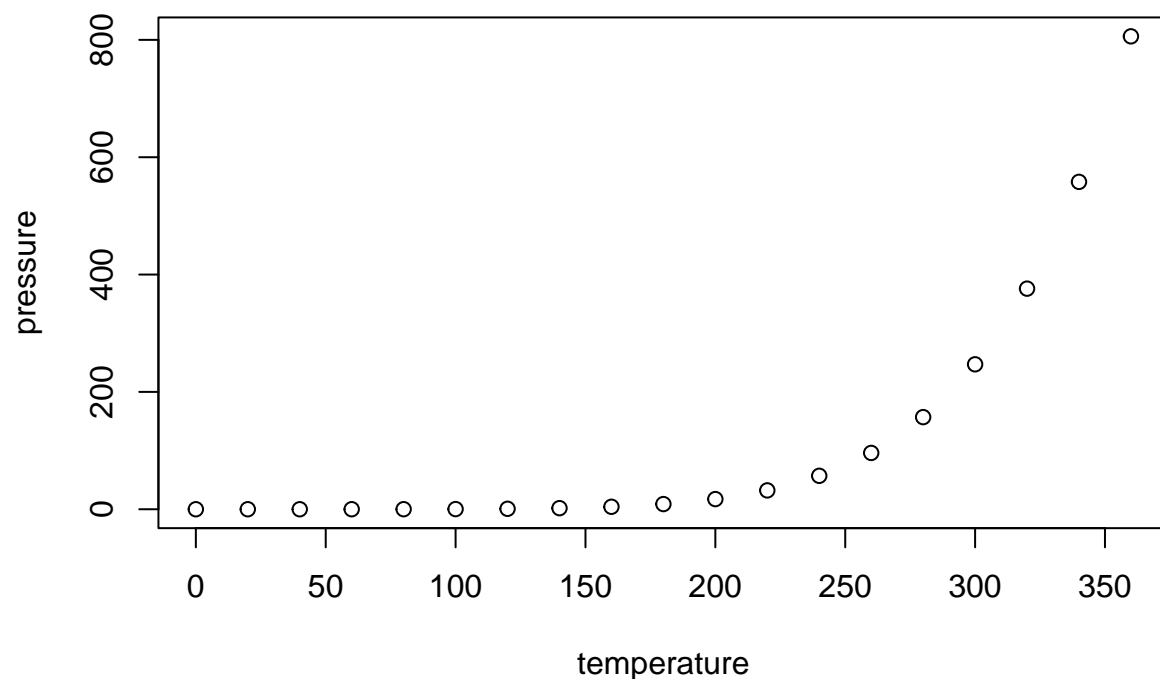
R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
summary(cars)
```

```
##      speed      dist
##  Min.   : 4.0    Min.   :  2.00
##  1st Qu.:12.0    1st Qu.: 26.00
##  Median :15.0    Median : 36.00
##  Mean   :15.4    Mean   : 42.98
##  3rd Qu.:19.0    3rd Qu.: 56.00
##  Max.   :25.0    Max.   :120.00
```



Select top predictive proteins

```
top_n <- 20

load('../data/biomarker-clean.RData')

library(dplyr)

dat <- biomarker_clean
dat$group <- factor(dat$group)

proteins <- setdiff(names(dat), c('group','ados'))

tt_res <- sapply(proteins, function(p){
  x <- dat[[p]]
  grp <- dat$group
  ok <- !is.na(x) & !is.na(grp)
  if(sum(ok) < 3) return(NA)
  t <- try(t.test(x[ok] ~ grp[ok]), silent=TRUE)
  if(inherits(t,'try-error')) return(NA)
  t$p.value
})
tt_df <- tibble::tibble(protein = proteins,
                        pvalue = as.numeric(tt_res)) %>%
```

```

    arrange(pvalue) %>%
    slice_head(n = top_n)

if(!requireNamespace('randomForest', quietly = TRUE)) {
  stop("Package 'randomForest' required. Please install it before running this chunk.")
}
rf_dat <- dat %>% dplyr::select(dplyr::all_of(proteins))
rf_resp <- dat$group
rf_fit <- randomForest::randomForest(x = rf_dat, y = rf_resp, ntree = 1000, importance = TRUE)
rf_imp_mat <- randomForest::importance(rf_fit, type = 2)
rf_imp_val <- if(is.matrix(rf_imp_mat)) rf_imp_mat[,1] else rf_imp_mat
rf_df <- tibble::tibble(protein = names(rf_imp_val), importance = as.numeric(rf_imp_val)) %>%
  arrange(desc(importance)) %>%
  slice_head(n = top_n)

if(!requireNamespace('glmnet', quietly = TRUE)) {
  stop("Package 'glmnet' required. Please install it before running this chunk.")
}
X <- as.matrix(rf_dat)
y <- as.numeric(dat$group) - 1
cv <- glmnet::cv.glmnet(X, y, family = 'binomial', alpha = 1, nfolds = 5)
# Get coefficients using the generic coef() function
coef_min <- as.matrix(coef(cv, s = 'lambda.min'))
# Remove intercept (first row) and get coefficients
coefs <- coef_min[-1,1]
lasso_df <- tibble::tibble(protein = proteins, coef = as.numeric(coefs)) %>%
  mutate(abscoef = abs(coef)) %>%
  arrange(desc(abscoef)) %>%
  slice_head(n = top_n)

tt_top <- tt_df$protein
rf_top <- rf_df$protein
lasso_top <- lasso_df$protein

cat(glue::glue("Selected top {top_n} proteins by each method:\n"))

```

Selected top 20 proteins by each method:

```
cat("- t-test (by p-value):\n")
```

- t-test (by p-value):

```
print(tt_top)
```

```
## [1] "DERM"          "RELT"          "FSTL1"         "C1QR1"
## [5] "Calcineurin"   "CXCL16, soluble" "IgD"          "MRC2"
## [9] "PTN"          "Cadherin-5"    "MAPK2"        "TGF- $\beta$  R III"
## [13] "DAF"          "MIA"           "Notch 1"      "gp130, soluble"
## [17] "MMP-2"        "ALCAM"         "ROR1"         "MATN2"
```

```
cat('\n- Random Forest (importance):\n')
```

```
##  
## - Random Forest (importance):
```

```
print(rf_top)
```

```
## [1] "DERM"          "MAPK14"        "IgD"           "TSP4"          "RELT"  
## [6] "FSTL1"         "Notch 1"       "TGF-b R III"  "ERBB1"         "eIF-4H"  
## [11] "PTN"           "MAPK2"         "M2-PK"         "MMP-2"         "TrkC"  
## [16] "CSK"           "SRCN1"         "Cadherin-5"   "SOST"          "IGFBP-1"
```

```
cat('\n- LASSO (coef magnitude):\n')
```

```
##  
## - LASSO (coef magnitude):
```

```
print(lasso_top)
```

```
## [1] "IgD"                "DERM"  
## [3] "14-3-3 protein zeta/delta" "Epo"  
## [5] "MAPK2"              "ENTP5"  
## [7] "Protein S"          "FSTL1"  
## [9] "IL-17 RC"           "SRCN1"  
## [11] "CD59"               "IL-6 sRa"  
## [13] "ITI heavy chain H4"  "TWEAKR"  
## [15] "PAI-1"              "PYY"  
## [17] "CSR3P3"             "FAM3D"  
## [19] "hnRNP K"            "HGFA"
```

```
unique_proteins <- unique(c(tt_top, rf_top, lasso_top))  
cat('\nTotal unique proteins across methods: ', length(unique_proteins), '\n')
```

```
##  
## Total unique proteins across methods: 45
```

```
selected_lists <- list(tt = tt_df, rf = rf_df, lasso = lasso_df)
```

```
selected_lists
```

```
## $tt  
## # A tibble: 20 x 2  
##   protein          pvalue  
##   <chr>          <dbl>  
## 1 DERM          0.00000000827  
## 2 RELT          0.00000000782  
## 3 FSTL1         0.0000000466  
## 4 C1QR1         0.0000000479  
## 5 Calcineurin   0.0000000537
```

```

## 6 CXCL16, soluble 0.000000875
## 7 IgD 0.000000933
## 8 MRC2 0.00000103
## 9 PTN 0.00000135
## 10 Cadherin-5 0.00000175
## 11 MAPK2 0.00000204
## 12 TGF-b R III 0.00000330
## 13 DAF 0.00000397
## 14 MIA 0.00000483
## 15 Notch 1 0.00000500
## 16 gp130, soluble 0.00000530
## 17 MMP-2 0.00000552
## 18 ALCAM 0.00000664
## 19 ROR1 0.00000786
## 20 MATN2 0.00000799
##
## $rf
## # A tibble: 20 x 2
##   protein importance
##   <chr> <dbl>
## 1 DERM 0.991
## 2 MAPK14 0.656
## 3 IgD 0.613
## 4 TSP4 0.519
## 5 RELT 0.478
## 6 FSTL1 0.461
## 7 Notch 1 0.427
## 8 TGF-b R III 0.411
## 9 ERBB1 0.375
## 10 eIF-4H 0.358
## 11 PTN 0.355
## 12 MAPK2 0.343
## 13 M2-PK 0.336
## 14 MMP-2 0.322
## 15 TrkC 0.312
## 16 CSK 0.285
## 17 SRCN1 0.281
## 18 Cadherin-5 0.275
## 19 SOST 0.273
## 20 IGFBP-1 0.267
##
## $lasso
## # A tibble: 20 x 3
##   protein coef abscoef
##   <chr> <dbl> <dbl>
## 1 IgD 0.563 0.563
## 2 DERM 0.525 0.525
## 3 14-3-3 protein zeta/delta 0.319 0.319
## 4 Epo 0.276 0.276
## 5 MAPK2 0.272 0.272
## 6 ENTP5 -0.266 0.266
## 7 Protein S 0.231 0.231
## 8 FSTL1 0.229 0.229
## 9 IL-17 RC -0.229 0.229

```

| | | |
|--------------------------|--------|-------|
| ## 10 SRCN1 | 0.216 | 0.216 |
| ## 11 CD59 | -0.214 | 0.214 |
| ## 12 IL-6 sRa | -0.165 | 0.165 |
| ## 13 ITI heavy chain H4 | 0.154 | 0.154 |
| ## 14 TWEAKR | -0.147 | 0.147 |
| ## 15 PAI-1 | -0.145 | 0.145 |
| ## 16 PYY | 0.140 | 0.140 |
| ## 17 CSRP3 | -0.139 | 0.139 |
| ## 18 FAM3D | -0.124 | 0.124 |
| ## 19 hnRNP K | 0.119 | 0.119 |
| ## 20 HGFA | 0.111 | 0.111 |