

Pertussis Challenge 1.2

Runqi Zhang

2024-11-21

```
# Load necessary libraries
suppressPackageStartupMessages({
  library(dplyr)      # Data manipulation
  library(tidyr)      # Data reshaping
  library(readr)      # Reading TSV files
  library(glmnet)     # Regularized regression (LASSO, Ridge)
  library(agua)       # H2O AutoML integration
  library(tibble)
  library(impute)
})
```

```
## Warning: package 'tidyr' was built under R version 4.3.3
```

```
## Warning: package 'readr' was built under R version 4.3.3
```

```
## Warning: package 'glmnet' was built under R version 4.3.3
```

```
## Warning: package 'agua' was built under R version 4.3.3
```

```
## Warning: package 'parsnip' was built under R version 4.3.3
```

```
# Set working directory and initialize H2O
workDir <- "C:/Users/zhang/Desktop/cmi-pb-3rd-public-challenge-data-prep/Runqi/CMI-PB"
options(readr.show_col_types = FALSE)
agua::h2o_start()
```

```
## Warning: JAVA not found, H2O may take minutes trying to connect.
```

```
## Warning in h2o.clusterInfo():
```

```
## Your H2O cluster version is (11 months) old. There may be a newer version available.
```

```
## Please download and install the latest version from: https://h2o-release.s3.amazonaws.com/h2o/latest.
```

```
# Function to read data files for a given year
read_data <- function(year, type = "LD") {
  list(
    pts = read_tsv(file.path(workDir, paste0("data/", year, type, "_subject.tsv"))),
    sample = read_tsv(file.path(workDir, paste0("data/", year, type, "_specimen.tsv"))),
    ab = read_tsv(file.path(workDir, paste0("data/", year, type, "_plasma_ab_titer.tsv")))
  )
}
```

```

}

# Read datasets for 2020, 2021, 2022, and 2023
data2020 <- read_data("2020")
data2021 <- read_data("2021")
data2022 <- read_data("2022")
data2023 <- read_data("2023", type = "BD") # Note: Different file naming for 2023

```

1 Challenge1.2

```

# Calculate Fold Change (FC) for IgG-PT from baseline (Day 0) to Day 14
yDF <- data2020$ab %>%
  filter(isotype == "IgG", antigen == "PT") %>%
  inner_join(data2020$sample, by = "specimen_id") %>%
  inner_join(data2020$pts, by = "subject_id") %>%
  filter(planned_day_relative_to_boost %in% c(0, 14)) %>%
  dplyr::select(subject_id, planned_day_relative_to_boost, MFI_normalised) %>%
  pivot_wider(names_from = planned_day_relative_to_boost, values_from = MFI_normalised, names_prefix =
  mutate(MFI_FC = MFI_14 / MFI_0) # Compute Fold Change (FC)

```

```

# Helper function to prepare predictors
prepare_predictors <- function(ab_data, sample_data, pts_data) {
  ab_data %>%
    inner_join(sample_data, by = "specimen_id") %>%
    inner_join(pts_data, by = "subject_id") %>%
    filter(planned_day_relative_to_boost == 0, grepl("IgG", isotype)) %>%
    mutate(cname = make.names(paste0(isotype, "_", antigen))) %>%
    dplyr::select(subject_id, cname, MFI_normalised) %>%
    distinct() %>%
    pivot_wider(names_from = cname, values_from = MFI_normalised) %>%
    column_to_rownames(var = "subject_id")
}

```

```

# Prepare predictors for each year
xDF <- prepare_predictors(data2020$ab, data2020$sample, data2020$pts)
x2DF <- prepare_predictors(data2021$ab, data2021$sample, data2021$pts)
x3DF <- prepare_predictors(data2022$ab, data2022$sample, data2022$pts)
x4DF <- prepare_predictors(data2023$ab, data2023$sample, data2023$pts)

```

```

# Align column names across datasets
common_cols <- Reduce(intersect, list(colnames(xDF), colnames(x2DF), colnames(x3DF), colnames(x4DF)))
xDF <- xDF[, common_cols]
x2DF <- x2DF[, common_cols]
x3DF <- x3DF[, common_cols]
x4DF <- x4DF[, common_cols]

```

```

# Scale each dataset
xDF <- scale(xDF)
x2DF <- scale(x2DF)
x3DF <- scale(x3DF)
x4DF <- scale(x4DF)

```

2 2022 test set

```
# Prepare response variable (MFI_FC) for 2020, 2021, 2022
yobs1 <- data2021$ab %>%
  filter(isotype == "IgG", antigen == "PT") %>%
  inner_join(data2021$sample, by = "specimen_id") %>%
  inner_join(data2021$pts, by = "subject_id") %>%
  filter(planned_day_relative_to_boost %in% c(0, 14)) %>%
  dplyr::select(subject_id, planned_day_relative_to_boost, MFI_normalised) %>%
  pivot_wider(names_from = planned_day_relative_to_boost, values_from = MFI_normalised, names_prefix =
  mutate(MFI_FC = MFI_14 / MFI_0) %>%
  slice(match(rownames(x2DF), subject_id))

yobs2 <- data2022$ab %>%
  filter(isotype == "IgG", antigen == "PT") %>%
  inner_join(data2022$sample, by = "specimen_id") %>%
  inner_join(data2022$pts, by = "subject_id") %>%
  filter(planned_day_relative_to_boost %in% c(0, 14)) %>%
  dplyr::select(subject_id, planned_day_relative_to_boost, MFI_normalised) %>%
  pivot_wider(names_from = planned_day_relative_to_boost, values_from = MFI_normalised, names_prefix =
  mutate(MFI_FC = MFI_14 / MFI_0) %>%
  slice(match(rownames(x3DF), subject_id))

# Combine predictors and response
trainDF <- rbind(xDF, x2DF, x3DF) %>%
  as.data.frame() %>%
  mutate(MFI_FC = c(yDF$MFI_FC, yobs1$MFI_FC, yobs2$MFI_FC))

# Apply k-NN imputation
train_matrix <- as.matrix(trainDF) # Convert to matrix
imputed_matrix <- impute.knn(train_matrix)$data
trainDF <- as.data.frame(imputed_matrix) # Convert back to dataframe

set.seed(3)
auto_fit <- auto_ml() %>%
  set_engine("h2o", max_runtime_secs = 5) %>%
  set_mode("regression") %>%
  fit(MFI_FC ~ ., data = trainDF)

# Predict on training data
train_predictions <- predict(auto_fit, new_data = trainDF)$pred

# Calculate correlations
pearson_cor <- cor(train_predictions, trainDF$MFI_FC, method = "pearson")
spearman_cor <- cor(train_predictions, trainDF$MFI_FC, method = "spearman")

# Display correlation results
cat("Pearson Correlation: ", pearson_cor, "\n")
```

```
## Pearson Correlation: 0.9782773
```

```
cat("Spearman Correlation: ", spearman_cor, "\n")
```

```
## Spearman Correlation:  0.9841374
```

```
# Create a correlation plot  
library(ggplot2)
```

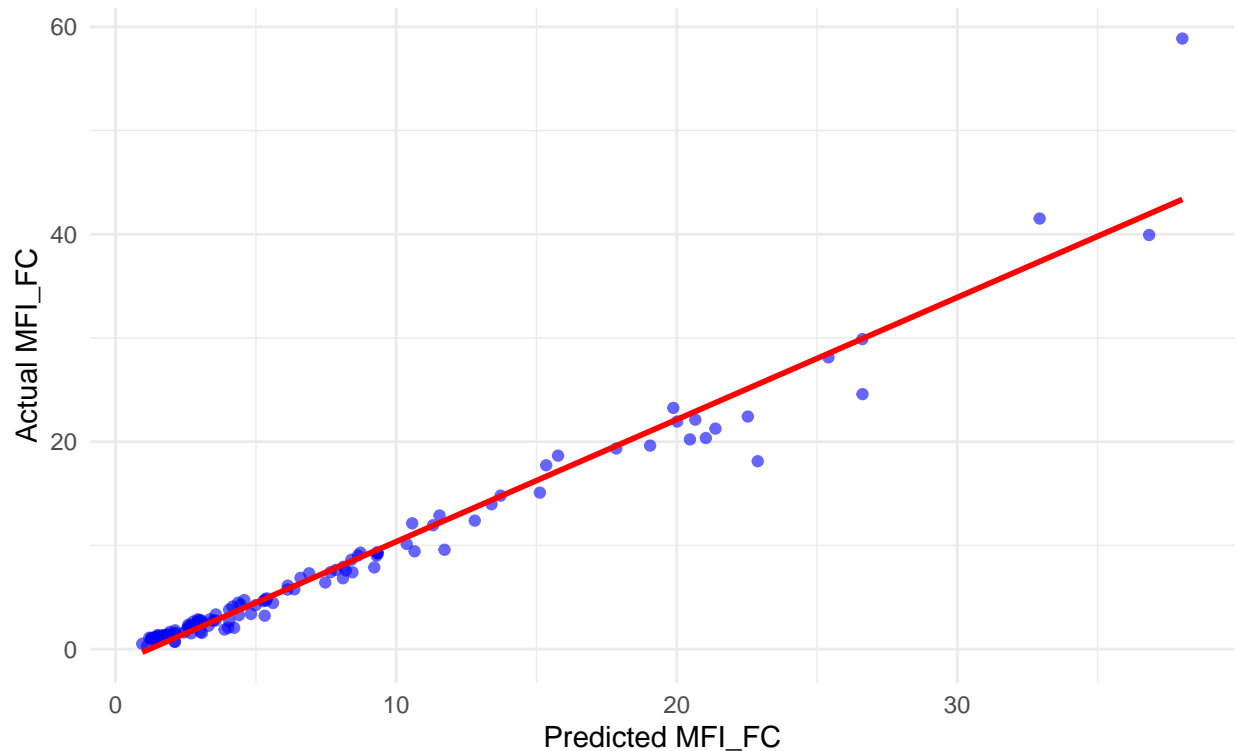
```
## Warning: package 'ggplot2' was built under R version 4.3.3
```

```
correlation_plot <- ggplot(data = data.frame(  
  Predicted = train_predictions,  
  Actual = trainDF$MFI_FC  
) , aes(x = Predicted, y = Actual)) +  
  geom_point(alpha = 0.6, color = "blue") + # Scatter plot  
  geom_smooth(method = "lm", color = "red", se = FALSE) + # Regression line  
  labs(  
    title = "Model Validation: Predicted vs Actual",  
    subtitle = paste0("Pearson: ", round(pearson_cor, 2),  
                      " | Spearman: ", round(spearman_cor, 2)),  
    x = "Predicted MFI_FC",  
    y = "Actual MFI_FC"  
  ) +  
  theme_minimal()  
  
# Display the plot  
print(correlation_plot)
```

```
## 'geom_smooth()' using formula = 'y ~ x'
```

Model Validation: Predicted vs Actual

Pearson: 0.98 | Spearman: 0.98



```
# Combine predictors and response
trainDF <- rbind(xDF, x2DF, x3DF) %>%
  as.data.frame() %>%
  mutate(MFI_FC = c(yDF$MFI_FC, yobs1$MFI_FC, yobs2$MFI_FC))

# Apply k-NN imputation
train_matrix <- as.matrix(trainDF)      # Convert to matrix
imputed_matrix <- impute.knn(train_matrix)$data
trainDF <- as.data.frame(imputed_matrix) # Convert back to dataframe

trainDF$MFI_FC <- sqrt(trainDF$MFI_FC)

set.seed(3)
auto_fit <- auto_ml() %>%
  set_engine("h2o", max_runtime_secs = 5) %>%
  set_mode("regression") %>%
  fit(MFI_FC ~ ., data = trainDF)

# Predict on training data
train_predictions <- predict(auto_fit, new_data = trainDF)$pred

# Calculate correlations
pearson_cor <- cor(train_predictions, trainDF$MFI_FC, method = "pearson")
spearman_cor <- cor(train_predictions, trainDF$MFI_FC, method = "spearman")

# Display correlation results
```

```
cat("Pearson Correlation: ", pearson_cor, "\n")
```

```
## Pearson Correlation: 0.9964828
```

```
cat("Spearman Correlation: ", spearman_cor, "\n")
```

```
## Spearman Correlation: 0.9969249
```

```
# Create a correlation plot
```

```
library(ggplot2)
```

```
correlation_plot <- ggplot(data = data.frame(  
  Predicted = train_predictions,  
  Actual = trainDF$MFI_FC  
) , aes(x = Predicted, y = Actual)) +  
  geom_point(alpha = 0.6, color = "blue") + # Scatter plot  
  geom_smooth(method = "lm", color = "red", se = FALSE) + # Regression line  
  labs(  
    title = "Model Validation: Predicted vs Actual",  
    subtitle = paste0("Pearson: ", round(pearson_cor, 2),  
                      " | Spearman: ", round(spearman_cor, 2)),  
    x = "Predicted MFI_FC",  
    y = "Actual MFI_FC"  
  ) +  
  theme_minimal()
```

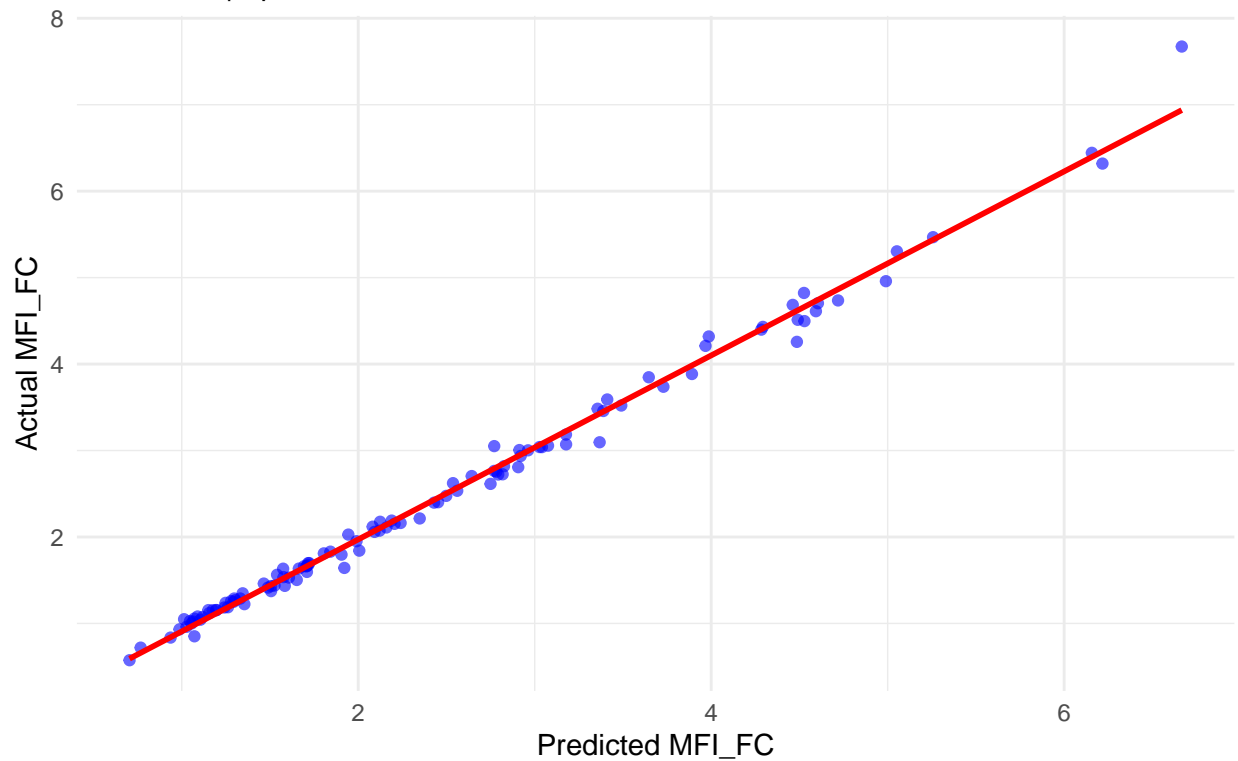
```
# Display the plot
```

```
print(correlation_plot)
```

```
## 'geom_smooth()' using formula = 'y ~ x'
```

Model Validation: Predicted vs Actual

Pearson: 1 | Spearman: 1



```
yhat <- predict(auto_fit, new_data = x4DF)$pred
rhat <- rank(-1 * yhat) # Rank predictions in descending order
print(cbind(rownames(x4DF), rhat))
```

```
##          rhat
## [1,] "142" "5"
## [2,] "146" "40"
## [3,] "163" "2"
## [4,] "124" "22"
## [5,] "134" "12"
## [6,] "170" "16"
## [7,] "132" "49"
## [8,] "140" "27"
## [9,] "155" "39"
## [10,] "172" "14"
## [11,] "148" "30"
## [12,] "135" "33"
## [13,] "123" "25"
## [14,] "128" "42"
## [15,] "164" "9"
## [16,] "159" "31"
## [17,] "167" "8"
## [18,] "158" "37"
## [19,] "151" "1"
## [20,] "150" "17"
```

```
## [21,] "133" "13"
## [22,] "126" "20"
## [23,] "125" "32"
## [24,] "130" "45"
## [25,] "145" "10"
## [26,] "122" "54"
## [27,] "138" "28"
## [28,] "157" "36"
## [29,] "119" "47"
## [30,] "136" "3"
## [31,] "149" "50"
## [32,] "165" "46"
## [33,] "131" "21"
## [34,] "169" "4"
## [35,] "160" "38"
## [36,] "168" "34"
## [37,] "141" "29"
## [38,] "154" "6"
## [39,] "153" "7"
## [40,] "147" "53"
## [41,] "156" "51"
## [42,] "121" "26"
## [43,] "120" "35"
## [44,] "144" "52"
## [45,] "143" "43"
## [46,] "171" "19"
## [47,] "127" "15"
## [48,] "129" "23"
## [49,] "162" "44"
## [50,] "166" "24"
## [51,] "152" "48"
## [52,] "161" "41"
## [53,] "137" "11"
## [54,] "139" "18"
```

```
# Read submission file
submission_file <- file.path(workDir, "3rdChallengeSubmissionTemplate_revised.tsv")
data <- read_tsv(submission_file)

# Update rankings for Challenge 1.2
ranking_df <- data.frame(
  SubjectID = as.numeric(rownames(x4DF)),
  `1.2) IgG-PT-D14-FC-Rank` = rhat,
  check.names = FALSE # Prevent automatic renaming of column names
)

data <- data %>%
  mutate(
    `1.2) IgG-PT-D14-FC-Rank` = ifelse(
      SubjectID %in% ranking_df$SubjectID,
      ranking_df[`1.2) IgG-PT-D14-FC-Rank`[match(SubjectID, ranking_df$SubjectID)],
      `1.2) IgG-PT-D14-FC-Rank`
    )
  )
```



```
# Save updated file
write_tsv(data, submission_file)
```

```
agua::h2o_end()
sessionInfo()
```

```
## R version 4.3.1 (2023-06-16 ucrt)
## Platform: x86_64-w64-mingw32/x64 (64-bit)
## Running under: Windows 11 x64 (build 22631)
##
## Matrix products: default
##
##
## locale:
## [1] LC_COLLATE=English_United States.utf8
## [2] LC_CTYPE=English_United States.utf8
## [3] LC_MONETARY=English_United States.utf8
## [4] LC_NUMERIC=C
## [5] LC_TIME=English_United States.utf8
##
## time zone: America/Los_Angeles
## tzcode source: internal
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  methods   base
##
## other attached packages:
## [1] ggplot2_3.5.1  impute_1.76.0  tibble_3.2.1  agua_0.1.4    parsnip_1.2.1
## [6] glmnet_4.1-8   Matrix_1.5-4.1 readr_2.1.5   tidyr_1.3.1   dplyr_1.1.3
##
## loaded via a namespace (and not attached):
## [1] tidysselect_1.2.1    timeDate_4041.110  farver_2.1.2
## [4] bitops_1.0-9         fastmap_1.2.0      RCurl_1.98-1.16
## [7] digest_0.6.37        rpart_4.1.19       timechange_0.3.0
## [10] lifecycle_1.0.4      yardstick_1.3.1    survival_3.5-5
## [13] magrittr_2.0.3        compiler_4.3.1     rlang_1.1.1
## [16] tools_4.3.1          utf8_1.2.3         yaml_2.3.10
## [19] data.table_1.16.2    knitr_1.49         labeling_0.4.3
## [22] bit_4.5.0            curl_6.0.1         DiceDesign_1.10
## [25] withr_3.0.2          purrr_1.0.2        workflows_1.1.4
## [28] h2o_3.44.0.3         nnet_7.3-19        grid_4.3.1
## [31] tune_1.2.1           fansi_1.0.4        colorspace_2.1-0
## [34] future_1.34.0        globals_0.16.3     scales_1.3.0
## [37] iterators_1.0.14     MASS_7.3-60        cli_3.6.1
## [40] crayon_1.5.3         rmarkdown_2.29     generics_0.1.3
## [43] rstudioapi_0.17.1    future.apply_1.11.3 tzdb_0.4.0
## [46] splines_4.3.1        dials_1.3.0        parallel_4.3.1
## [49] vctrs_0.6.3          hardhat_1.4.0      jsonlite_1.8.9
## [52] hms_1.1.3            bit64_4.5.2        listenv_0.9.1
## [55] foreach_1.5.2        gower_1.0.1        recipes_1.1.0
## [58] glue_1.6.2           parallelly_1.39.0  codetools_0.2-19
## [61] rsample_1.2.1        lubridate_1.9.3    shape_1.4.6.1
## [64] gtable_0.3.6         munsell_0.5.1      GPfit_1.0-8
```

## [67]	pillar_1.9.0	furrr_0.3.1	htmltools_0.5.8.1
## [70]	ipred_0.9-15	lava_1.8.0	R6_2.5.1
## [73]	lhs_1.2.0	vroom_1.6.5	evaluate_1.0.1
## [76]	lattice_0.21-8	class_7.3-22	Rcpp_1.0.11
## [79]	nlme_3.1-162	prodlim_2024.06.25	mgcv_1.8-42
## [82]	xfun_0.48	pkgconfig_2.0.3	