

Prognostic of NASA Turbofan Jet Engine

Fernando Chafim

05/01/2021

Contribuciones	Firma
Investigación previa	Fernando Chafim
Redacción de las respuestas	Fernando Chafim
Desarrollo código	Fernando Chafim

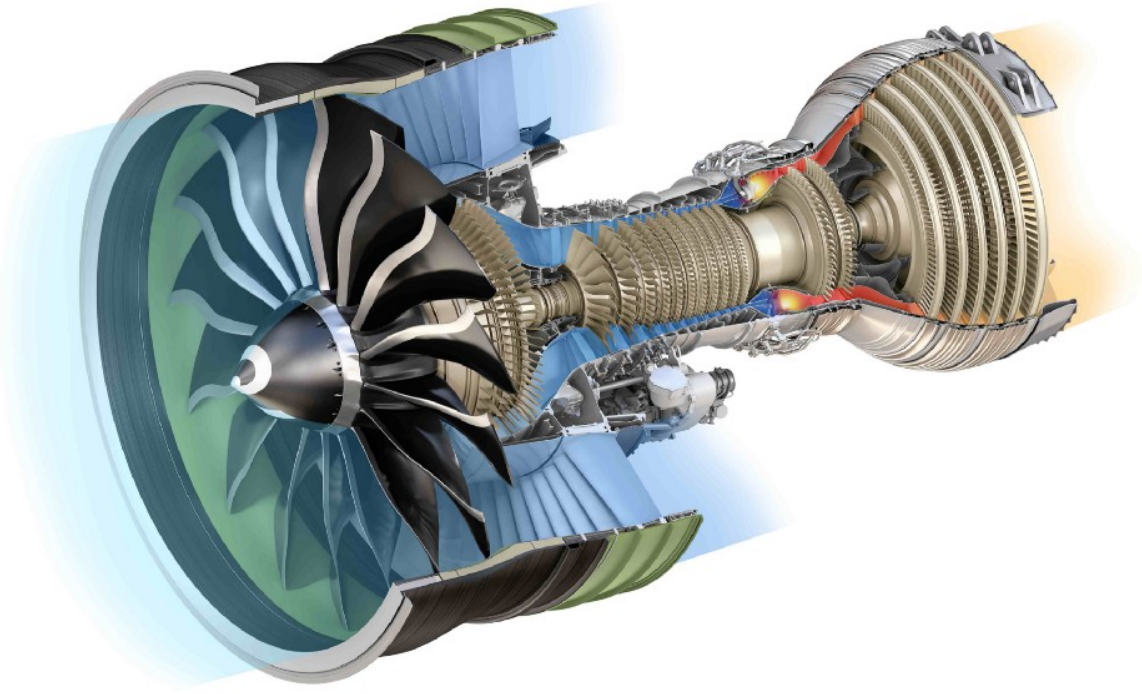
NASA TURBOFAN JET ENGINE

We define prognostics here exclusively as the estimation of remaining useful component life. The remaining useful life (RUL) estimates are in units of time (e.g., hours or cycles).

End-of-life can be subjectively determined as a function of operational thresholds that can be measured. These thresholds depend on user specifications to determine safe operational limits.

Prognostics is currently at the core of systems health management. Reliably estimating remaining life holds the promise for considerable cost savings (for example by avoiding unscheduled maintenance and by increasing equipment usage) and operational safety improvements. Remaining life estimates provide decision makers with information that allows them to change operational characteristics (such as load) which in turn may prolong the life of the component. It also allows planners to account for upcoming maintenance and set in motion a logistics process that supports a smooth transition from faulty equipment to fully functional.

```
knitr::include_graphics("img/turbofan1.jpeg")
```



2. Dataset

Los pronósticos y la gestión de la salud son un tema importante en la industria para predecir el estado de los activos para evitar tiempos de inactividad y fallas. Este conjunto de datos es la versión de Kaggle del muy conocido conjunto de datos públicos para el modelado de degradación de activos de la NASA. Incluye datos simulados Run-to-Failure de motores a reacción con turboventilador.

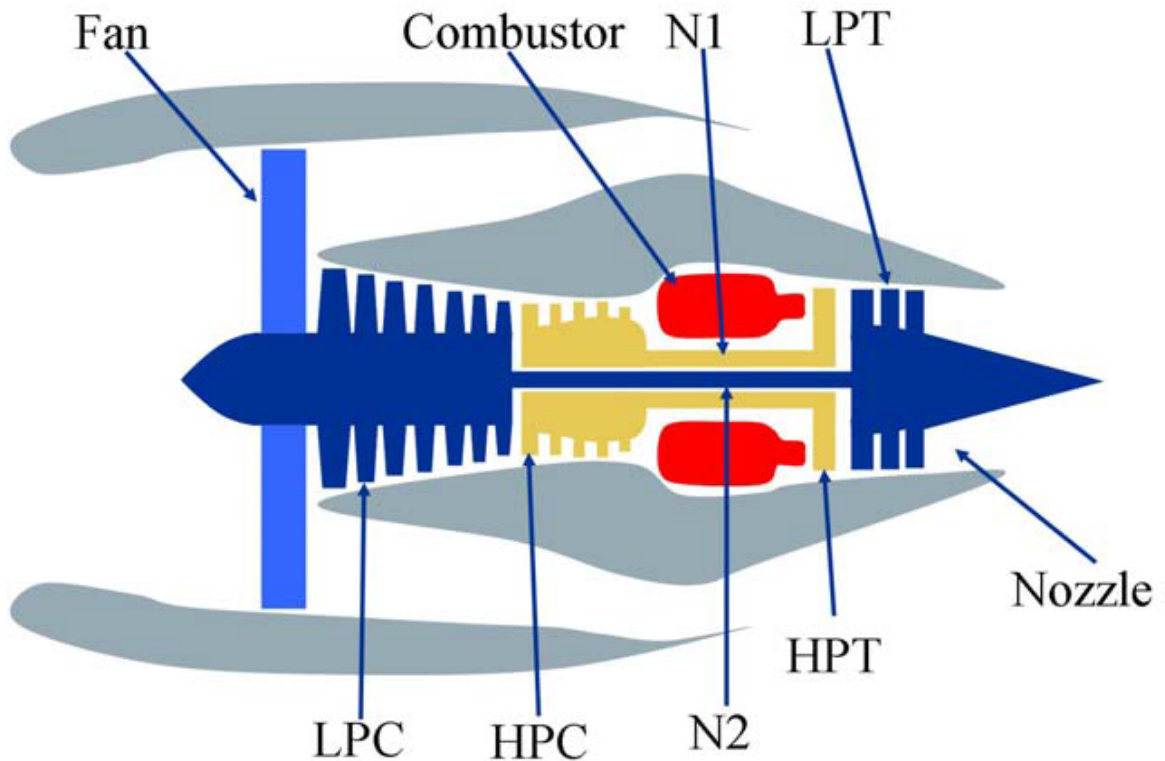
La simulación de la degradación del motor se llevó a cabo utilizando C-MAPSS. Se simularon cuatro conjuntos diferentes bajo diferentes combinaciones de condiciones operativas y modos de falla. Registra varios canales de sensores para caracterizar la evolución de la falla. El conjunto de datos fue proporcionado por el CoE de Pronósticos de NASA Ames.

Descripción

Symbol	Description	Units
Unit	-	-
Time	-	t
Altitude	Altitude	ft
Mach_N	Mach Number	M
SeaTemp	Sea Level temperature	°F
T2	Total temperature at fan inlet	°R
T2	Total temperature at fan inlet	°R
T24	Total temperature at LPC outlet	°R
T30	Total temperature at HPC outlet	°R
T50	Total temperature at LPT outlet	°R
P2	Pressure at fan inlet	psia
P15	Total pressure in bypass-duct	psia
P30	Total pressure at HPC outlet	psia

Symbol	Description	Units
Nf	Physical fan speed	rpm
Nc	Physical core speed	rpm
epr	Engine pressure ratio (P50/P2)	–
Ps30	Static pressure at HPC outlet	psia
phi	Ratio of fuel flow to Ps30	pps/psi
NRf	Corrected fan speed	rpm
NRc	Corrected core speed	rpm
BPR	Bypass Ratio	–
farB	Burner fuel-air ratio	–
htBleed	Bleed Enthalpy	–
Nf_dmd	Demanded fan speed	rpm
PCNfR_dmd	Demanded corrected fan speed	rpm
W31	HPT coolant bleed	lbm/s
W32	LPT coolant bleed	lbm/s
T48	(EGT) Total temperature at HPT outlet	°R
SmFan	Fan stall margin	–
SmLPC	LPC stall margin	–
SmHPC	HPC stall margin	–

```
knitr::include_graphics("img/turbofan2.png")
```



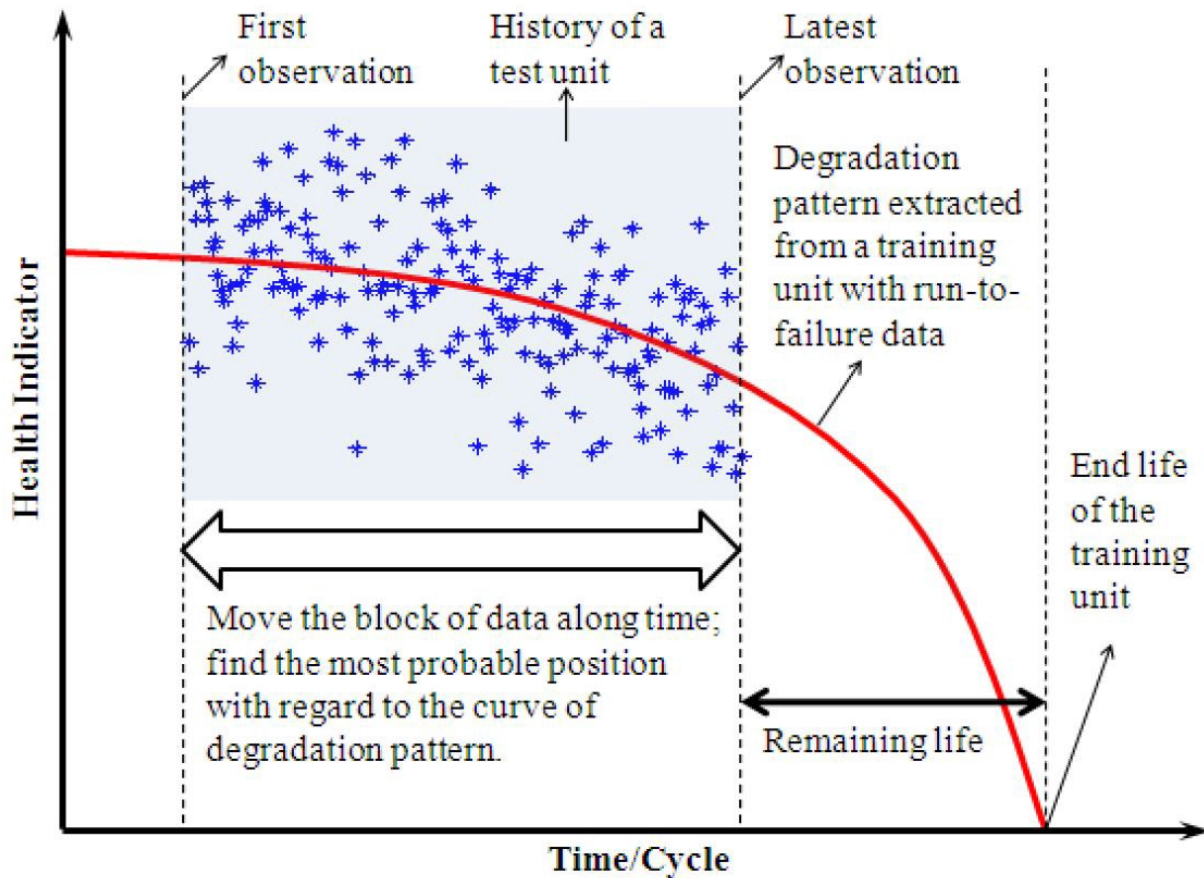
Importancia y objetivos de los análisis

Los pronósticos y la gestión de la salud son un tema importante en la industria para predecir el estado de los activos y evitar tiempos de inactividad y fallas. Este conjunto de datos es la versión de Kaggle del conjunto de datos públicos para el modelado de degradación de activos de la NASA. Incluye datos simulados Run-to-Failure de motores a reacción con turboventilador.

La simulación de la degradación del motor se llevó a cabo utilizando C-MAPSS. Se simularon cuatro conjuntos diferentes bajo diferentes combinaciones de condiciones operativas y modos de fallo. Registra varios canales de sensores para caracterizar la evolución de la falla. El conjunto de datos fue proporcionado por Prognostics CoE en NASA Ames.

En este conjunto de datos, el objetivo es predecir la **vida útil restante (RUL)** de cada motor. La vida útil restante (RUL) es el período de tiempo que es probable que funcione una máquina antes de que requiera reparación o reemplazo. Al tener en cuenta RUL, los ingenieros pueden programar el mantenimiento, optimizar la eficiencia operativa y evitar tiempos de inactividad no planificados. Por esta razón, estimar el RUL es una prioridad máxima en los programas de mantenimiento predictivo.

```
knitr::include_graphics("img/turbofan3.jpeg")
```



2. Integración y selección de los datos de interés a analizar.

Integración

La integración o fusión de los datos consiste en la combinación de datos procedentes de múltiples fuentes, con el fin de crear una estructura de datos coherente y única que contenga mayor cantidad de información.

```

library(dplyr)
library(e1071)
library(zeallot)
library(stringr)
library(data.table)
library(DT)
library(DataExplorer)
library(ggcorrplot)
library(plotly)
library(anomalize)
library(zoo)
library(tibbletime)
library(h2o)
library(isoform)
library(boot)
library(mice)
library(ggpubr)
library(MASS)
library(TTR)
library(caret)
library(FNN)
library(dbSCAN)
library(outliers)
library(car)

```

```

trainset <- fread("D:/UOC/Tipología y ciclo de vida de los datos/PRA2/CMaps/train_FD001.txt")
names(trainset) <- c('unit_number', 'time_in_cycles', 'altitude', 'mach_number', 'TRA', 'T2', 'T24', 'T30', 'T50', 'P2', 'Nc', 'epr', 'Ps30', 'phi', 'NRf', 'NRc', 'BPR', 'farB', 'htBleed', 'Nf_dmd', 'PCNfR_dmd', 'W31', 'W32')

testset <- fread("D:/UOC/Tipología y ciclo de vida de los datos/PRA2/CMaps/test_FD001.txt")
names(testset) <- c('unit_number', 'time_in_cycles', 'altitude', 'mach_number', 'TRA', 'T2', 'T24', 'T30', 'T50', 'P2', 'Nc', 'epr', 'Ps30', 'phi', 'NRf', 'NRc', 'BPR', 'farB', 'htBleed', 'Nf_dmd', 'PCNfR_dmd', 'W31', 'W32')

y_testset <- fread("D:/UOC/Tipología y ciclo de vida de los datos/PRA2/CMaps/RUL_FD001.txt") %>% pull()

read_logs <- function(path = "D:/UOC/Tipología y ciclo de vida de los datos/PRA2/CMaps"){

  files <- c("FD001.txt", "FD002.txt", "FD003.txt", "FD004.txt")
  sets <- c("train", "test")
  files <- apply(expand.grid(sets, files), 1, paste, collapse="_")
  f <- file.path(path, files)
  d <- lapply(f, fread)
  names(d) <- str_remove(files, ".txt")

  columns = c('unit_number', 'time_in_cycles', 'altitude', 'mach_number', 'TRA', 'T2', 'T24', 'T30', 'T50', 'P2', 'Nc', 'epr', 'Ps30', 'phi', 'NRf', 'NRc', 'BPR', 'farB', 'htBleed', 'Nf_dmd', 'PCNfR_dmd', 'W31', 'W32')

  temp_func = function(df){
    names(df) <- columns
    return(df)
  }

  d <- lapply(d, temp_func)
}

```

```

    return(d)
}

#dfs <- read_logs()

#trainset <- bind_rows(
#  dfs$train_FD001#,
#  dfs$train_FD002,
#  dfs$train_FD003,
#  dfs$train_FD004
# )

#testset <- bind_rows(
#  dfs$test_FD001#,
#  dfs$test_FD002,
#  dfs$test_FD003,
#  dfs$test_FD004
# )

#trainset$set <- "train"
#testset$set <- "test"

df <- trainset

```

Selección

Features Selection 1

```

FindOutliers <- function(x) {
  if(class(x)%in% c("numeric", "integer")){
    lowerq = quantile(x, na.rm = TRUE)[2]
    upperq = quantile(x, na.rm = TRUE)[4]
    iqr = upperq - lowerq #Or use IQR(data)
    # we identify extreme outliers
    extreme.threshold.upper = (iqr * 3) + upperq
    extreme.threshold.lower = lowerq - (iqr * 3)
    result <- which(x > extreme.threshold.upper | x < extreme.threshold.lower)
    output <- length(result)
  } else {
    output <- 0
  }

  return(output)
}

```

```

DataProfiling <- function(df){

```

```

  column_names <- colnames(df)
  column_classes <- as.vector(sapply(df, function(x) class(x)))
  column_mean <- as.vector(sapply(df, function(x) if(class(x)%in% c("numeric", "integer")){mean(x, na.rm = TRUE)} else {NA})))
  column_sd <- as.vector(sapply(df, function(x) if(class(x)%in% c("numeric", "integer")){sd(x, na.rm = TRUE)} else {NA})))
  column_median <- as.vector(sapply(df, function(x) if(class(x)%in% c("numeric", "integer")){median(x, na.rm = TRUE)} else {NA})))
  column_max <- as.vector(sapply(df, function(x) if(class(x)%in% c("numeric", "integer")){max(x, na.rm = TRUE)} else {NA})))

```

```

column_min <- as.vector(sapply(df, function(x) if(class(x)%in% c("numeric", "integer")){min(x, na.rm = TRUE)}))
column_nunique <- as.vector(sapply(df, function(x) if(class(x)%in% c("numeric", "integer")){length(unique(x))})))
column_quantile0 <- as.vector(sapply(df, function(x) if(class(x)%in% c("numeric", "integer")){quantile(x, 0)})))
column_quantile25 <- as.vector(sapply(df, function(x) if(class(x)%in% c("numeric", "integer")){quantile(x, 0.25)})))
column_quantile50 <- as.vector(sapply(df, function(x) if(class(x)%in% c("numeric", "integer")){quantile(x, 0.5)})))
column_quantile75 <- as.vector(sapply(df, function(x) if(class(x)%in% c("numeric", "integer")){quantile(x, 0.75)})))
column_quantile100 <- as.vector(sapply(df, function(x) if(class(x)%in% c("numeric", "integer")){quantile(x, 1)})))
column_interquartile_range <- as.vector(sapply(df, function(x) if(class(x)%in% c("numeric", "integer")){quantile(x, 0.75) - quantile(x, 0.25)})))
column_skewness <- as.vector(sapply(df, function(x) if(class(x)%in% c("numeric", "integer")){skewness(x)})))
column_kurtosis <- as.vector(sapply(df, function(x) if(class(x)%in% c("numeric", "integer")){kurtosis(x)})))
column_na <- as.vector(sapply(df, function(x) sum(is.na(x)))))
column_zero <- as.vector(sapply(df, function(x) sum(x==0))))
column_outliers <- as.vector(sapply(df, FindOutliers))

df_table <- tibble(
  names = column_names,
  classes = column_classes,
  min = column_min,
  max = column_max,
  n_unique = column_nunique,
  mean = column_mean,
  sd = column_sd,
  median = column_median,
  quantile_0 = column_quantile0,
  quantile_25 = column_quantile25,
  quantile_50 = column_quantile50,
  quantile_75 = column_quantile75,
  quantile_100 = column_quantile100,
  interquartile_range = column_interquartile_range,
  skewness = column_skewness,
  kurtosis = column_kurtosis,
  MissingValues = column_na,
  n_zero = column_zero,
  Outliers = column_outliers
)

return(df_table)
}

temp <- DataProfiling(df)

```

```

unique_columns <- temp$names[temp$n_unique == 1 & !is.na(temp$n_unique)]
print(unique_columns)

```

```

## [1] "TRA"      "T2"      "P2"      "epr"      "farB"     "Nf_dmd"
## [7] "PCNfR_dmd"

```

```

df <- df %>% dplyr::select(-one_of(unique_columns))
trainset <- trainset %>% dplyr::select(-one_of(unique_columns))
testset <- testset %>% dplyr::select(-one_of(unique_columns))

```


Preprocesado de los datos

Feature Engineering 1

Remaining Useful Life

```
add_remaining_useful_life <- function(df){  
  
  df <- df %>%  
    group_by(unit_number) %>%  
    mutate(max_time_in_cycles = max(time_in_cycles)) %>%  
    ungroup() %>%  
    mutate(  
      RUL = max_time_in_cycles - time_in_cycles  
    ) %>%  
    dplyr::select(-max_time_in_cycles)  
  
  return(df)  
}
```

Cumulative terms

```
add_cumulative_features <- function(df, columns){  
  df <- df %>%  
    as_tibble() %>%  
    group_by(unit_number) %>%  
    mutate(  
      across(  
        all_of(columns),  
        list(  
          cumsum = cumsum,  
          cummin = cummin,  
          cummax = cummax  
        ),  
        .names = "{.fn}_{.col}"  
      )  
    )  
  return(df)  
}
```

Lag terms

```
add_lag <- function(df, columns, num = 1){  
  
  for (i in 1:length(num)) {  
    num_temp <- num[i]  
    df <- df %>%  
      as_tibble() %>%  
      group_by(unit_number) %>%  
      mutate(  
        across(all_of(columns), ~lag(.x, n = num_temp, default = NA), .names = paste0("lag", num_temp,
```



```

    )
  }
  return(df)
}

```

Moving Average terms

```

add_moving_avg_features <- function(df, columns, num = 5){

  for (i in 1:length(num)) {
    num_temp <- num[i]
    df <- df %>%
      as_tibble() %>%
      group_by(unit_number) %>%
      mutate(
        across(
          all_of(columns),
          list(
            rollmean = ~rollapplyr(.x, num_temp, FUN = mean, fill=NA),
            rollsd = ~rollapplyr(.x, num_temp, FUN = sd, fill=NA),
            rolliqr = ~rollapplyr(.x, num_temp, FUN = IQR, fill=NA)
          ),
          .names = paste0("{.fn}", num_temp, "_{.col}")
        )
      )
  }
  return(df)
}

```

Interaction terms

```

add_interaction_terms <- function(df, columns){
  df_1 <- df %>% dplyr::select(all_of(columns))
  df_2 <- df %>% dplyr::select(!all_of(columns))
  df_interaction <- as.data.frame(model.matrix(~ .^2-1, df_1))
  df_output <- bind_cols(df_2, df_interaction)
  names(df_output) <- gsub(x = names(df_output), pattern = ":", replacement = "_")
  return(df_output)
}

```

Todos juntos

```

columns <- c("altitude", "mach_number", "T24", "T30", "T50", "P15", "P30", "Nf", "Nc", "Ps30", "phi", "NI")

trainset2 <- trainset %>%
  add_remaining_useful_life() %>%
  add_interaction_terms(columns) %>%
  add_lag(columns=columns, num = 1:5) %>%
  add_moving_avg_features(columns = columns, num = c(3, 5)) %>%
  add_cumulative_features(columns=columns)

```

```

y_trainset <- trainset2$RUL

testset2 <- testset %>%
  add_interaction_terms(columns) %>%
  add_lag(columns=columns, num = 1:5) %>%
  add_moving_avg_features(columns = columns, num = c(3, 5)) %>%
  add_cumulative_features(columns=columns)

```

Benchmark

```

bechmark_performance <- function(x_train, y_train, x_test, y_test){

  df <- bind_cols(x_train, y= y_train)
  model_bench <- lm(y~., data=df)

  x_test <- x_test %>%
    group_by(unit_number) %>%
    filter(time_in_cycles == max(time_in_cycles))

  suppressWarnings({
    y_hat <- predict.lm(model_bench, x_test)
  })
  result <- postResample(y_hat, y_test)

  return(result)
}

bechmark_performance(x_train = trainset,
  y_train = y_trainset,
  x_test = testset,
  y_test = y_testset)

```

```

##          RMSE    Rsquared      MAE
## 31.7404750  0.6525414 26.1234550

```

Análisis estadístico descriptivo

```

temp <- DataProfiling(trainset2)

#datatable(temp, class = 'cell-border stripe', rownames = FALSE, filter = 'top',
#           options = list(
#             pageLength = 10,
#             autoWidth = TRUE,
#             scrollX = TRUE,
#             scrollY = TRUE
#           )
#           )
knitr::kable(temp, "pipe")

```

[illegible]

[illegible]

names	class	chain	max	n	uniqu	mean	sd	median	quantile	deaf	ftdeaf	deaf50	deaf75	deaf100	skewness	Missing	Vz	Outliers
T24_P80me	353993	575011003335560382366429707525298005355906623029589731050348394500	-	0	0	0												
T24_Nfume	1531195	393506070580078424087834792328190029917930232955650606009738039002763	0	0	0													
T24_Nume	579964954252000619000668488821157293528065682015520357050513286872523968708151	0	806															
T24_Ps30me	30089371240004700053.081.0781223725089094420303302001678073000253408070902595	0	0	0														
T24_phume	3381636290505700002.27284336553.3381630688005114333363027945227040500	-	0	0	0													
T24_NRfme	1531215.392410029500783288213746923.1855589910970325985308934756096032035936	0	0	0														
T24_NRme	521034642050000200836504683702271176502055207633076860220009892261593332898	0	808															
T24_BRRme	5345.05278861462025.62808072522.85348.00475.88426.48448.36628.4861360950288360	0	0	0														
T24_htBleed	189485758800000000709.12792925774890050850208072053002858000080000088765	0	0	0														
T24_W31me	1539235700008200046.114726846632568004874000450056100331080338040000	-	0	0	0													
T24_W32me	1734.0920294722643967.62442214951.70834.600274407355958411.697628478233665	-	0	0	0													
T30_T50me	182503.820070589040029.11832364908256202026326848209372007528065003087041	0	0	0														
T30_Pi5me	33934.40400025100070.8925833352369.004027084166200645030571025.6893300472056713	0	0															
T30_P80me	870196783900008800402571.48800678737908783489008808783076790008870380808159219	0	0															
T30_Nfume	75156859500038798023.7106.178723473.063089520702240053860500458.40360000637482	0	0															
T30_Nume	12238745808807620867283584890092.287086000000004408370800888067620900207259103	0	124															
T30_Ps30me	73869.78032007000016.662538538273830.0813607254070034008332876.697095073266	0	0	0														
T30_phume	8203743328200158000125402182925620376807298202356209783952839725700950043285287	0	0															
T30_NRfme	751643625000000298022.710437083370509878812507039500986602304042008500243188	0	0															
T30_NRme	12789233.634880032950834.710209890780220007800279586023607882085002090939242	0	146															
T30_BRRme	13450.10793905230027.5748602556132151.10319368210722489.18793.87552034681888	0	0	0														
T30_htBleed	1270564287000020004143327686353.827060202800006680600237600900000706656	0	0	0														
T30_W31me	60902.623600872300037.2322395674561912405780607256089400566007248835007452	0	0	0														
T30_W32me	6515.370274206037042.337410005423694130927930782939318755625044007005875	0	0	0														
T50_Pi5me	99870312250499500046.78463766527228700203507832730283803615509875004091055	0	0	0														
T50_P80me	667067944000029000652160306726396670007675007355582260020755090013007366	0	0	0														
T50_Nfume	3300730402510006800020.52338521037012300658072533081230007300068000404020	0	0	0														
T50_Nume	12518662.04206007070086.11322735525.887200422005032800742004209070397870255089	0	16															
T50_Ps30me	65048.6850708239500984.373326871266268064006091100246700947502450013334684	0	0	0														
T50_phume	7290749732008700633888568432225196707075033225303600735094410003704134	0	0	0														
T50_NRfme	330073742700067000029.5253302683106320980289052308320037030070009360232	0	0	0														

names	class	chain	max	n_uniq	mean	sd	median	quantile	mean	quantile	mean	quantile	mean	quantile	mean	quantile	mean	quantile	mean	quantile	missing	Vz	Outliers
T50_NRR	mer	12525518087566074706485612035339252740003450903775037108083330879127000733128	0	18																			
T50_BRR	mer	1589.720112005432894.6801217937680282408012407288026248384751.03205253798	0	0	0																		
			0.1532927																				
T50_htB	mer	5390775236000056001537982531833910705000956901059723297236080800000000	0	0	0																		
			0.1415006																				
T50_W31	mer	53824.5556000220500882386350586543220385305705865024023560002600000000	0	0	0																		
			0.0262919																				
T50_W32	mer	32253.334392706832812.0266249302323595678093284123299958406568682501035	0	0	0																		
			0.0558001																				
P15_P30	mer	11882.228560956600958.11670240069.188820095600249060871002006257669000	-	0	0	0																	
			0.405411470229																				
P15_Nf	mer	51579.5072007816600062887001606.525490002055000650397070516078600000	32.2419137	0	410																		
			5.3061050																				
P15_Nnum	er	194959987567659996898.7765805728.949500953020795948958002773958.002256091.3733818	0	806																			
P15_Ps30	mer	11.90008.023300027.3553273726.60110900200326000130.7074807333350004636958	0	0	0																		
			0.1754078																				
P15_phin	mer	1208.800000011800067.65431162269.11828800297002600062800023102011578000	-	0	0	0																	
			0.453633277055																				
P15_NRf	mer	51578.5160007816600062879320006.525490002055000650397070516078600000	32.0831441	0	410																		
			5.2854979																				
P15_NRR	mer	175039.792270089205084.82364859161.758300707510399150008300322828.622070978722532	0	690																			
P15_BRR	mer	179.9011889072348082.43381231022.364729011859438236482941181.6075287883820207	0	0	0																		
			0.1236243																				
P15_htB	mer	8380.880000000008097.28553118432.78380080700080020850003860000002200000081487	0	0	0																		
			0.0435012																				
P15_W31	mer	824.20520082780838.81389488237.118240083609073000083006362008232000000	-	0	0	0																	
			0.36207324674																				
P15_W32	mer	494.743062700877603.2853290503.4642474506522301.7552499522207043702085	-	0	0	0																	
			0.35347641685																				
P30_Nf	mer	134321377826632140002583076292473.83215300202806713899033288333027439000	-	0	0	0																	
			0.391011698022																				
P30_Nnum	er	1969458230200458000008062306324965050008000027023000000231502600040801088352	0	105																			
P30_Ps30	mer	15971.287600231020007.55290726267259700622402728703370027757028400093792444	0	0	0																		
			0.1206391																				
P30_phin	mer	18570290741008728863327090258672871080800308801080108007800832394000	-	0	0	0																	
			0.478035776701																				
P30_NRf	mer	1343193278670019400025833328473.064002097539037020010080808070.8470000	-	0	0	0																	
			0.391011528048																				
P30_NRR	mer	1605768110000950008209974270694360570099920098950875448071200098007594579203	0	51																			
P30_BRR	mer	164.24738.400511071.5859509970.7463620655303707.70681507322801440012235246	0	0	0																		
			0.0507883																				
P30_htB	mer	150722870000000006864297862562.85072080000056006802007038800300007093545623	0	1																			
P30_W31	mer	11031.21880008820079.728688745622350006480004650725720088000389003500	-	0	0	0																	
			0.41884731213																				
P30_W32	mer	12621.333313037220887.86106031824.11278482851.8250468878483107320327203140	-	0	0	0																	
			0.41962235358																				
Nf_Nc	mer	154622881115689044062726210414395.4620000206902205642000723769014705004678765	0	816																			
Nf_Ps30	mer	1878.26892057100632.0192085260.587300087007950031400580807440000095011	0	0	0																		
			0.1737205																				
Nf_phin	mer	12387924935000024008573206801372137320003920036400432002834006400500	-	0	0	0																	
			0.43954451769																				
Nf_NRf	mer	57021540572000240003261295202297027000675002987027860582200624000874229682	0	6																			

names	class	min	max	n_unique	mean	sd	median	quantile	mean	ftile	mean	50	mean	75	tile	100	skewness	Missing	Vz	Outliers
Nf_NR	mean	193454	198699	6699	192080	6876	188206	189342	190009	190003	194080	198000	197094	192005	194066	6410	0	704		
Nf_BPR	mean	19879	20592	1820	19960	86	19628	1972	19879	19795	19552	19921	19502	1820	19807	389	1463	0	0	0
																	0.1234488			
Nf_ht	Bleed	192655	195526	6699	193900	2371	191356	192365	193000	193452	194000	195000	196000	197000	198000	650	0	0	0	
																	0.0425408			
Nf_W31	mean	91087	90228	2934	90096	429	89779	89720	90038	90092	90093	90120	90115	89900	-	0	0	0	0	
																	0.357501	34505		
Nf_W32	mean	14677	16740	5738	15618	262	15397	15665	16173	16645	16932	16801	16697	16000	160860	-	0	0	0	0
																	0.349603	29823		
Nc_Ps30	mean	12455	14691	1820	13007	1000	12895	13000	13073	13089	13090	13132	13106	11200	13000	950	0	23		
Nc_phin	mean	16853	33808	3003	17007	3032	15536	16883	17000	17000	17000	17000	17000	17000	17000	17000	17000	0	131	
Nc_NR	mean	11546	12707	5980	11860	627	11806	11785	11829	11800	11800	11800	11800	11800	11800	11800	11800	0	814	
Nc_NR	mean	73096	74672	13000	73806	525	73763	73900	74000	74000	74000	74000	74000	74000	74000	74000	74000	0	800	
Nc_BPR	mean	75406	78817	2058	7630	433	76046	75740	76232	76423	76732	76865	76904	76000	76000	76000	76000	0	56	
Nc_ht	Bleed	13517	7936	6600	14000	668	13339	13568	14205	15000	15000	15000	15000	15000	15000	15000	15000	0	119	
Nc_W31	mean	14495	15475	1408	14508	77	14789	14623	14700	14800	14800	14800	14800	14800	14800	14800	14800	0	0	
																	0.3580910			
Nc_W32	mean	12671	12325	1417	12052	630	12064	12228	12170	12222	12353	12310	12300	12300	12300	12300	12300	0	0	
																	0.3584811			
Ps30_phin	mean	14479	25230	7067	17008	418	16309	17026	17000	17000	17000	17000	17000	17000	17000	17000	17000	0	0	0
																	0.1262953			
Ps30_NR	mean	11878	17370	2601	13063	349	15724	14891	16870	16770	16785	16665	16500	16000	16000	16000	16000	0	0	0
																	0.1736164			
Ps30_NR	mean	1381	5389	3763	3002	1300	624	3920	6369	3817	5885	4380	6330	6228	6000	598	222	0	29	
Ps30_BPR	mean	1391	5241	1000	5844	001	357	2847	0071	9030	5024	800	900	500	874	205	049	766	235	372
																	0.1637307			
Ps30_ht	Bleed	14236	13200	6020	10893	192	13741	12751	13606	13572	13667	13687	13600	13600	13600	13600	13600	0	0	0
																	0.1489774			
Ps30_W31	mean	1822	4187	1000	4508	45	3300	4085	35	3052	4085	4000	4050	4000	4000	4000	4000	0	0	0
																	0.0483855			
Ps30_W32	mean	1090	4222	2800	7400	7	1395	635	1027	1108	1402	2828	2880	900	1881	881	2315	457	632	5159
phi_NR	mean	12387	6047	5700	1000	857	3595	625	4238	800	970	1100	730	1000	1000	1000	1000	0	0	0
																	0.439980	55827		
phi_NR	mean	12052	13523	6000	12000	60	1082	1000	1035	12000	1030	1030	1030	1030	1030	1030	1030	0	78	
phi_BPR	mean	1349	4423	3000	4001	827	2386	604	957	430	233	300	109	957	730	608	520	796	85	208
																	0.0502864			
phi_ht	Bleed	1261	2080	5000	1000	2432	792	1989	920	2600	600	800	600	600	600	600	600	0	1	
phi_W31	mean	1981	3205	2000	2003	9	1262	1302	1300	1238	1300	1300	1300	1300	1300	1300	1300	0	0	0
																	0.421016	85255		
phi_W32	mean	1898	5236	2000	2004	3	1260	650	649	1789	800	102	122	189	1870	800	1800	150	-	0
																	0.411009	63896		
NRf_NR	mean	19345	19870	6699	19208	6876	18820	18934	19000	19000	19408	19800	19709	19200	19406	6410	0	705		
NRf_BPR	mean	19879	20592	1820	19960	86	19628	1972	19879	19795	19552	19921	19502	1820	19807	389	1463	0	0	0
																	0.1234894			
NRf_ht	Bleed	19265	19552	6699	19390	2371	19135	19236	19300	19345	19400	19500	19600	19700	19800	650	0	0	0	
																	0.0425452			
NRf_W31	mean	91088	90228	2934	90096	429	89779	89720	90038	90092	90093	90120	90115	89900	-	0	0	0	0	
																	0.357501	345286		
NRf_W32	mean	14676	16740	5738	15618	262	15397	15665	16173	16645	16932	16801	16697	16000	160860	-	0	0	0	0
																	0.349603	2981525		
NRc_BPR	mean	67745	78752	2058	7630	433	76046	75740	76232	76423	76732	76865	76904	76000	76000	76000	76000	0	62	
NRc_ht	Bleed	16099	32700	2000	16000	1621	18309	13550	17000	18000	18000	18000	18000	18000	18000	18000	18000	0	129	

[illegible]

names	class	main	max	n_uniq	mean	sd	median	quantile	coef	theta	beta	gamma	delta	epsilon	zeta	eta	theta	iota	kappa	lambda	mu	nu	xi	missing	V	obs
lag2_T50me	1	38.25	468.69	0650	008.72	5798	27477.91	88225	4020	4007	9000	4368	5020	0000	4025	248	200	NA	0	0.2081566						
lag2_P15me	1	41.60	2000	0000	021.60	9800	3325	6102	0002	0200	0200	0200	0200	0000	0000	45.33	2026	NA	406	6.8801151						
lag2_P30me	1	49.85	5600	0050	053.38	7827	6523	4500	8555	2835	5845	5500	0255	6060	0900	000	-	200	NA	0	0.3579127	2003				
lag2_Nfume	2	387.92	8803	5700	0088.09	5299	2388.02	8709	2888	02888	02888	02888	0300	0000	0087	856	5271	NA	0							
lag2_Nrume	9	621.73	2008	8308	0064.92	7249	9055.63	0207	3058	0906	0006	0027	2088	6100	0206	18.79	6702	NA	749							
lag2_Ps80me	4	6c85	0480	8005	047.53	4923	5655.51	0468	5047	9404	0750	0470	0048	4800	3600	000	427	2748	200	NA	0	0.2236239				
lag2_phume	5	18.69	2008	0000	0621.42	9740	3634	4908	6952	0975	2049	5209	6520	6809	9900	000	-	200	NA	0	0.4017328	8299				
lag2_NRfume	2	387.83	8803	5700	0088.09	4689	2388.02	8709	2888	02888	02888	02888	0300	0000	0389	684	2254	NA	0							
lag2_NRume	8	699.98	2708	8007	08043.51	6339	85957.52	0009	8408	0904	0804	0827	0860	9400	0275	8.27	5215	NA	628							
lag2_BRRme	3	249	6848	0801	8.44	1336	678.73	8402	4984	4784	3840	6488	5848	0050	1003	556	420	200	NA	0	0.1393215					
lag2_htBleed	3	8.00	0000	0000	00893.17	9584	7363.00	8800	8020	0080	0080	0004	0000	0020	0000	0031	60489	200	NA	0	0.0842956					
lag2_W31me	3	1400	0300	0000	038.82	0136	7038	5308	0408	0700	8803	0809	0804	3000	0200	0000	-	200	NA	0	0.3178050	4882				
lag2_W32me	2	2c89	4200	8460	023.29	1962	5123.29	9209	4202	2250	9203	0672	0081	0431	500	-	200	NA	0	0.3172412	3220					
lag3_ahume	1	0.0087	058	-	0.0021	860	00000	-	0.0000	0015	0087	0030	000	-	300	NA	0	0.0087000	0.0000073	0.0087000	15000	0.0295352	4056			
lag3_machume	1	0.0006	040	0.0000	0250	2933	00000	-	0.0000	0003	0006	0005	0008	7298	300	NA	0	0.0006000	0.0000020	0.0006000	2000	1.1340397				
lag3_T24me	4	1d.21	6040	5000	00642.66	5787	6647.64	0016	2642	3262	6462	0986	0450	0660	000	2711	679	300	NA	0	0.1391064					
lag3_T30me	4	571.01	6000	9295	00590.35	9359	12770.01	6700	4586	0959	0005	0405	5169	7965	000	2698	408	300	NA	0	0.0127268					
lag3_T50me	1	38.25	468.69	0650	008.68	7472	93907.83	8225	4020	4007	8840	4220	8850	1930	0089	6036	300	NA	0	0.2092889						
lag3_P15me	1	41.60	2000	0000	021.60	9800	3329	6102	0002	0200	0200	0200	0200	0000	0000	45.09	13998	NA	406	6.8621995						
lag3_P30me	1	49.85	5600	0000	053.39	6357	5552.46	0085	5528	4558	4655	0025	5606	0800	000	-	300	NA	0	0.3412361	0770					
lag3_Nfume	2	387.92	8803	5700	0088.09	4548	2388.02	8709	2888	02888	02888	02888	0300	0000	0385	828	3330	NA	0							
lag3_Nrume	9	621.73	2004	7625	79064.72	1958	2035.61	0207	3058	0906	0006	0920	4766	0800	0208	3280	3173	NA	727							
lag3_Ps80me	4	6c85	0480	8005	047.53	1925	7574.51	0468	5047	9404	0750	0470	0048	4800	3500	000	406	3944	300	NA	0	0.2494484				
lag3_phume	5	18.95	2008	0000	0621.43	8732	8523.49	0809	5520	9852	0495	2096	5206	6809	9800	000	-	300	NA	0	0.3807224	7523				
lag3_NRfume	2	387.83	8803	5700	0088.09	4692	2388.02	8709	2888	02888	02888	02888	0300	0000	0369	693	7338	NA	0							
lag3_NRume	8	699.98	2704	6502	08043.40	7938	02058000	9840	8809	0408	0827	0465	0800	0203	1961	9353	NA	602								
lag3_BRRme	3	249	6848	0801	8.44	0919	8639	4382	0249	8404	4684	3820	6428	5848	0096	0042	4328	300	NA	0	0.1415849					
lag3_htBleed	3	8.00	0000	0000	00893.16	5273	6323.00	8800	8020	0080	0080	0004	0000	0020	0000	0030	54984	300	NA	0	0.0847392					
lag3_W31me	3	1400	0300	0000	038.82	1947	5328	3088	0408	0700	8803	0809	0804	3000	0200	0000	-	300	NA	0	0.3067723	2180				
lag3_W32me	2	2c89	4200	8460	023.29	3030	5233	3002	0942	3002	5703	0028	0682	0081	0423	000	-	300	NA	0	0.3050502	4219				

[illegible]

names	classes	n	min	max	n_unique	mean	sd	median	quantile	skewness	kurtosis	entropy	missing	NaN	Outliers		
lag5_NRMSE	float64	387	8.838800000000000	88.0926416328	2.0687083880023880	0.26880023880023880	0.0000000000000000	0.365176	500	NA	0	0.0101762					
lag5_NRMSE	float64	699	9.926800000000000	43.177633821456	4.000000000000000	0.9816008000000000	0.0805268000000000	0.7300000000000000	17.5233108	NA	548						
lag5_BRMSE	float64	324	9.067800000000000	8.440200357834	1.7780324900000000	4.3503778000000000	0.4634800000000000	0.9050000000000000	177248	500	NA	0	0.1434169				
lag5_hbRMSE	float64	8.000000000000000	0.00893.1354830	1.97.0088000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.2808471	500	NA	0	0.0992936					
lag5_WRMSE	float64	1408	0.4300000000000000	38.82501672338	4.4088000000000000	0.6884000000000000	0.0000000000000000	0.2400000	500	NA	0	0.2779845	2075				
lag5_WRMSE	float64	2289423	0.0084500000000000	23.295162335320	12.08942302282330	12.0000000000000000	0.0000000000000000	0.2500	500	NA	0	0.2763797	7333				
rollmeam3_merit	float64	0.0046670	-	0.0013626000000	-	0.0000000000000000	0.0000000000000000	0.3346667	17000	-	200	NA	1	0.0063000	0.0000087	0.0063000008667	0.02625093285
rollsd3_nation	float64	0.00000000647331	0.0019010010076181	0.0000000011790181	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.2594264733	1415	0.00334812058	NA	0					
rolliqr3_nation	float64	0.0000000063507	0.0018534009667175	0.0000000011000175	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.2500000000000000	0.635001400619301520234	NA	0						
rollmeam3_merit	float64	0.0005493	0.00000026017	0.0000000	-	0.0000000000000000	0.0000000000000000	0.330533802667	263404	200	NA	0	0.0005333	0.000533301333	0.3356897		
rollsd3_nmerit	float64	0.000000005774	0.00026460120002646	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.350205774001780	245260	200	NA	0	0.7656307				
rolliqr3_nmerit	float64	0.000000005400	0.00025020104002500	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.35000500002000	-	200	NA	0	0.0143798	8741			
rollmeam3_merit	float64	5.215038833962867	4.67734298206	6.3600000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3383306567667	4466303	200	NA	0	0.1897775				
rollsd3_nmerit	float64	0.00000000365726	0.268217398125	1.0640000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.359190307794291	45622497203	NA	0						
rolliqr3_nmerit	float64	0.00000000750650	0.2562682421234000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.350008500067104338	6200	NA	0						
rollmeam3_merit	float64	3.3043.302066590.4832885789.8956766758333383.336666663333336666700	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.099217	200	NA	0	0.0140140				
rollsd3_nmerit	float64	0.00000000362052B.54606786373636620000200743.3366203392373619314802522488270	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rolliqr3_nmerit	float64	0.00000000740000D3.38495.78773.180000000020010000800000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rollmeam3_merit	float64	3.3043.302066590.4832885789.8956766758333383.336666663333336666700	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.099217	200	NA	0	0.0140140				
rollsd3_nmerit	float64	0.00000000362052B.54606786373636620000200743.3366203392373619314802522488270	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rolliqr3_nmerit	float64	0.00000000740000D3.38495.78773.180000000020010000800000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rollmeam3_merit	float64	3.3043.302066590.4832885789.8956766758333383.336666663333336666700	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.099217	200	NA	0	0.0140140				
rollsd3_nmerit	float64	0.00000000362052B.54606786373636620000200743.3366203392373619314802522488270	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rolliqr3_nmerit	float64	0.00000000740000D3.38495.78773.180000000020010000800000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rollmeam3_merit	float64	3.3043.302066590.4832885789.8956766758333383.336666663333336666700	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.099217	200	NA	0	0.0140140				
rollsd3_nmerit	float64	0.00000000362052B.54606786373636620000200743.3366203392373619314802522488270	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rolliqr3_nmerit	float64	0.00000000740000D3.38495.78773.180000000020010000800000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rollmeam3_merit	float64	3.3043.302066590.4832885789.8956766758333383.336666663333336666700	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.099217	200	NA	0	0.0140140				
rollsd3_nmerit	float64	0.00000000362052B.54606786373636620000200743.3366203392373619314802522488270	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rolliqr3_nmerit	float64	0.00000000740000D3.38495.78773.180000000020010000800000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rollmeam3_merit	float64	3.3043.302066590.4832885789.8956766758333383.336666663333336666700	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.099217	200	NA	0	0.0140140				
rollsd3_nmerit	float64	0.00000000362052B.54606786373636620000200743.3366203392373619314802522488270	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rolliqr3_nmerit	float64	0.00000000740000D3.38495.78773.180000000020010000800000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rollmeam3_merit	float64	3.3043.302066590.4832885789.8956766758333383.336666663333336666700	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.099217	200	NA	0	0.0140140				
rollsd3_nmerit	float64	0.00000000362052B.54606786373636620000200743.3366203392373619314802522488270	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rolliqr3_nmerit	float64	0.00000000740000D3.38495.78773.180000000020010000800000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rollmeam3_merit	float64	3.3043.302066590.4832885789.8956766758333383.336666663333336666700	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.099217	200	NA	0	0.0140140				
rollsd3_nmerit	float64	0.00000000362052B.54606786373636620000200743.3366203392373619314802522488270	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rolliqr3_nmerit	float64	0.00000000740000D3.38495.78773.180000000020010000800000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rollmeam3_merit	float64	3.3043.302066590.4832885789.8956766758333383.336666663333336666700	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.099217	200	NA	0	0.0140140				
rollsd3_nmerit	float64	0.00000000362052B.54606786373636620000200743.3366203392373619314802522488270	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rolliqr3_nmerit	float64	0.00000000740000D3.38495.78773.180000000020010000800000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rollmeam3_merit	float64	3.3043.302066590.4832885789.8956766758333383.336666663333336666700	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.099217	200	NA	0	0.0140140				
rollsd3_nmerit	float64	0.00000000362052B.54606786373636620000200743.3366203392373619314802522488270	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rolliqr3_nmerit	float64	0.00000000740000D3.38495.78773.180000000020010000800000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.64662263286	NA	0						
rollmeam3_merit	float64	3.3043.302066590.4832885789.8956766758333383.336666663333336666700	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.3500000000000000	0.099217	200	NA	0	0.0140140				
rollsd3_nmerit	float64	0.00000000362052B.54606786373636620000200743.3366203392373619314802522488270	0.0000000000000000	0.0000000000000000	0.0000000000000000	0.000000000											

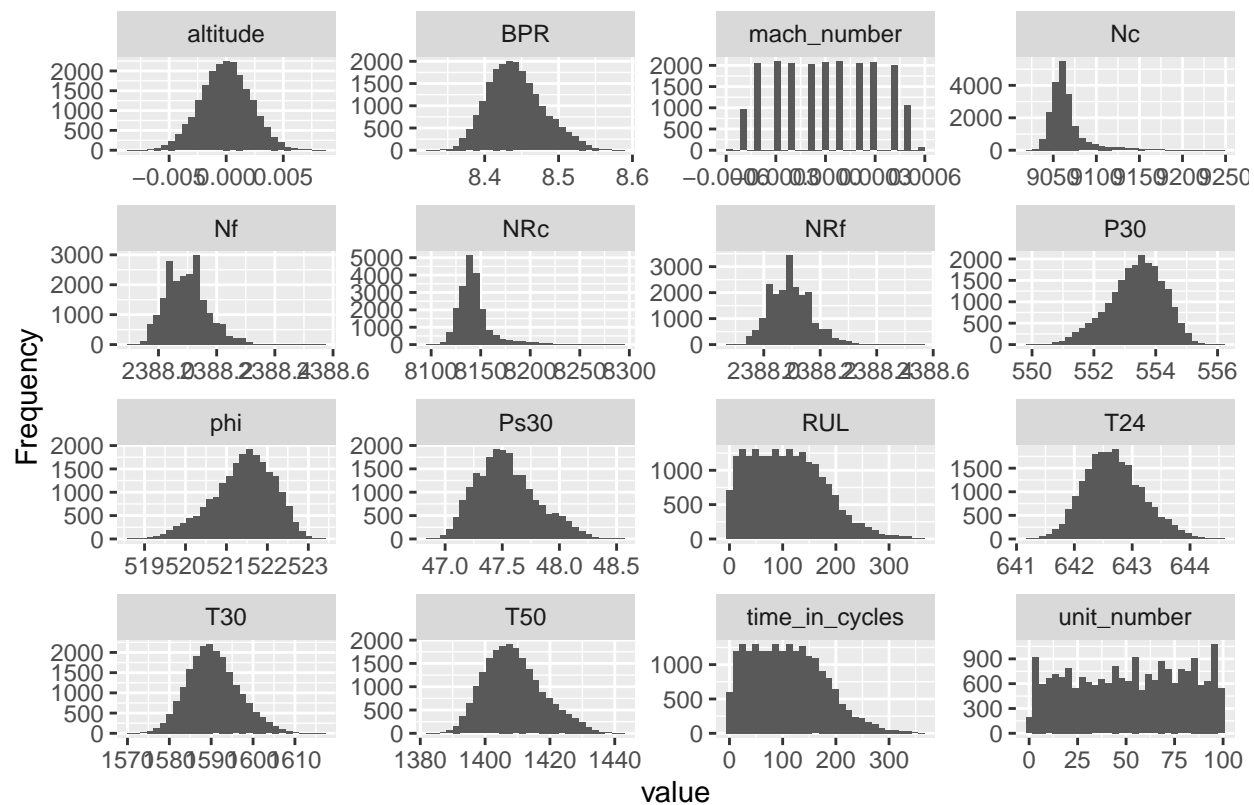
[illegible]

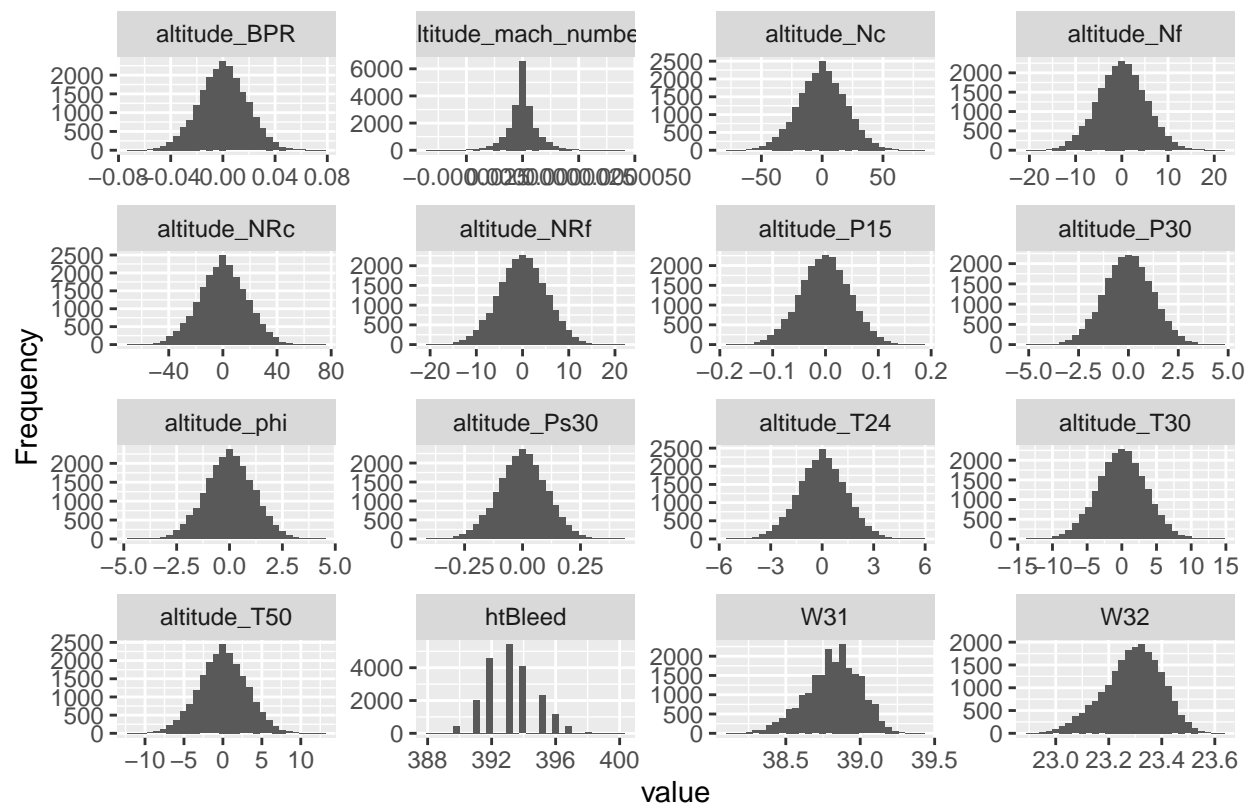
[illegible]

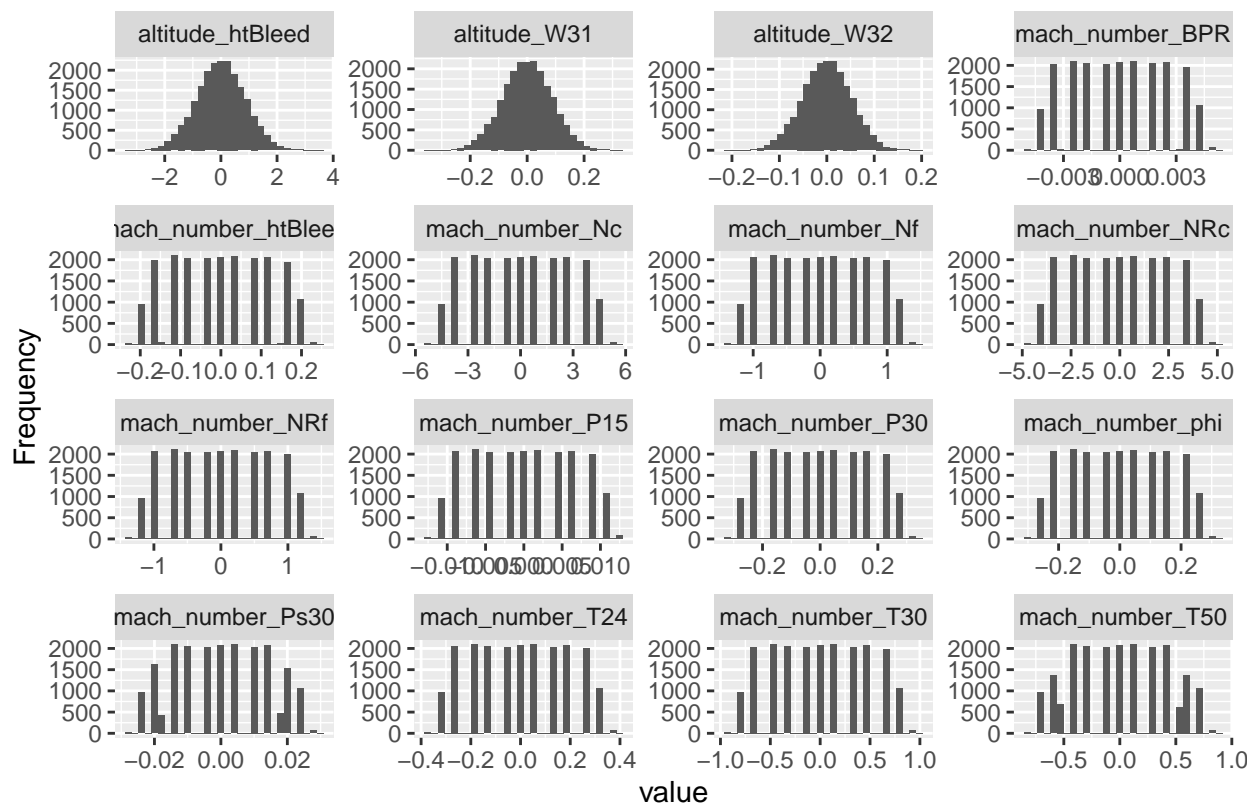
[illegible]

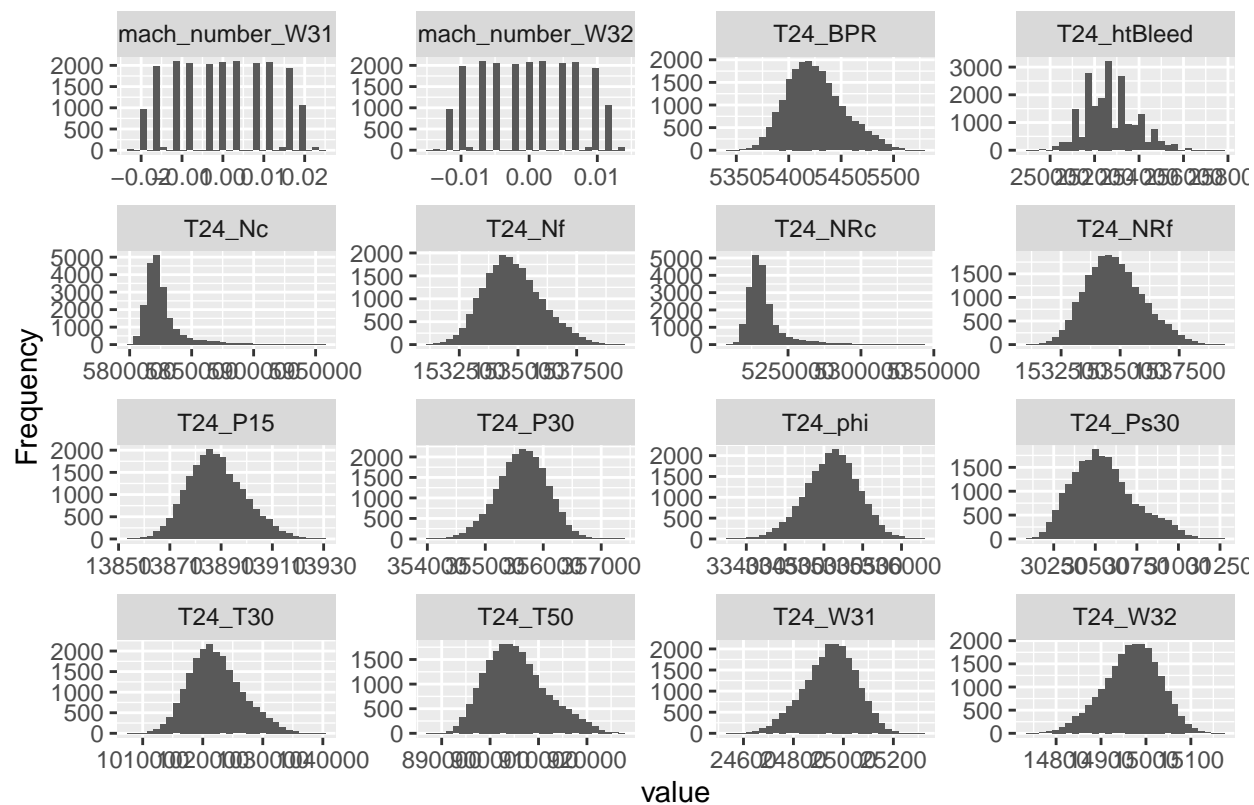
names	class	min	max	n_unique	mean	sd	median	quantile	skewness	kurtosis	levene	shapiro	jarque	bern	mc	missing	var	outliers
cumsum	NR	187.98	60000950	50840.64	4973.18	352357	90000702	42352005080400	20866000098976	0	0	0						
									0.2190301									
cummin	NR	187.83	880070000	2087.98	847024	177.92	87083870923870	2688023880070900000	0876883	0	0	0						
									0.9052250									
cummax	NR	187.92	880060000	2088.15	2659268	123870923880026880	2388026880268800	00900000	35550659047	0	4							
cumsum	NR	16.82	400007200	60705.23	8958.84	483568000270843040	2706200505492830000	007856	0	0	0							
									0.2185192									
cummin	NR	99.98	1008020080	29.82	76448831.08	000981260900008000	220680003000000	-	0	0	0							
									0.30953465097									
cummax	NR	16.88	2007070680	50.72	97384459.42	16880400800800	4805048200702034000	702502353929	0	721								
cumsum	BR	3030055.72	300006.91	83.57823	768.7300308	12875079370858755.78	8138005002753	0	0	0								
									0.2170711									
cummin	BR	24904875087	8.3731602	0586731802	498358080731808818	48750030100919331	0	0	3									
									0.1382855									
cummax	BR	30300648870	8.4801903	03848793803	0384573847938004808480043100	2170000	0	0	0									
									0.1640209									
cumsum	ht	100002020.54	802709.27	9611.85	532100000404	507300008085026200084000000099	0	0	0									
									0.2169563									
cummin	ht	1800060000000	90.2208372388	0088800000000000000000000000000000000000000	1362824090	0	2											
cummax	ht	1000000000000	94.9252233318	0080000004000000000000000000000000000000000	4722871082	0	0											
cumsum	W3	6100048.0000001.32	17252436.88	000202640060008098004804560200980709	0	0	0											
									0.2219665									
cummin	W3	14000270000	038.6255743736	4308004088040880308807208007000800000	0.0324637	0	1											
									0.0934173									
cummax	W3	61000430680039.15	670948355.508800	0800908015080230804300000000	0.2365971	0	3											
									0.0677189									
cumsum	W3	20187462.21070	0638.8425840522053	50010706.3282006940585320247402068983058	0	0	0											
									0.2211579									
cummin	W3	2894205349000	23.17776852504812	20942002120081200432053490002000	-	0	0	0										
									0.20791025366									
cummax	W3	20120008420	023.49385635834952	00120449204952053520008010859000	-	0	0	0										
									0.19342404125									

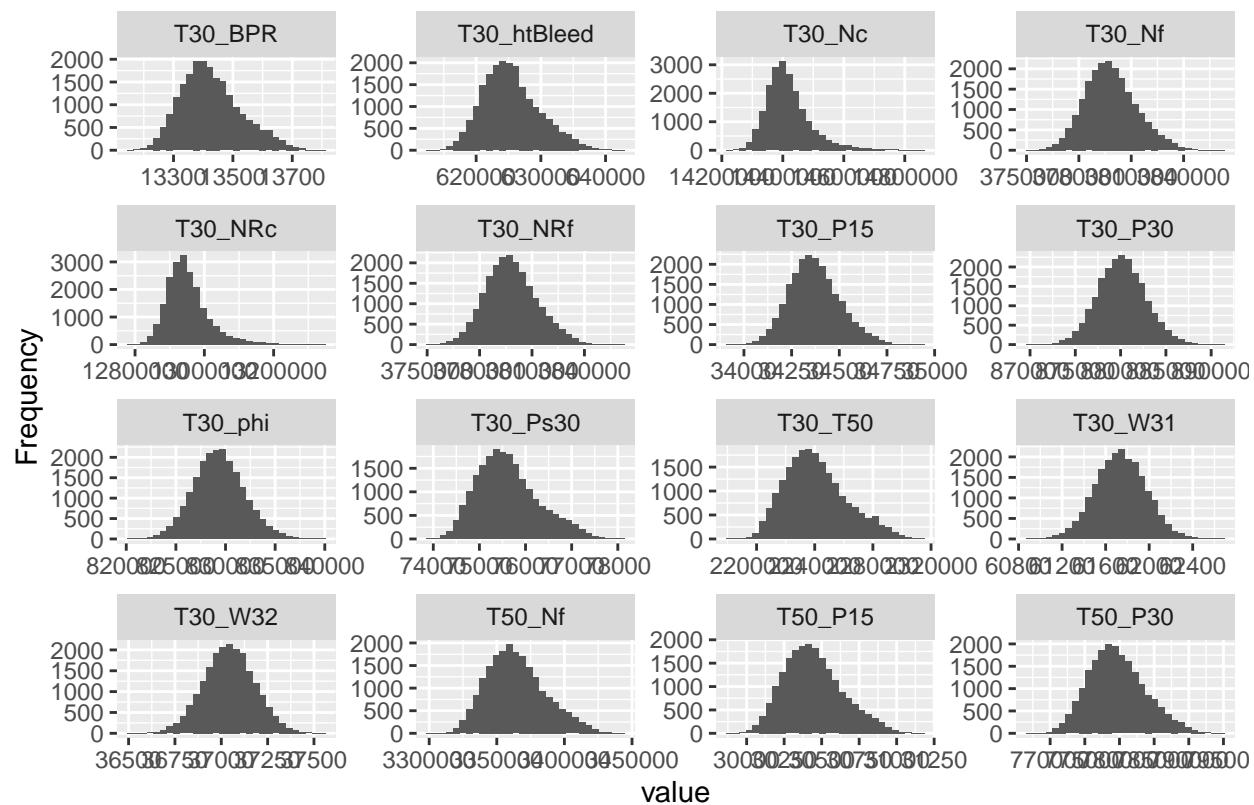
```
plot_histogram(trainset2)
```

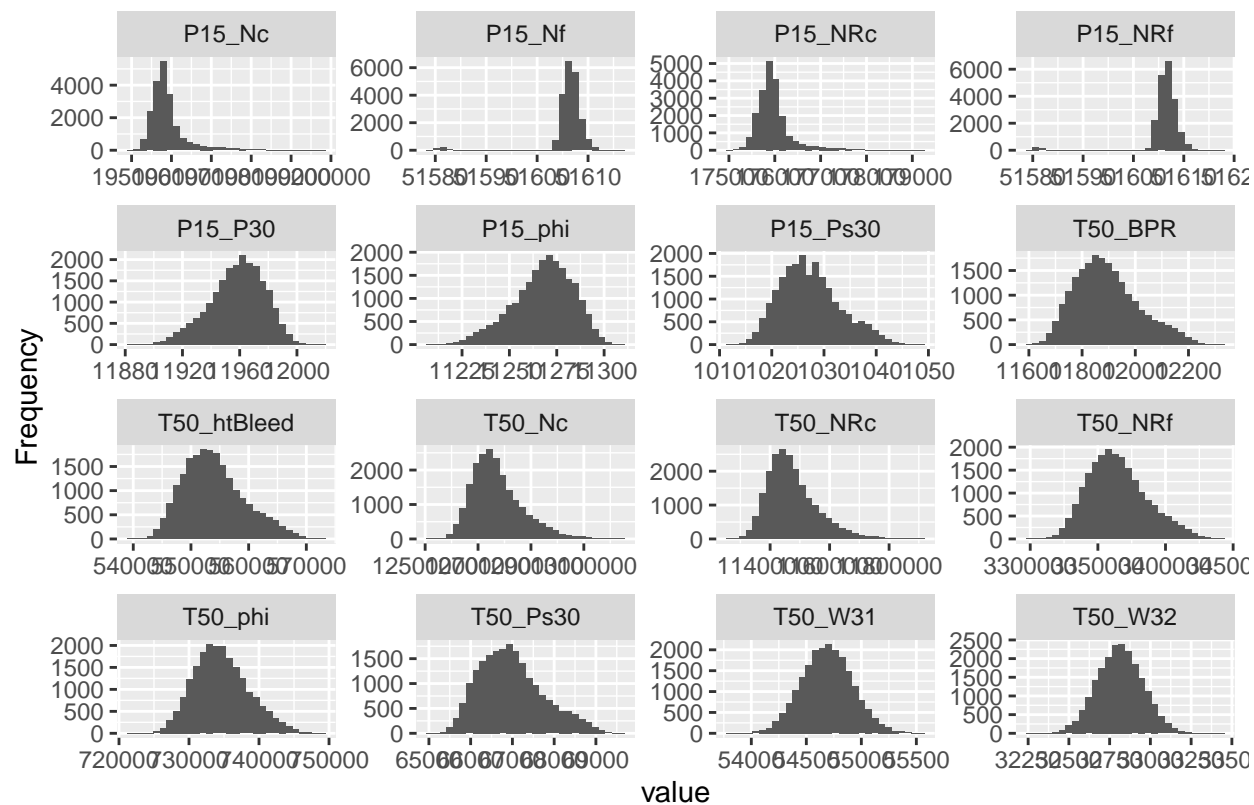


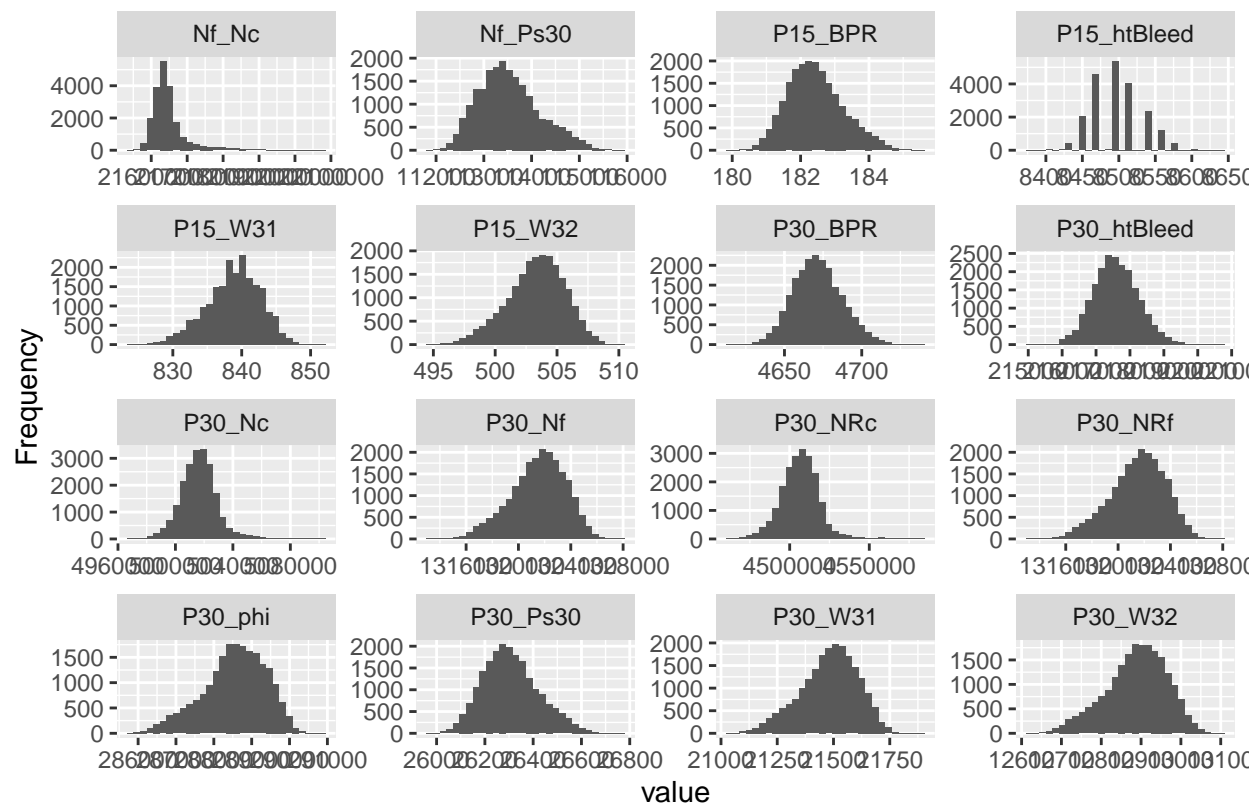


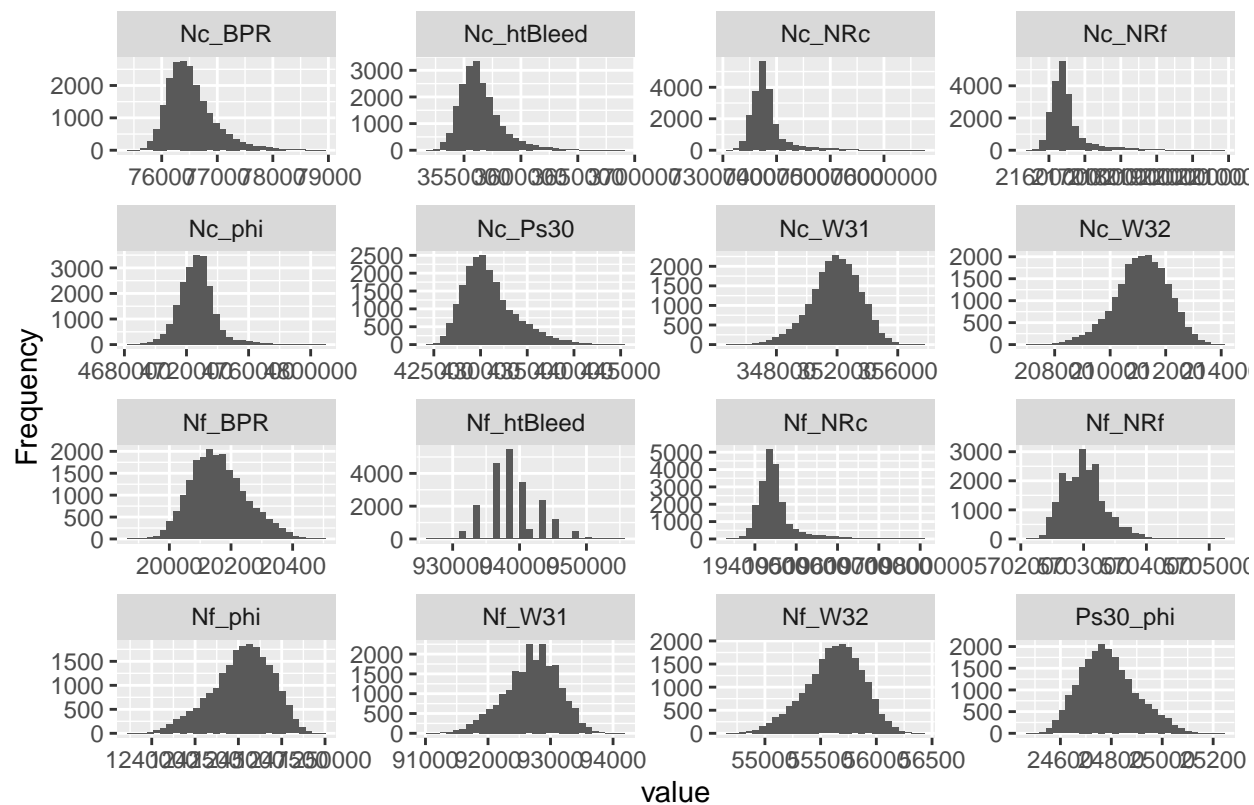


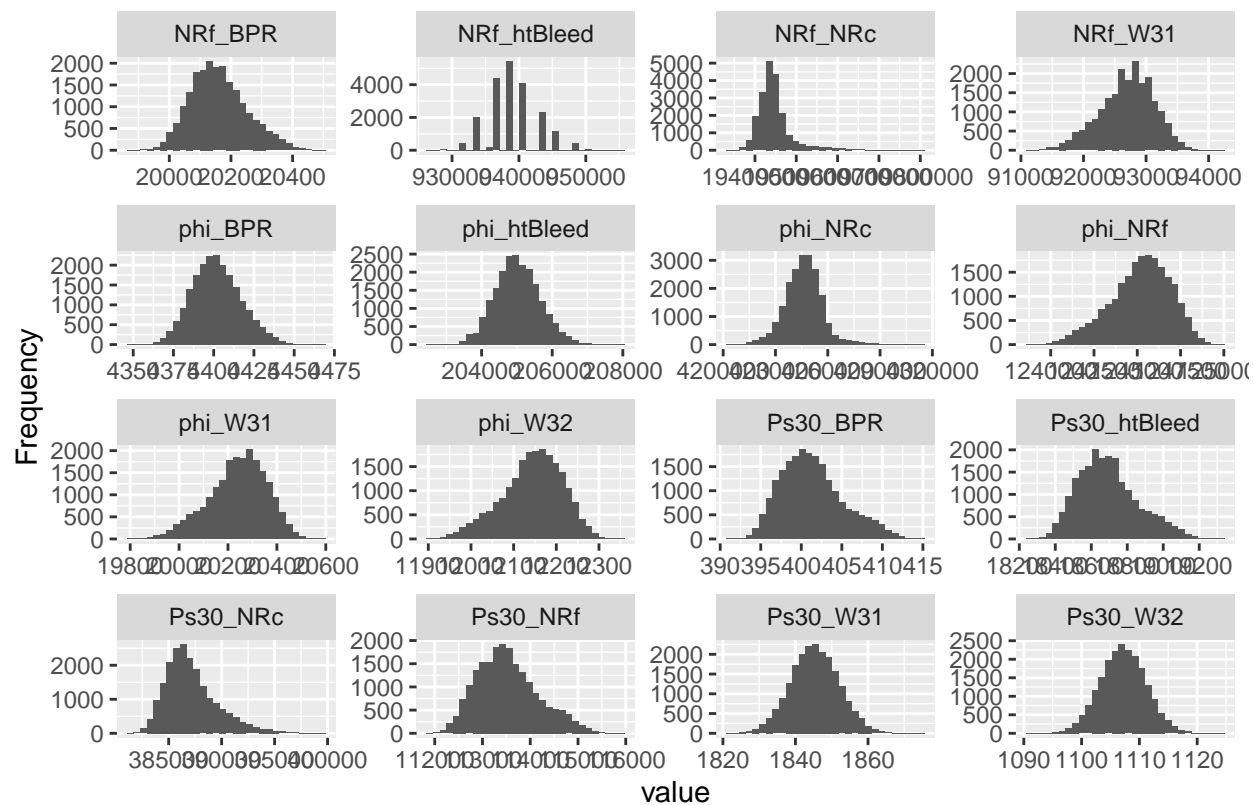


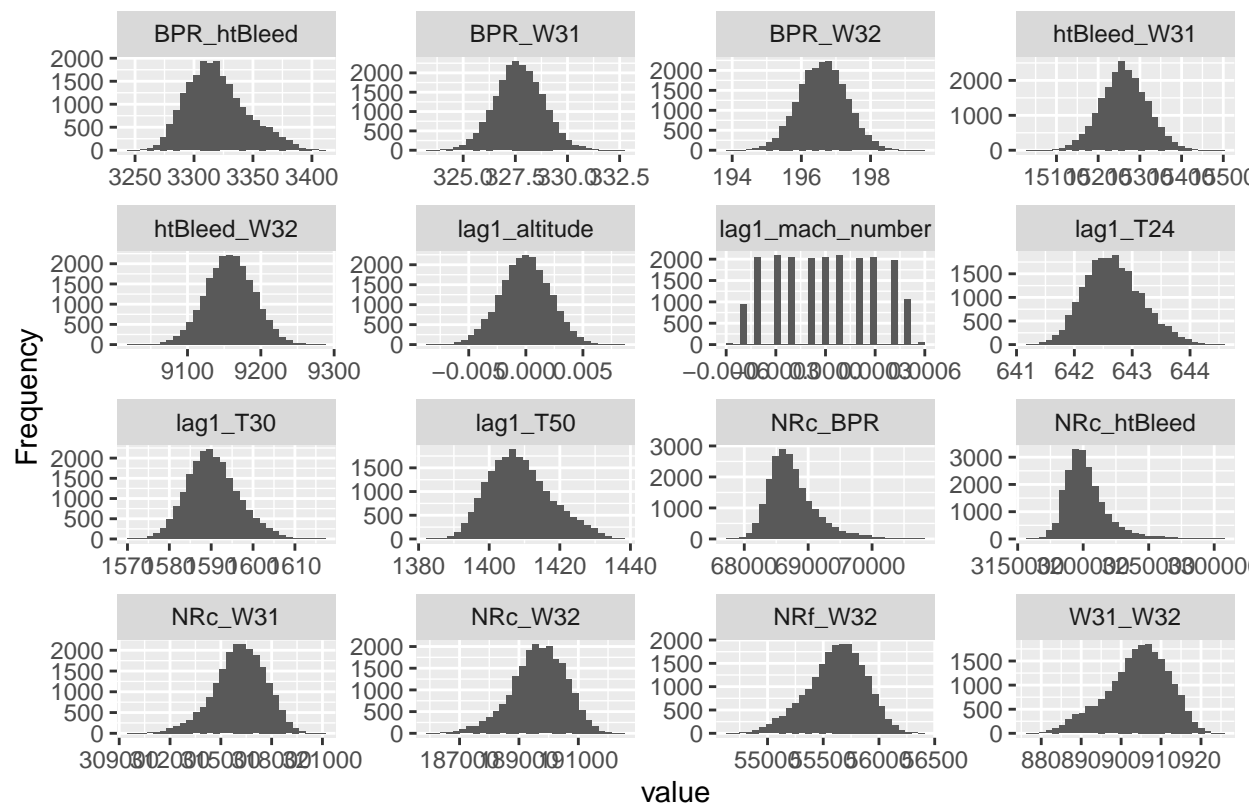


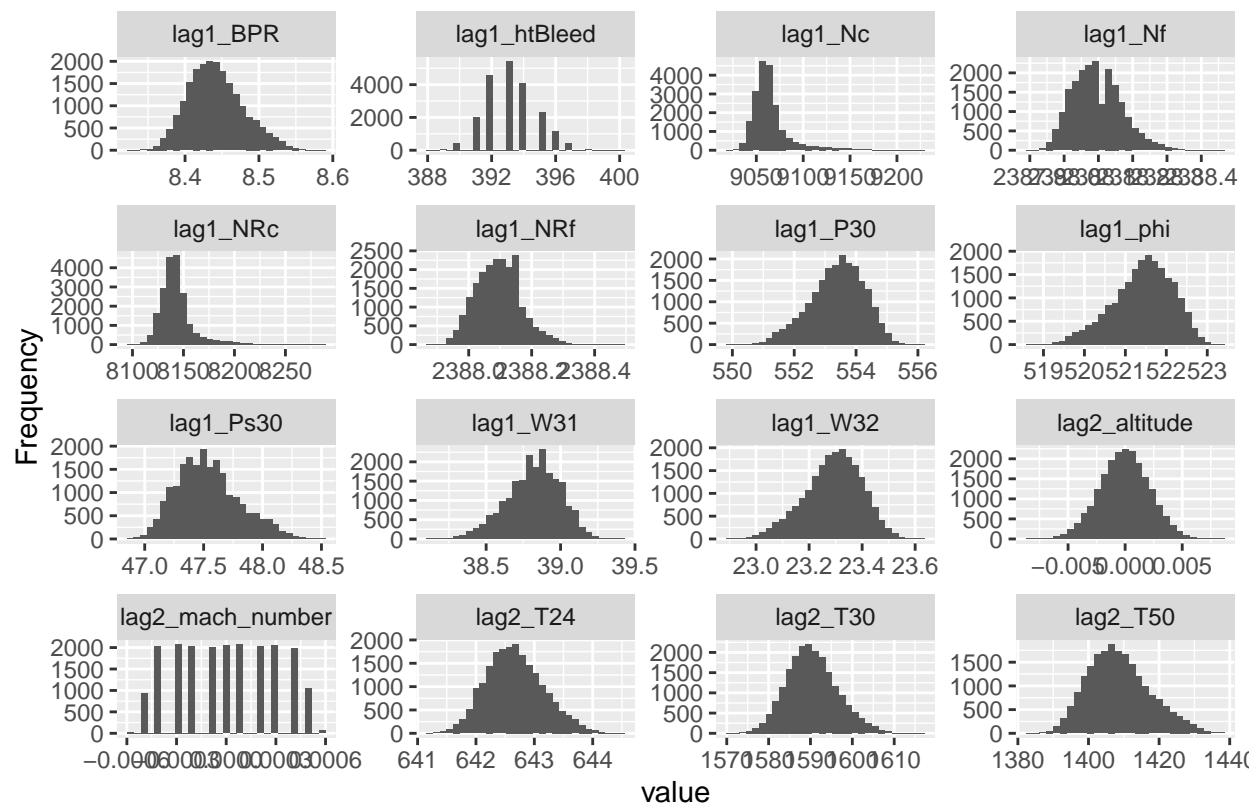


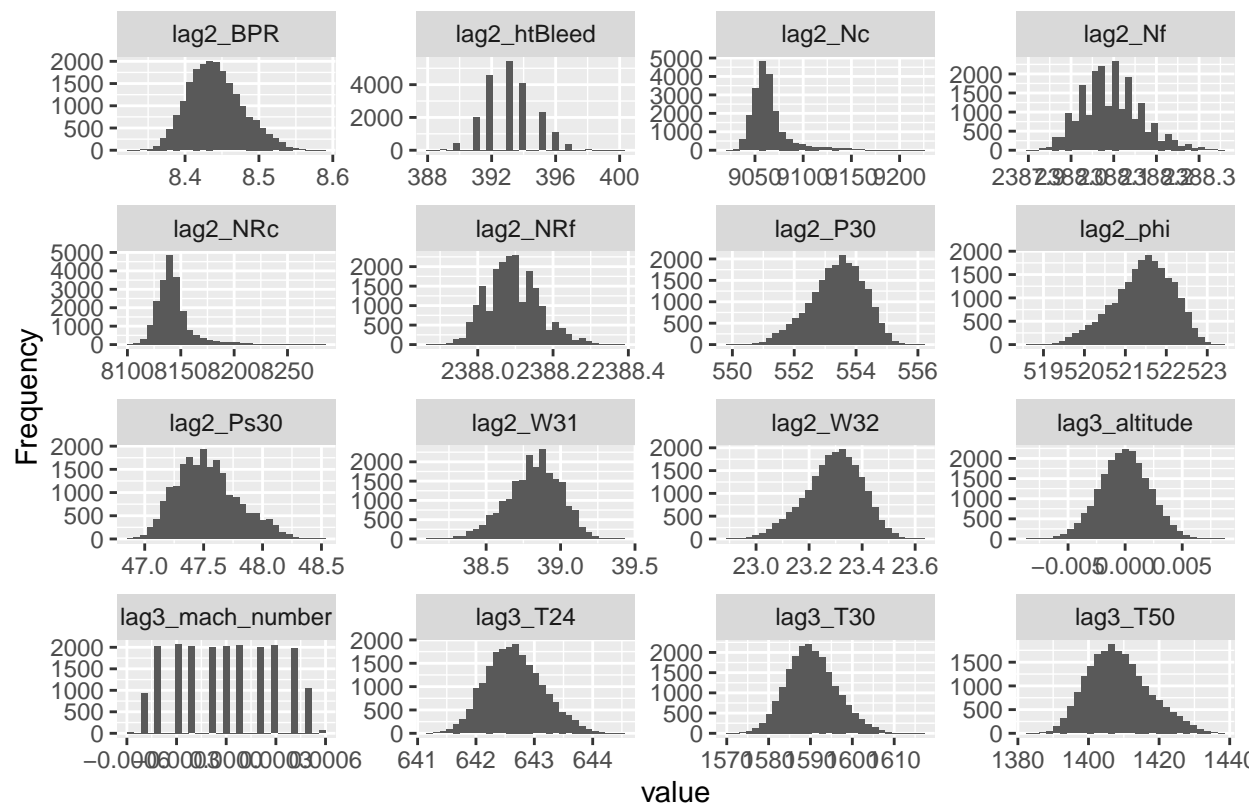


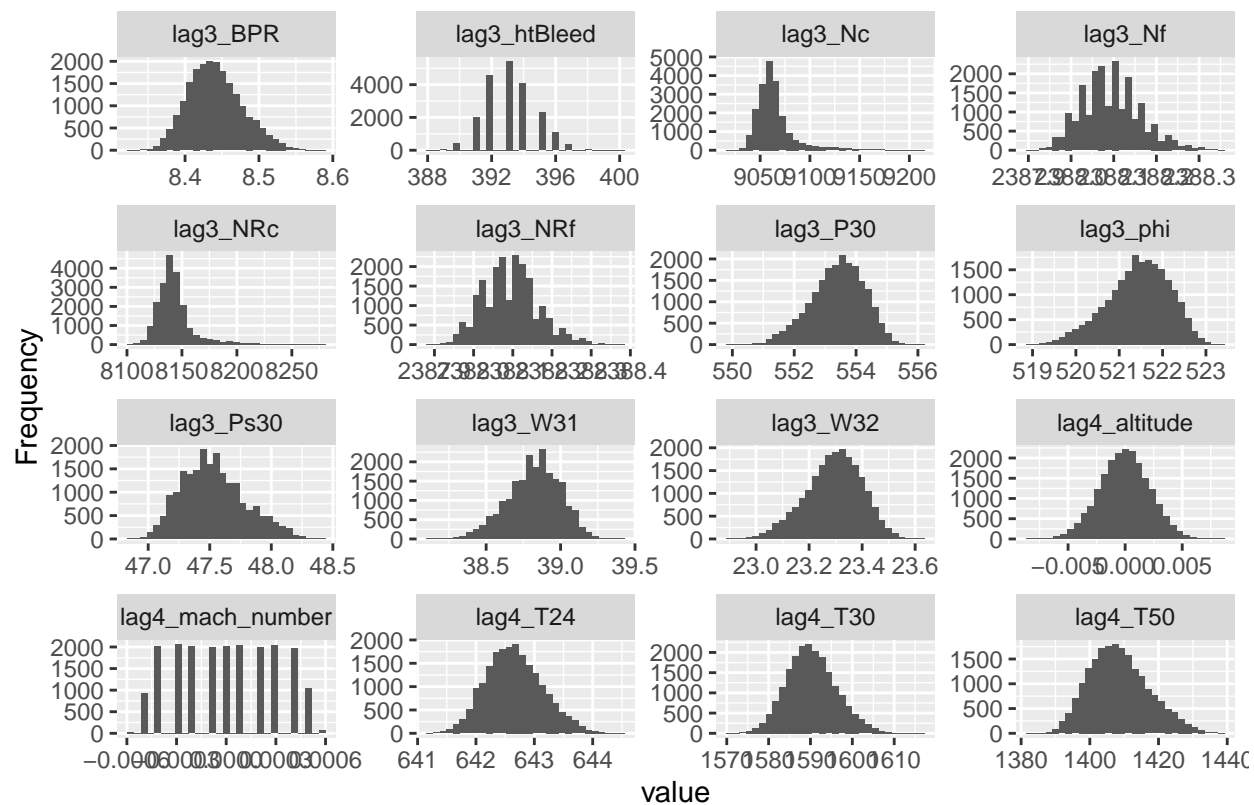


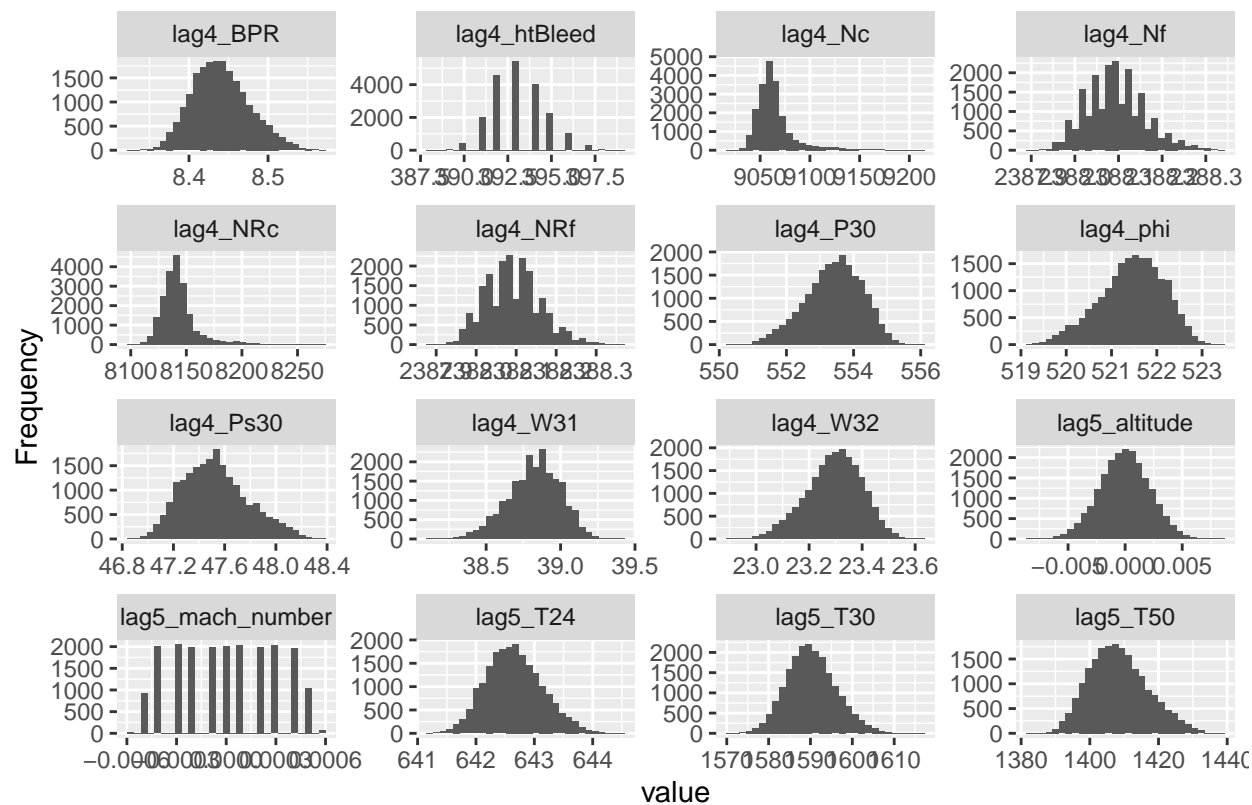


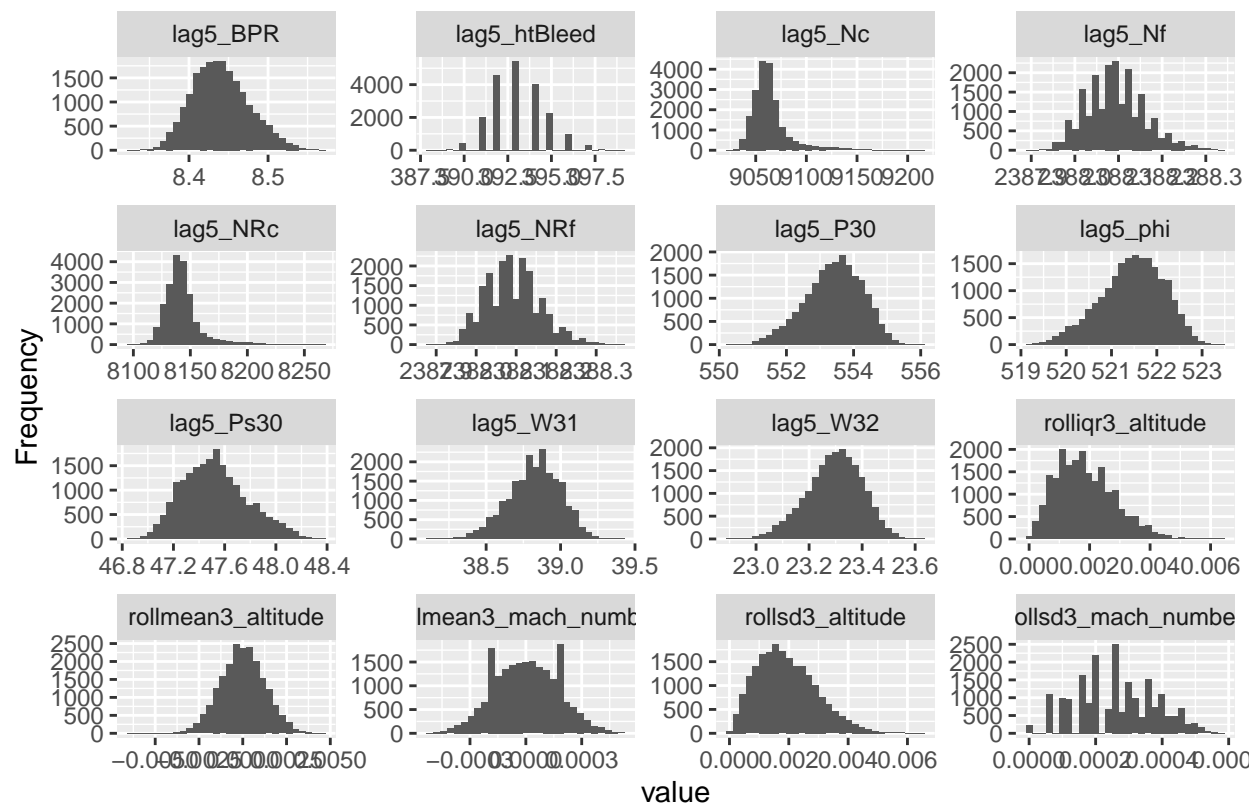


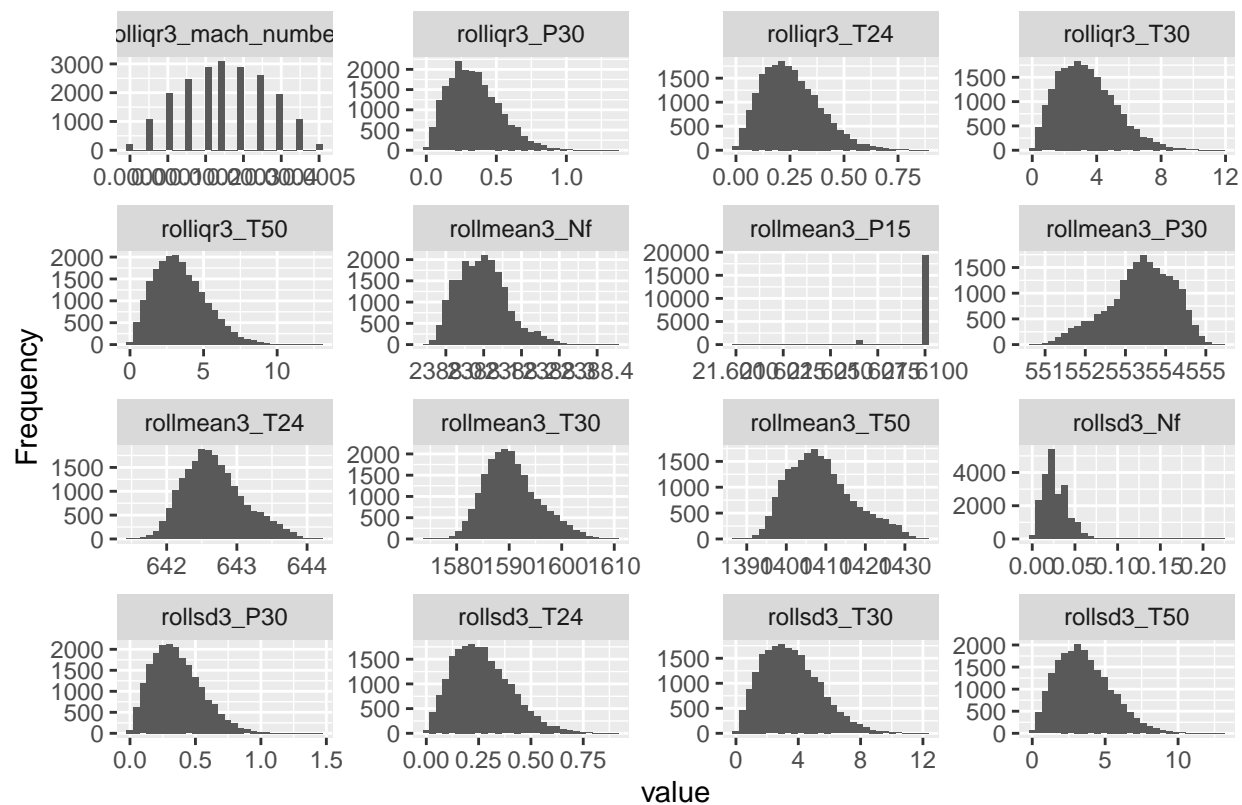


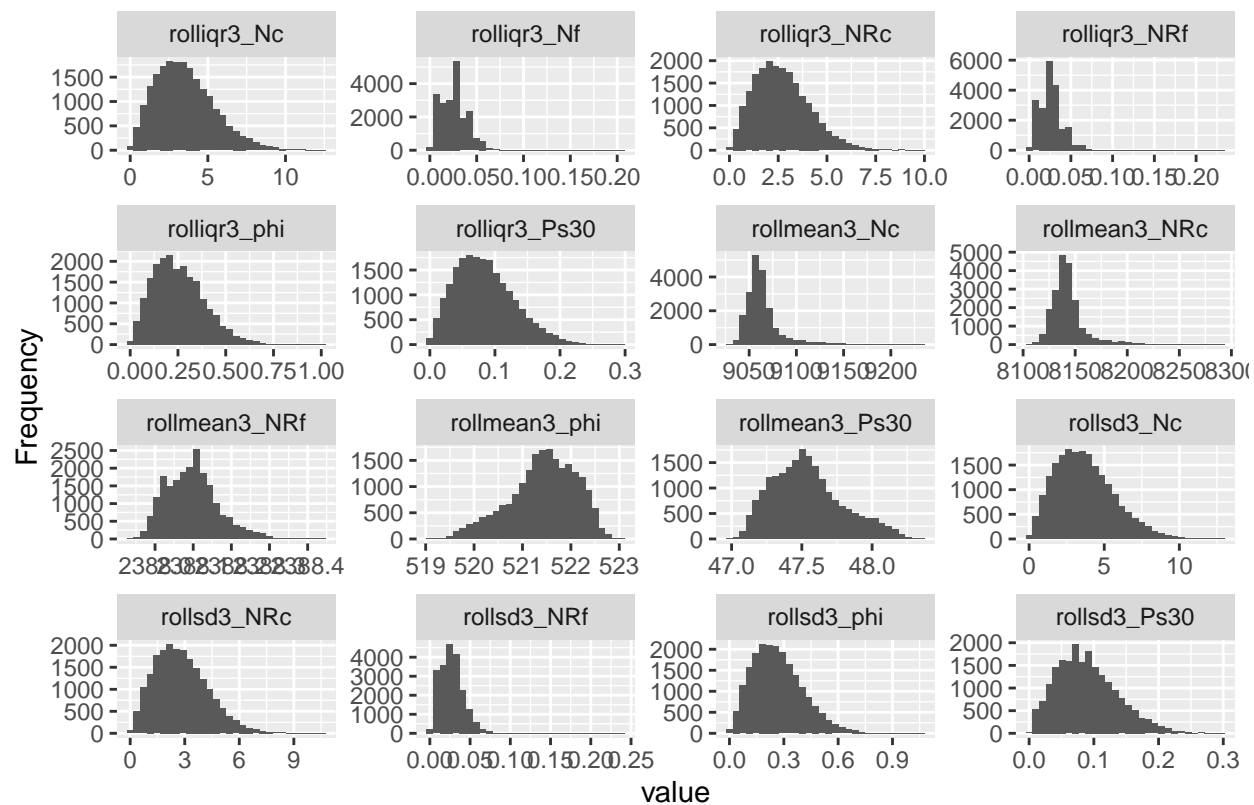


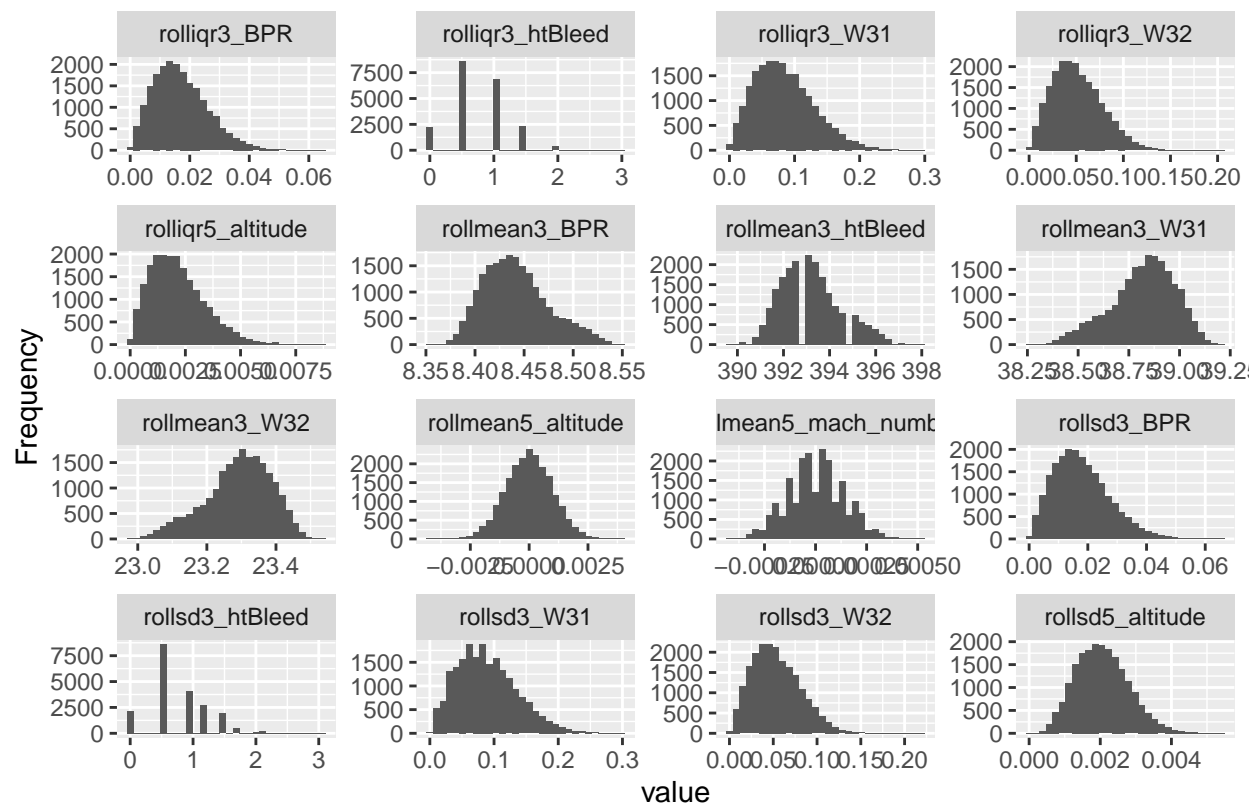


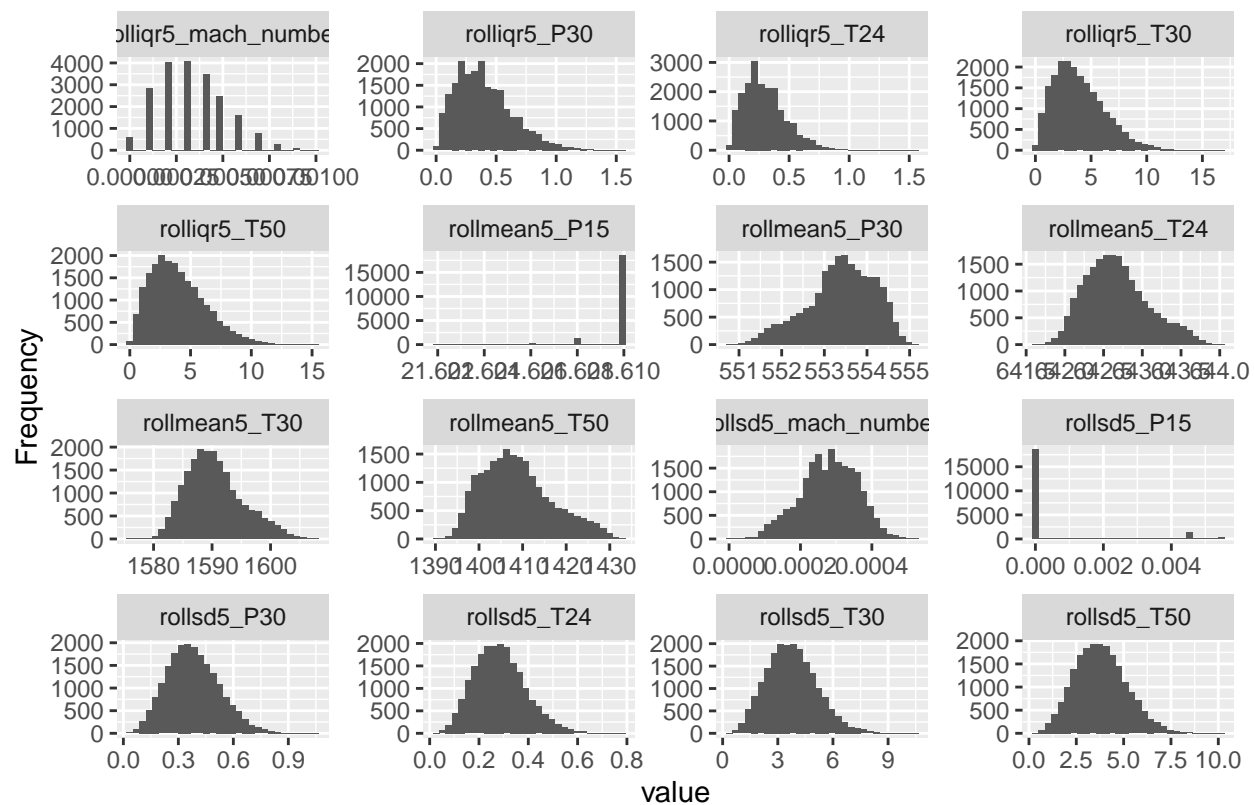


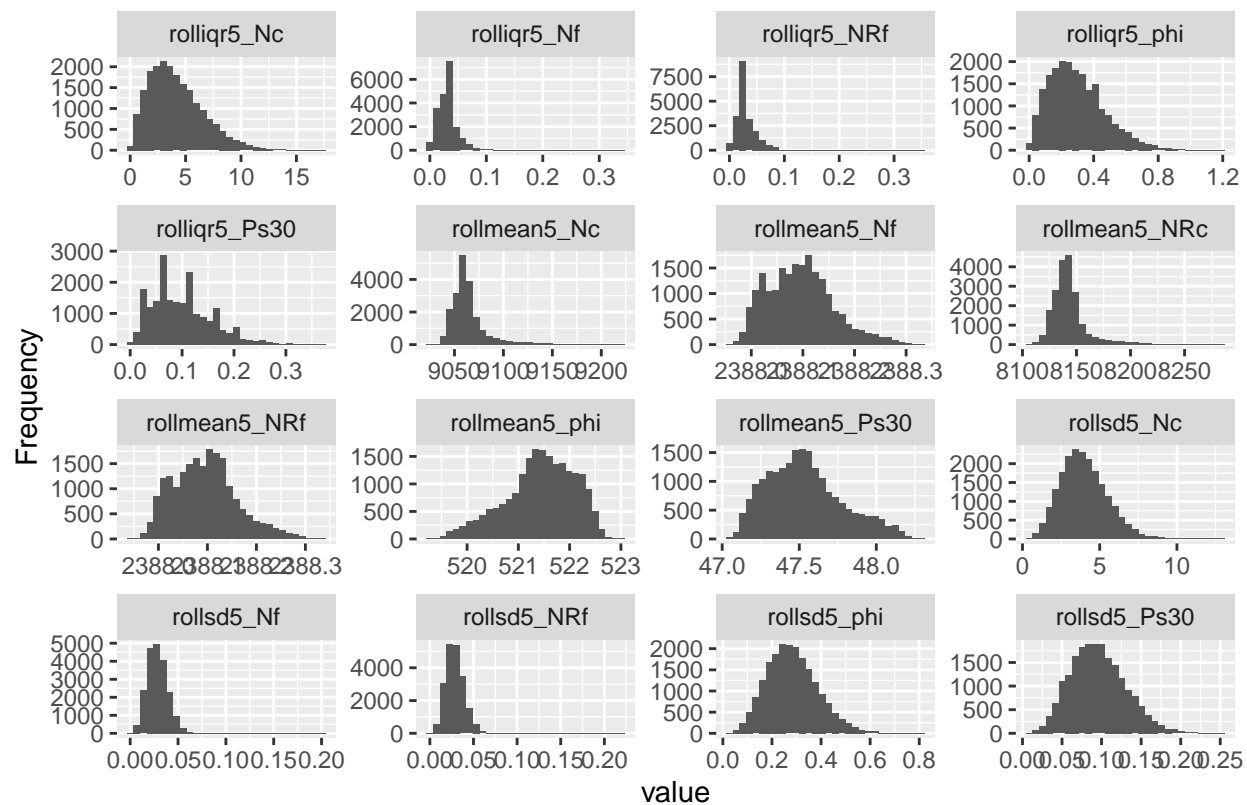


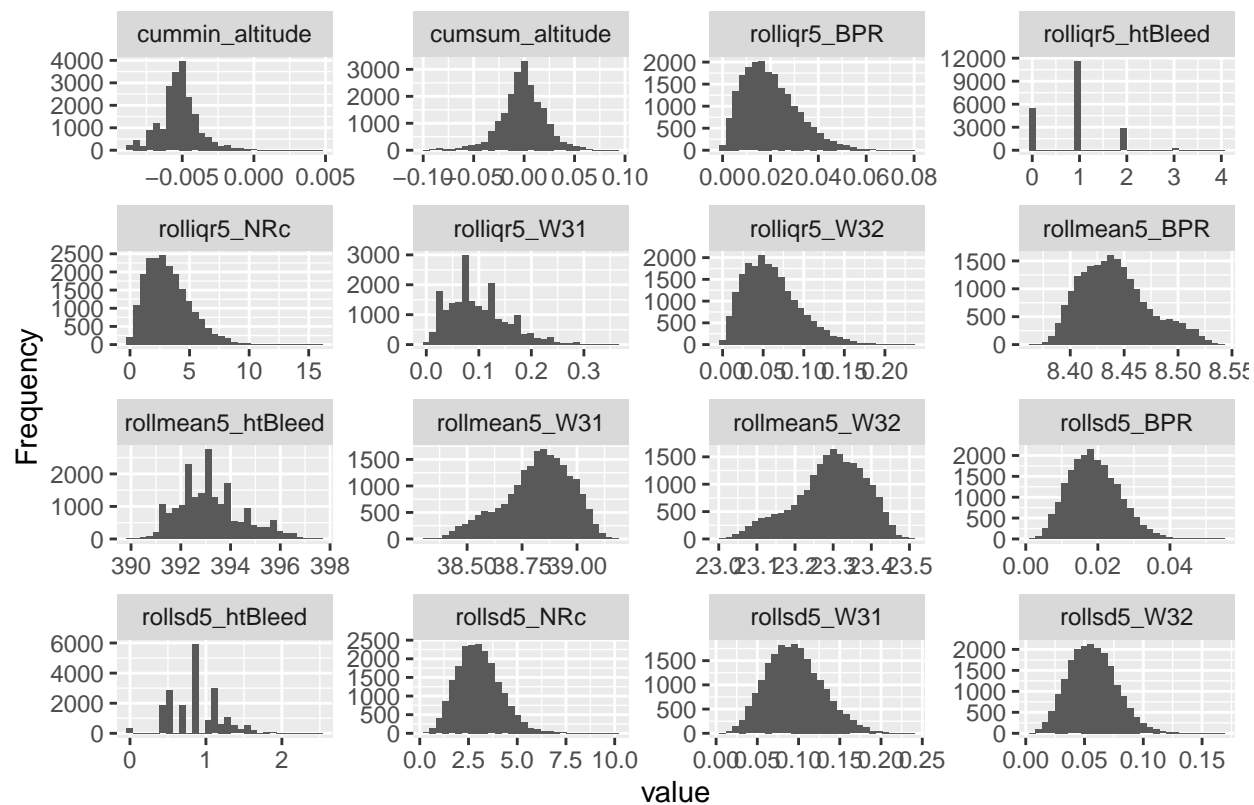


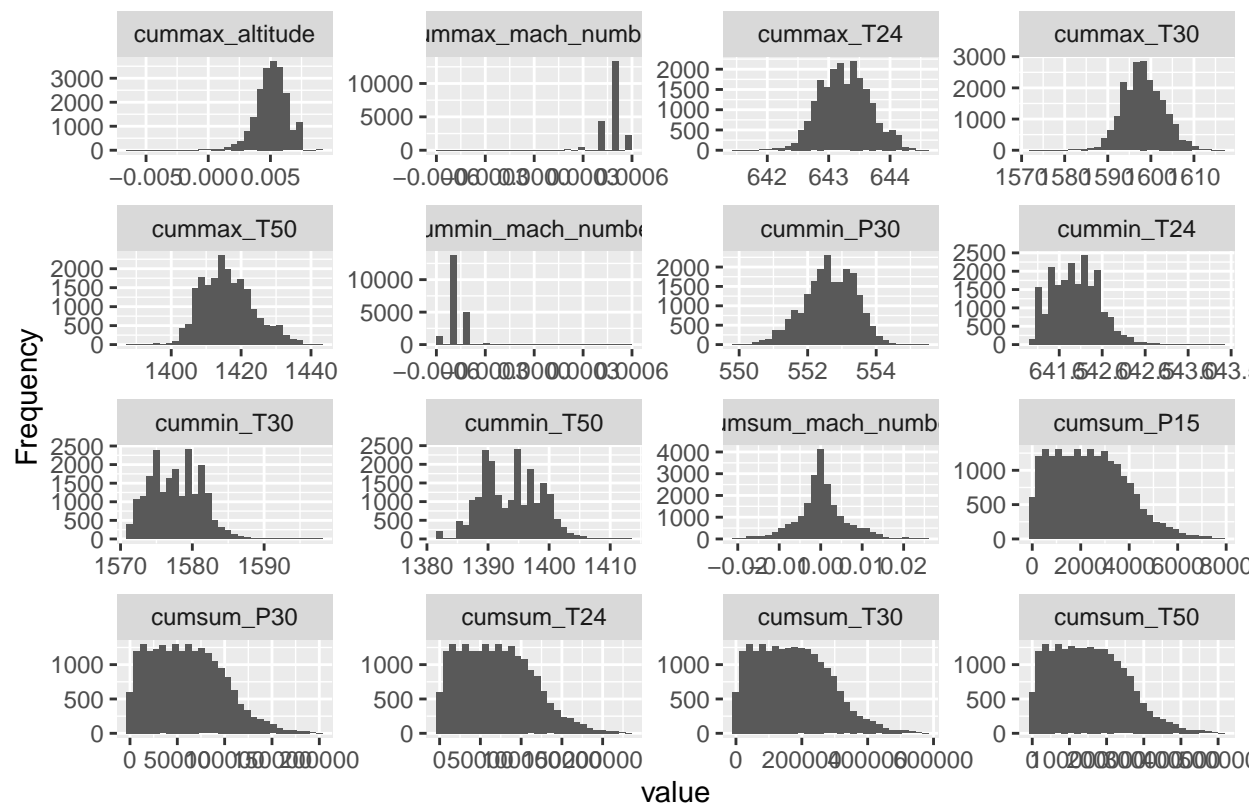


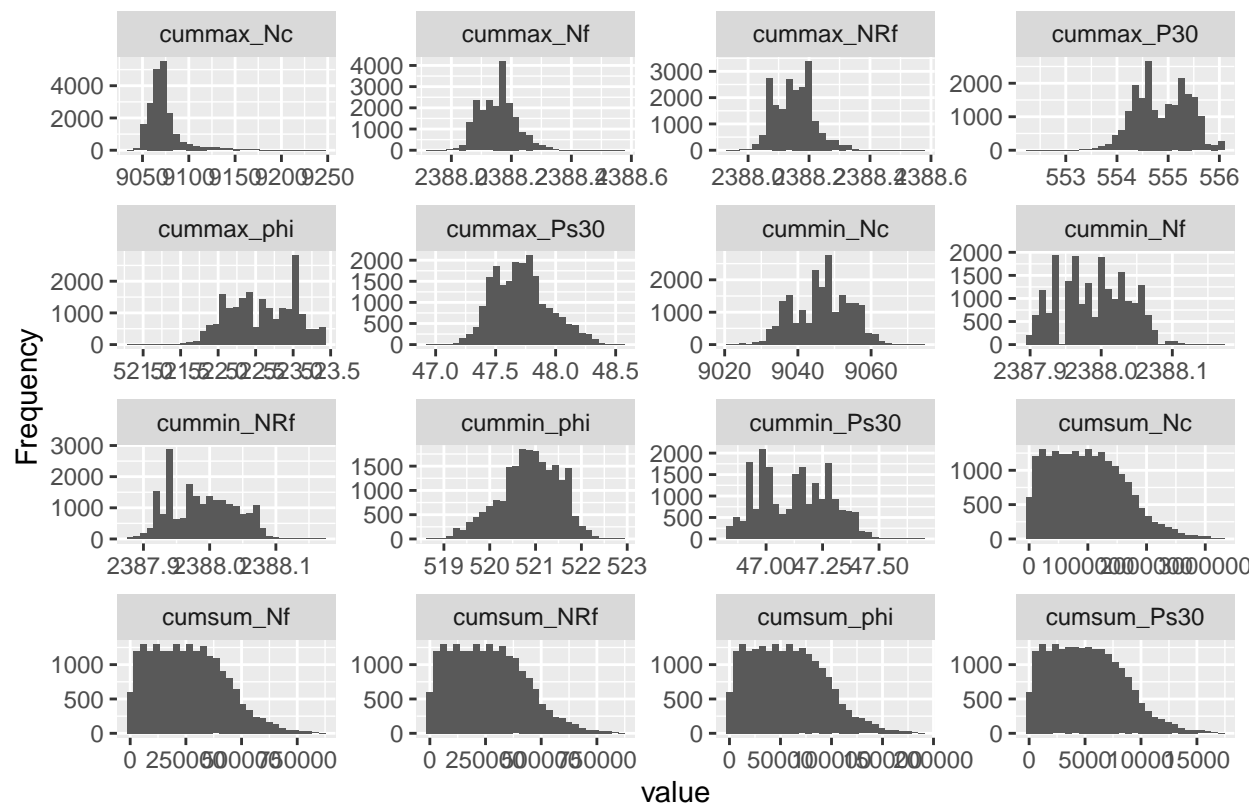


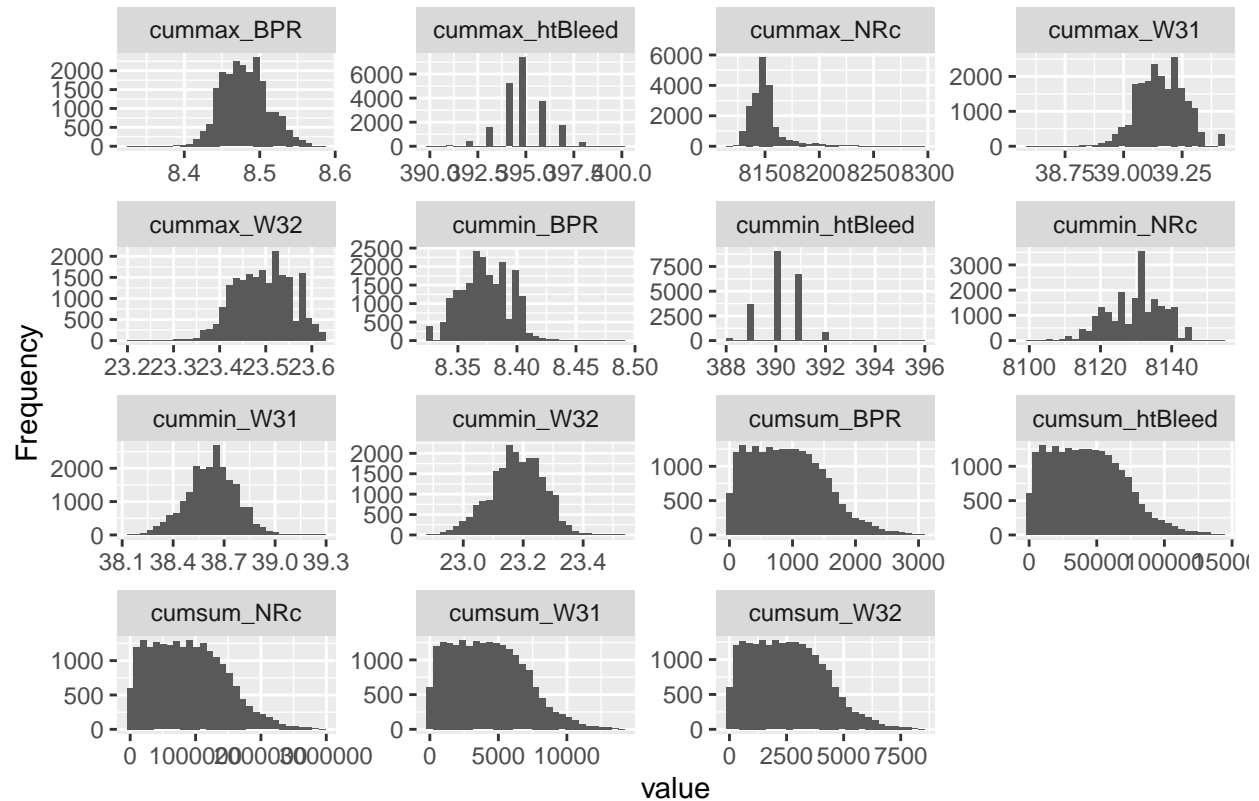












3. Limpieza de los datos

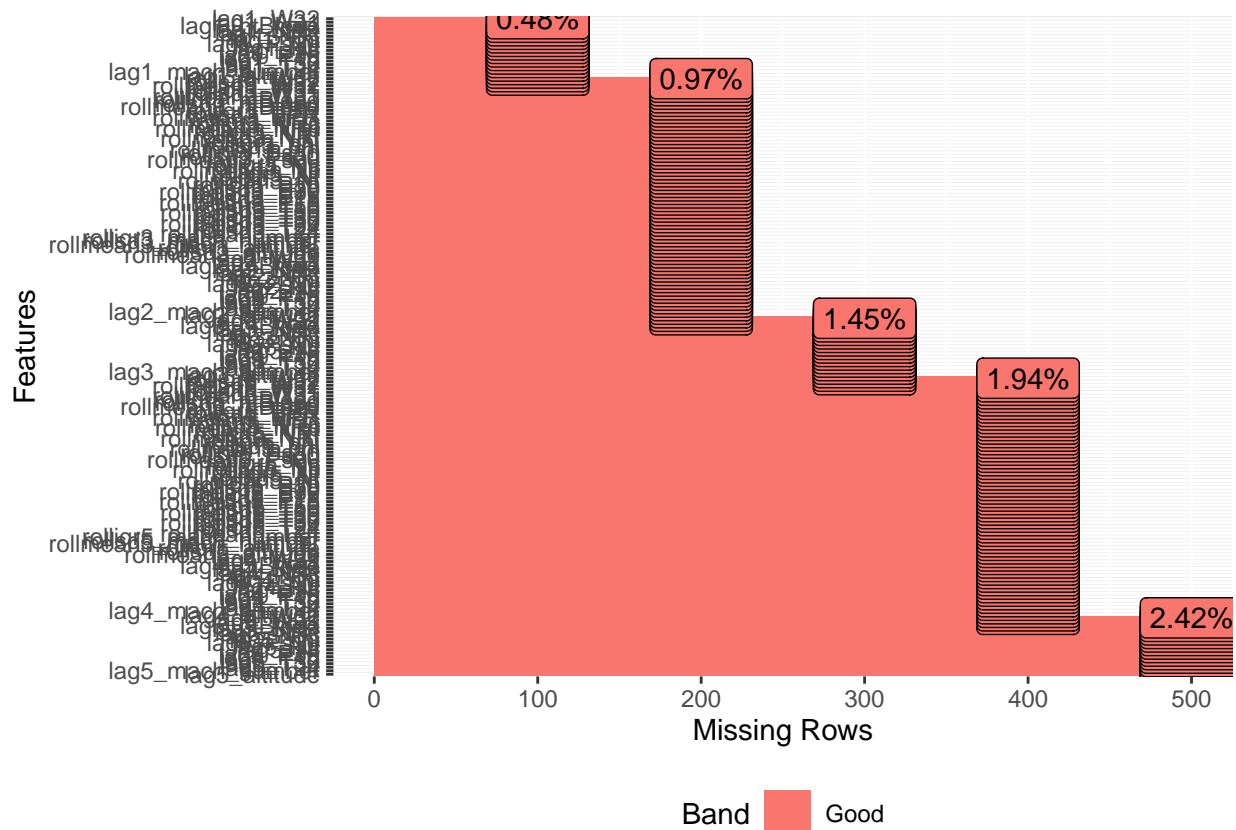
3.2. Identificación y tratamiento de valores extremos.

Nuestro dataset no tiene Missing values, pero hay valores extremos. Para

3.1. Missing Values o Ceros

Como tenemos muchas variables, filtramos las variables sin na y generamos un gráfica para conocer las columnas con NA y sus patrones.

```
columns_without_na <- colnames(trainset2[, sapply(trainset2, Negate(anyNA)), drop = FALSE])
df_columns_with_na <- trainset2[, -which(names(trainset2) %in% columns_without_na)]
plot_missing(df_columns_with_na)
```



Como se esperaba, las columnas con lag y roll contienen NAs. Pero esto tiene que ser hasta un `time_in_cycles` máximo de 5. Vamos a confirmar esto:

```
df_columns_with_na$time_in_cycles <- trainset2$time_in_cycles

borra <- df_columns_with_na %>% filter_all(any_vars(is.na(.)))

table(borra$time_in_cycles)
```

```
##
##  1  2  3  4  5
## 100 100 100 100 100
```

Merece la pena hacer algún tipo de imputación de valor?

Para contestar a esta pregunta, vamos a ver si hay alguna máquina en nuestros datos que tenga fallado en este período:

```
trainset2 %>% filter(RUL == 0) %>% pull(time_in_cycles) %>% sort() %>% unique()
```

```
## [1] 128 135 137 147 150 153 154 155 156 158 163 165 166 168 170 172 174 178 179
## [20] 180 181 185 188 189 191 192 193 194 195 196 198 199 200 201 202 207 208 209
## [39] 210 213 214 215 216 217 222 229 230 231 234 240 256 257 258 259 267 269 275
## [58] 276 278 283 287 293 313 336 341 362
```

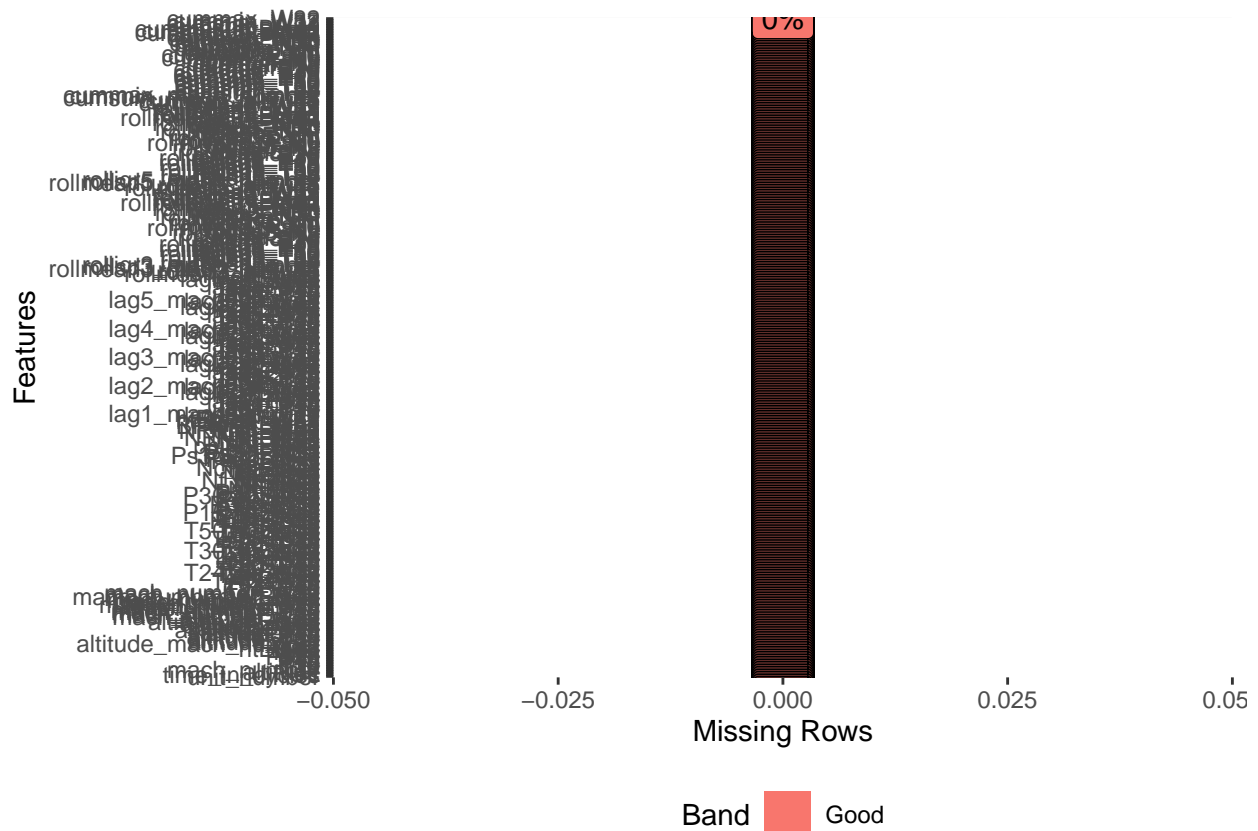
```
testset2 %>%
  group_by(unit_number) %>%
  filter(time_in_cycles == max(time_in_cycles)) %>%
  pull(time_in_cycles) %>%
  sort() %>%
  unique()
```

```
## [1] 31 34 37 39 46 48 49 50 54 55 56 68 71 72 73 74 76 78 83
## [20] 88 89 94 97 98 101 105 106 110 112 113 121 123 125 126 130 131 133 134
## [39] 135 136 137 140 143 144 145 146 147 148 150 152 155 156 158 159 160 162 164
## [58] 165 166 168 171 172 176 177 184 186 187 189 192 195 196 198 203 205 213 217
## [77] 232 234 244 303
```

Como se puede ver no existe ninguna necesidad de usar las observaciones de 1 a 5. Por lo tanto, vamos a eliminar estas observaciones con NA sin perder ninguna información importante para nuestro proyecto.

```
trainset2 <- trainset2 %>% na.omit()
testset2 <- testset2 %>% na.omit()

plot_missing(trainset2)
```



3.2. Identificación y tratamiento de valores extremos

El tratamiento de valores extremos es importante porque estos pueden sesgar / cambiar drásticamente las estimaciones y predicciones de ajuste. Estos datos se encuentran muy alejados de la distribución normal de

Nuestros datos son basados en simulación de sensores. Esto significa a priori que no hay sospechas de que las observaciones que se desvían tanto del resto fueron generadas mediante un mecanismo distinto.

Nc numeric NRc numeric P15 numeric Nf numeric NRf numeric T30 numeric

```
plot_missing(temp)
```



```
columns = c("altitude", "mach_number", "T24", "T30", "T50", "P15", "P30", "Nf", "Nc", "Ps30", "phi", "N")
```

```
OutlierAnomalyDetection <- function(df, columns = "all", method = "all"){
  if(any(class(df) == "grouped_df")){
    df <- df %>% ungroup()
  }
  if(any(columns == "all")){
    df_temp <- df
  } else {
    df_temp <- df %>% dplyr::select(one_of(columns))
  }

  if(method %in% c("all", "cook")){
    mod <- lm(RUL ~., data = df)
    cooks_d <- cooks.distance(mod)
    #influential <- as.numeric(names(cooks_d)[(cooks_d > (4/sample_size))])
    #df <- df[-influential, ]
    df$cook_score <- as.vector(cooks_d)
  }
  if(method %in% c("all", "mahalanobis")){
    score <- as.vector(mahalanobis(df_temp, colMeans(df_temp), cov(df_temp)))
    df$mahalanobis_score <- score
  }
  if(method %in% c("all", "PCA")){
    # Cálculo de PCA
    pca <- prcomp(
      x = df_temp,
      center = TRUE,
      scale. = TRUE
    )
    comp <- seq_along(pca$sdev)
    recon <- as.matrix(pca$x[, comp]) %*% t(as.matrix(pca$rotation[, comp]))
    # Se revierte la transformación centrado y escalado.
    recon <- scale(recon, center = FALSE, scale = 1/pca$scale)
    recon <- scale(recon, center = -1*pca$center, scale = FALSE)
    # Cálculo del error cuadrático medio de reconstrucción
    reconstruction_error <- apply(X = (recon - df_temp)^2, MARGIN = 1, FUN = mean)
    # Se añade el error de reconstrucción al dataframe original.
    df$reconstruction_error_pca <- reconstruction_error
  }
  if(method %in% c("all", "knn")){
    #
    df_scaled <- scale(df_temp)
    df_knn <- get.knn(data = df_scaled, k = 5)
    # Averaged distance to nearest neighbors
    df_score <- rowMeans(df_knn$nn.dist)
    df$knn_score <- df_score
  }
  if(method %in% c("all", "dbscan")){
    #
    df$dbscan_score <- lof(scale(df_temp), k = 5)
  }
}
```

```

}
# if(method %in% c("all", "oneclassssum")){
#   replications = 5
#   # function to obtain R-Squared from the data
#   oneclassssum <- function(data, indices) {
#     x <- data[indices,] # allows boot to select sample
#     y <- rep("oneclass", nrow(x))
#     model <- svm(x, type='one-classification') # train an one-classification model
#     # the output is Boolean
#     pred <- predict(model, x) # create predictions
#     return(pred)
#   }
#   # bootstrapping
#   results <- boot(
#     data=df_temp,
#     statistic=oneclassssum,
#     R=replications,
#     parallel = "multicore",
#     ncpus = 5
#   )
#   temp <- as.data.frame(results$t)
#   df$oneclassssum_score <- as.vector(colSums(temp))/replications
# }
if(method %in% c("all", "isolationforest")){
  #
  df_tree <- iForest(X=df_temp, nt = 100)
  # Scores near 1 indicate anomalies (small path length)
  df_score <- predict(df_tree, newdata = df_temp)
  df$isolationforest_score <- df_score
}

if(method %in% c("all", "autoencoder")){
  h2o.init(ip = "localhost",
    # Si emplea un único core los resultados son reproducibles.
    # Si se dispone de muchos datos, mejor emplear varios cores.
    nthreads = 1
  )
  # Se eliminan los datos del cluster por si ya había sido iniciado.
  h2o.removeAll()
  h2o.no_progress()
  df_h2o <- as.h2o(
    x = df_temp,
    destination_frame = "df_h2o"
  )
  # División de las observaciones en conjunto de entrenamiento y validación.
  df_h2o_split <- h2o.splitFrame(data = df_h2o, ratios = 0.8, seed = 123)
  df_h2o_train <- df_h2o_split[[1]]
  df_h2o_validacion <- df_h2o_split[[2]]
  # Se define la variables empleadas por el autoencoder
  predictores <- setdiff(h2o.colnames(df_h2o), "y")
  # Entrenamiento del modelo autoencoder con 2 neuronas en la capa oculta.
  autoencoder <- h2o.deeplearning(
    x = predictores,

```



```

training_frame = df_h2o_train,
validation_frame = df_h2o_validacion,
activation      = "Tanh",
autoencoder     = TRUE,
hidden          = c(2),
epochs          = 50,
ignore_const_cols = FALSE,
score_each_iteration = TRUE,
# Seed solo válido cuando se emplea un único core
seed = 999
)

reconstruction_error <- h2o.anomaly(
  object = autoencoder,
  data   = df_h2o[, predictores],
  per_feature = FALSE
)
reconstruction_error <- as.data.frame(reconstruction_error)
reconstruction_error <- reconstruction_error[, 1]
# Se añade el error de reconstrucción al dataframe original.
df$anomaly_score_autoencoder <- reconstruction_error

}

return(df)
}

temp <- OutlierAnomalyDetection(df = trainset2, columns, method = "all")

```

```

## Connection successful!
##
## R is connected to the H2O cluster:
##   H2O cluster uptime:      47 minutes 11 seconds
##   H2O cluster timezone:    Europe/Paris
##   H2O data parsing timezone: UTC
##   H2O cluster version:     3.32.0.1
##   H2O cluster version age:  2 months and 28 days
##   H2O cluster name:        H2O_started_from_R_ferna_wkw207
##   H2O cluster total nodes: 1
##   H2O cluster total memory: 7.07 GB
##   H2O cluster total cores: 8
##   H2O cluster allowed cores: 1
##   H2O cluster healthy:     TRUE
##   H2O Connection ip:       localhost
##   H2O Connection port:     54321
##   H2O Connection proxy:    NA
##   H2O Internal Security:   FALSE
##   H2O API Extensions:      Amazon S3, Algos, AutoML, Core V3, TargetEncoder, Core V4
##   R Version:                R version 4.0.3 (2020-10-10)

```

```
df_temp <- temp %>% ungroup() %>% dplyr::select(cook_score, mahalanobis_score, reconstruction_error_pca
```

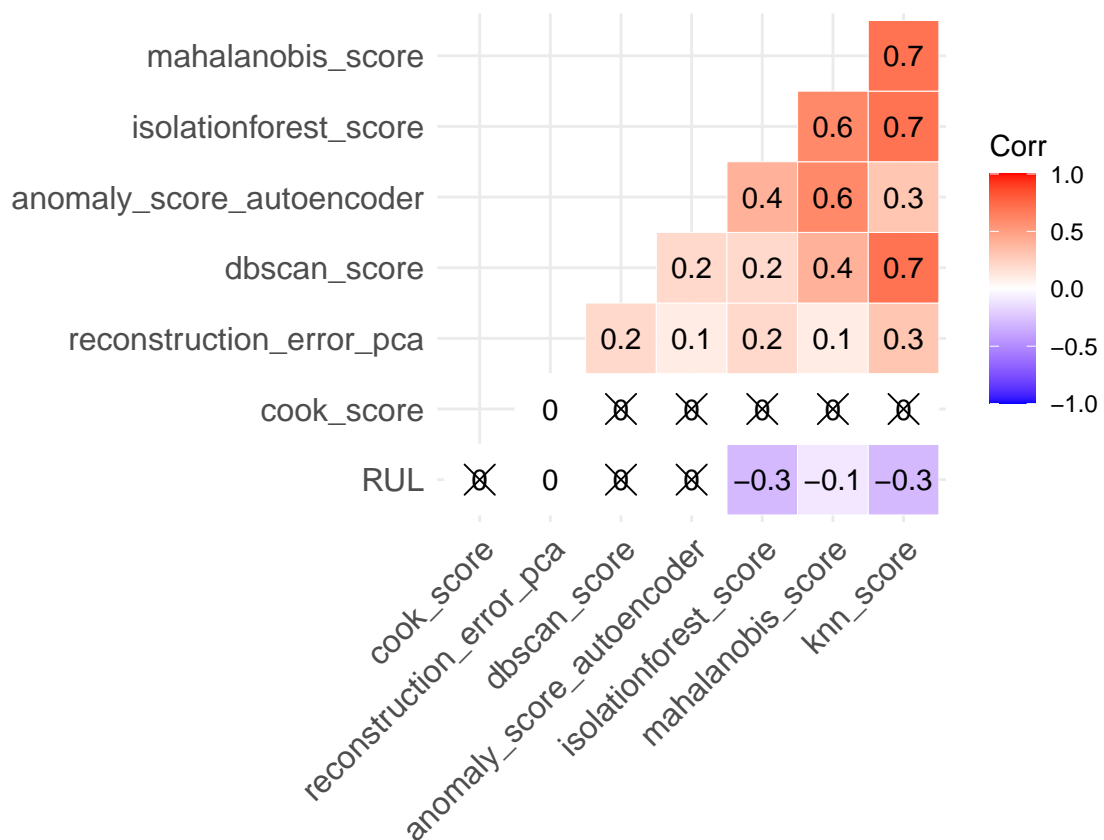
```

# Compute a correlation matrix
corr <- round(cor(df_temp), 1)
# Compute a matrix of correlation p-values
p.mat <- cor_pmat(df_temp)

# Visualize the lower triangle of the correlation matrix
# Barring the no significant coefficient
corr.plot <- ggcorrplot(
  corr, hc.order = TRUE, type = "lower", lab = TRUE, outline.col = "white",
  p.mat = p.mat
)

corr.plot

```

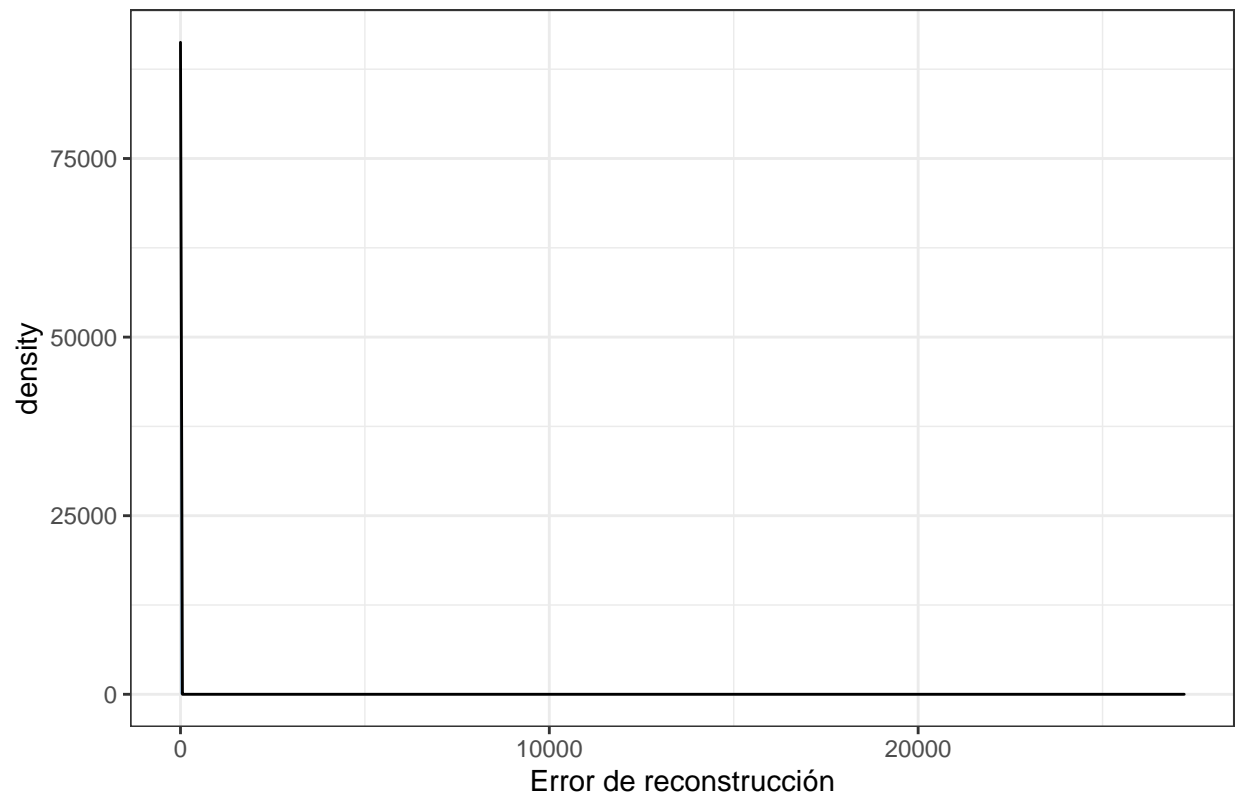


```

# Distribución del error de reconstrucción.
tibble(df_temp) %>%
  ggplot(aes(x = cook_score)) +
  geom_density(fill = "steelblue") +
  labs(title = "Distribución de las medidas de la Distancia de Cook",
        x = "Error de reconstrucción") +
  theme_bw()

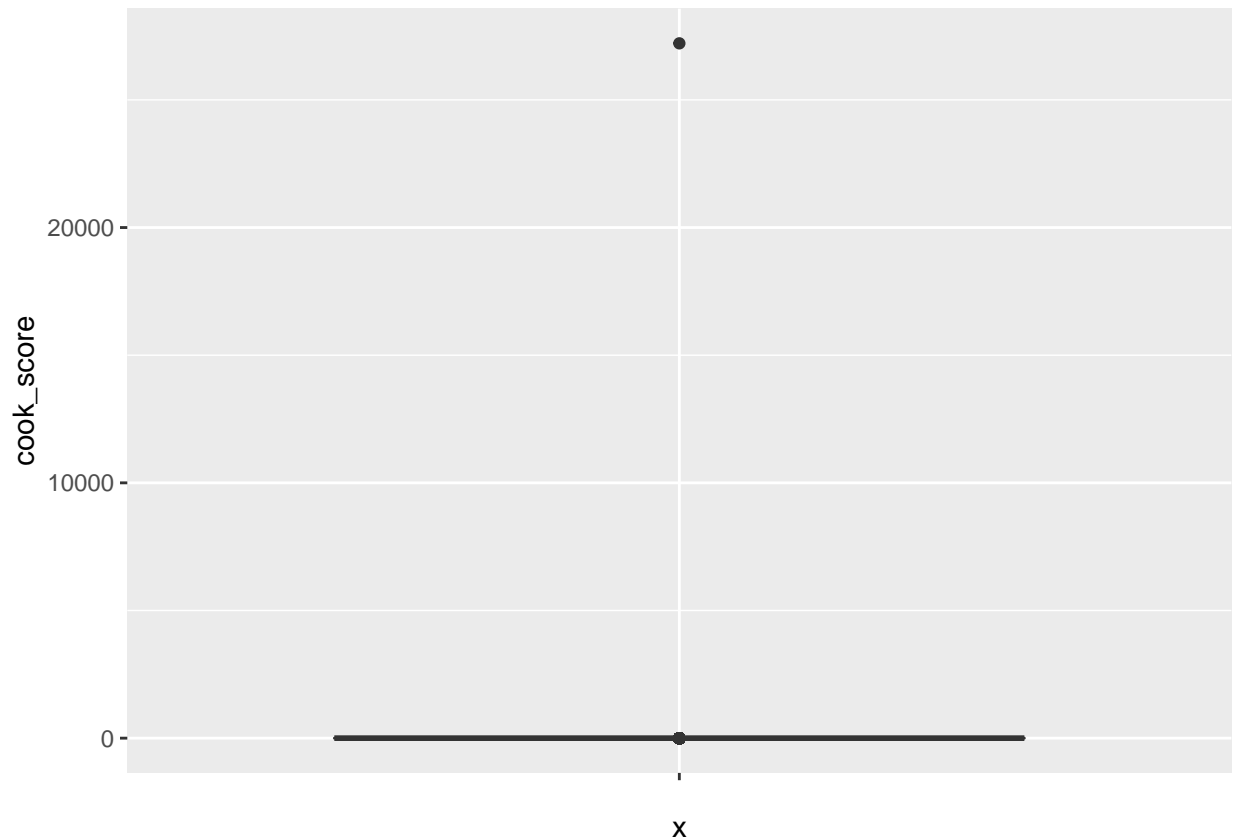
```

Distribución de las medidas de la Distancia de Cook



Según la distribución de la distancia de Cook, podemos ver claramente que existen outliers en los valores más grandes.

```
ggplot(data = df_temp, aes(x = "", y = cook_score)) +  
  geom_boxplot() #+
```



```
# coord_cartesian(ylim = c(0, 150)) # I set the y axis scale so the plot looks better.
```

Podemos aplicar la prueba de Grubbs. La prueba de Grubbs detecta un valor atípico a la vez (valor más alto o más bajo), por lo que las hipótesis nula y alternativa son las siguientes:

- H_0 : El valor más alto **no** es un outlier
- H_1 : El valor más alto es un outlier

El valor de $p < 0,05$ significa que al nivel de significancia del 5%, rechazamos la hipótesis de que el valor más no es un valor atípico.

Abajo generamos un loop dónde evaluamos en cada turno el valor más grande, si es outliers, se lo eliminamos y aplicamos otra vez la prueba hasta que no exista ningún outlier según Grubbs.

```
actual_p_value <- 0
p_value <- 0.05
outliers <- c()
df_no_cook_outlier <- df_temp

while (actual_p_value < p_value) {

  grubbs_test <- grubbs.test(df_no_cook_outlier$cook_score)
  actual_p_value <- grubbs_test$p.value
  value <- gsub(pattern = "highest value ", "", x = grubbs_test$alternative)
  value <- as.numeric(gsub(pattern = " is an outlier", "", x = value))
}
```

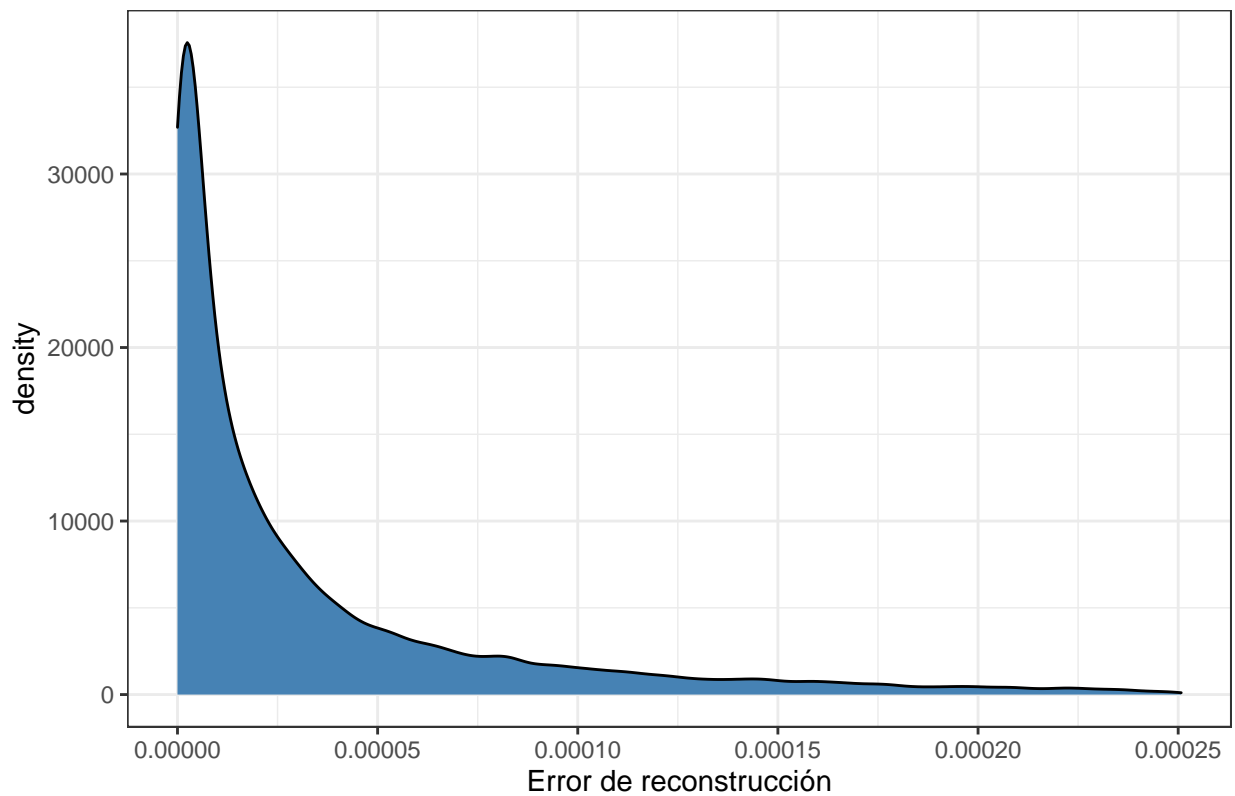
```

outliers <- c(outliers, value)
if(actual_p_value < p_value){
  df_no_cook_outlier <- df_no_cook_outlier %>% filter(cook_score < max(cook_score))
  #print(max(df_no_cook_outlier$cook_score))
  #print(paste0("The value ", value, " is an outlier and it's been removed."))
} else {
  break
}
}

tibble(df_no_cook_outlier) %>%
  ggplot(aes(x = cook_score)) +
  geom_density(fill = "steelblue") +
  labs(title = "Distribución de las medidas de la Distancia de Cook",
       x = "Error de reconstrucción") +
  theme_bw()

```

Distribución de las medidas de la Distancia de Cook



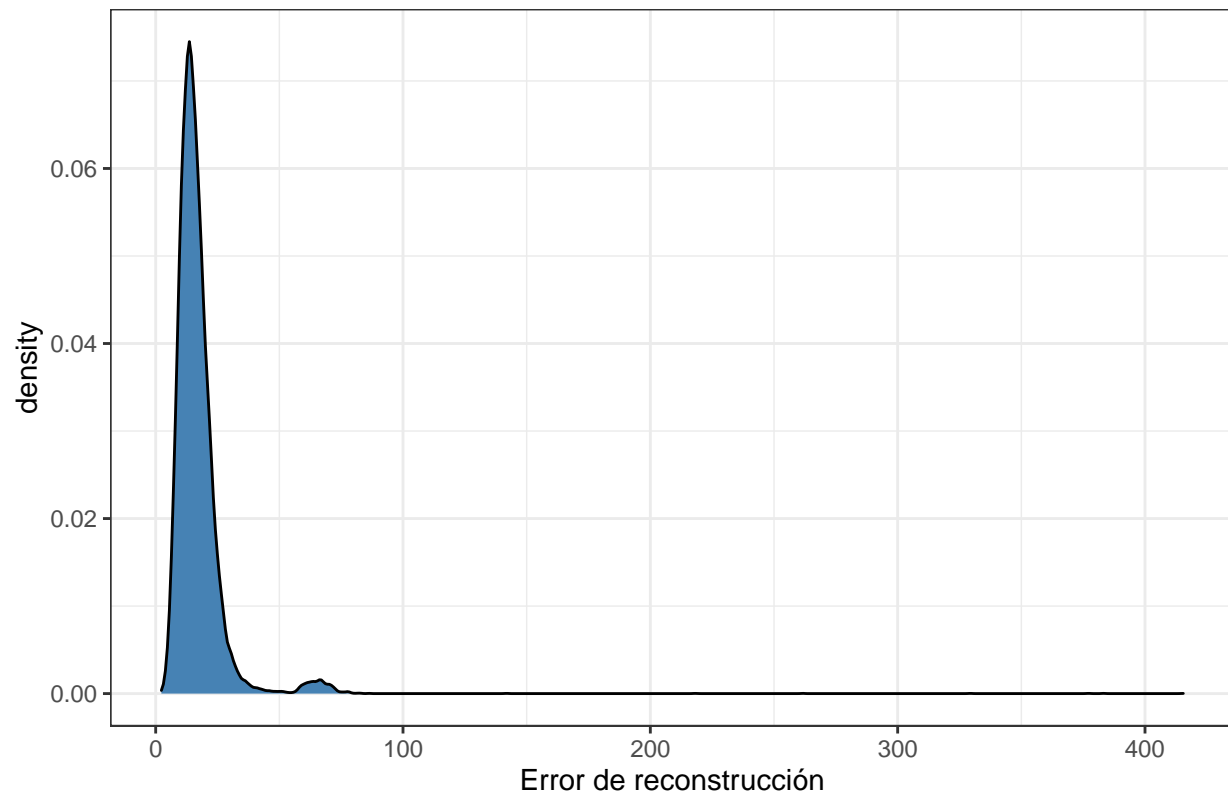
```

# Distribución del error de reconstrucción.
tibble(df_temp) %>%
  ggplot(aes(x = mahalanobis_score)) +
  geom_density(fill = "steelblue") +
  labs(title = "Distribución de las medidas de la Distancia de Mahalanobis",
       x = "Error de reconstrucción") +

```

```
theme_bw()
```

Distribución de las medidas de la Distancia de Mahalanobis



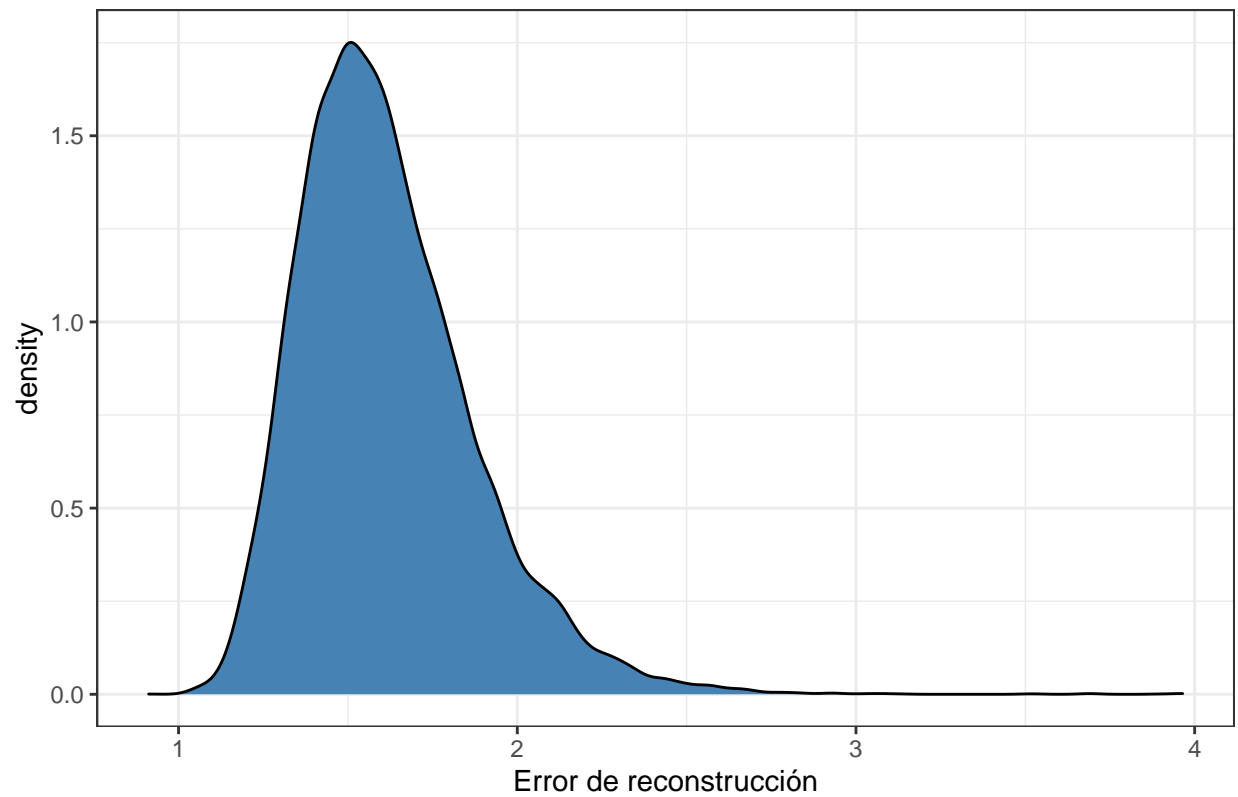
```
# Distribución del error de reconstrucción.  
tibble(df_temp) %>%  
  ggplot(aes(x = reconstruction_error_pca)) +  
  geom_density(fill = "steelblue") +  
  labs(title = "Distribución del error de reconstrucción de PCA",  
        x = "Error de reconstrucción") +  
  theme_bw()
```

A graph on a grid showing a function that starts at a high value on the y-axis and decays rapidly towards the x-axis. The curve is oscillatory, with several small peaks and troughs as it approaches zero. The area under the curve is shaded in blue.

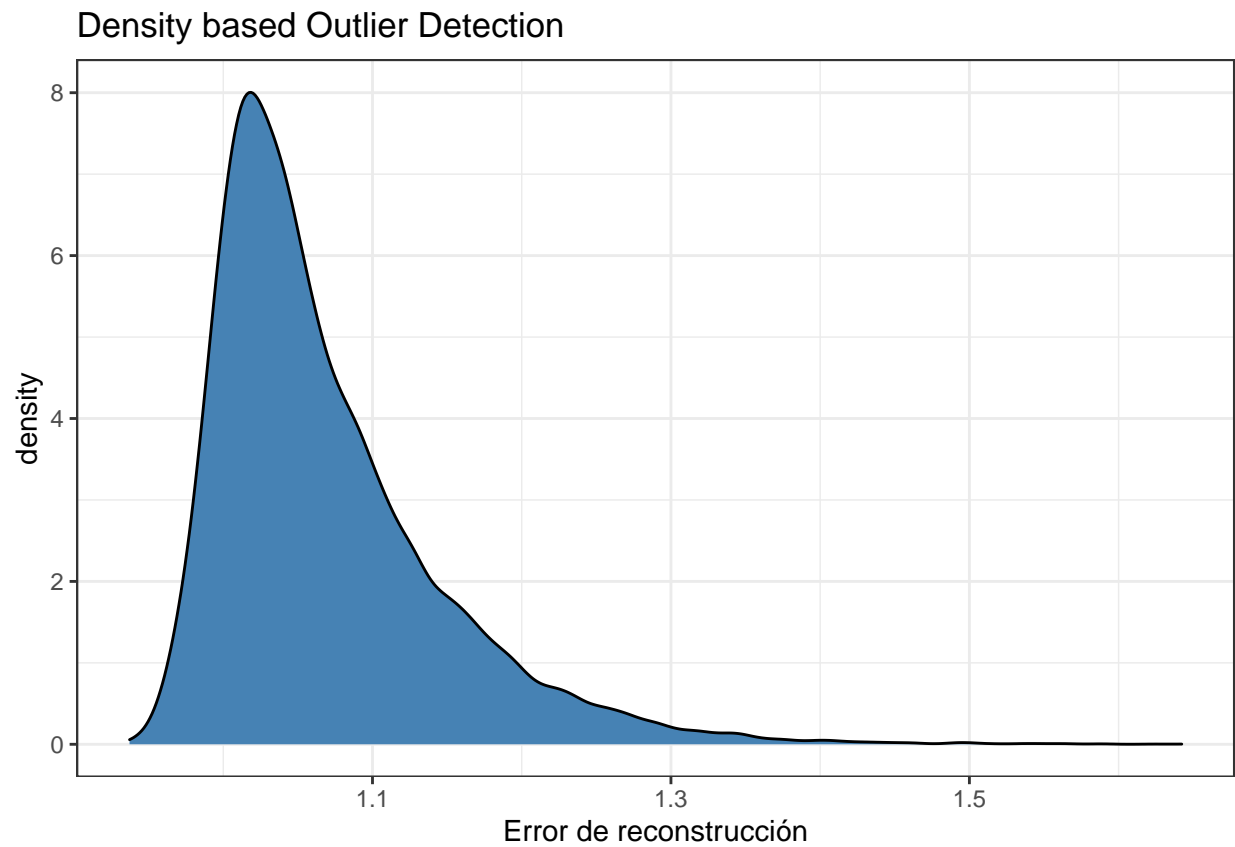
Error de reconstrucción

```
# Distribución del error de reconstrucción.
tibble(df_temp) %>%
  ggplot(aes(x = knn_score)) +
  geom_density(fill = "steelblue") +
  labs(title = "Distance based Outlier Detection",
       x = "Error de reconstrucción") +
  theme_bw()
```

Distance based Outlier Detection

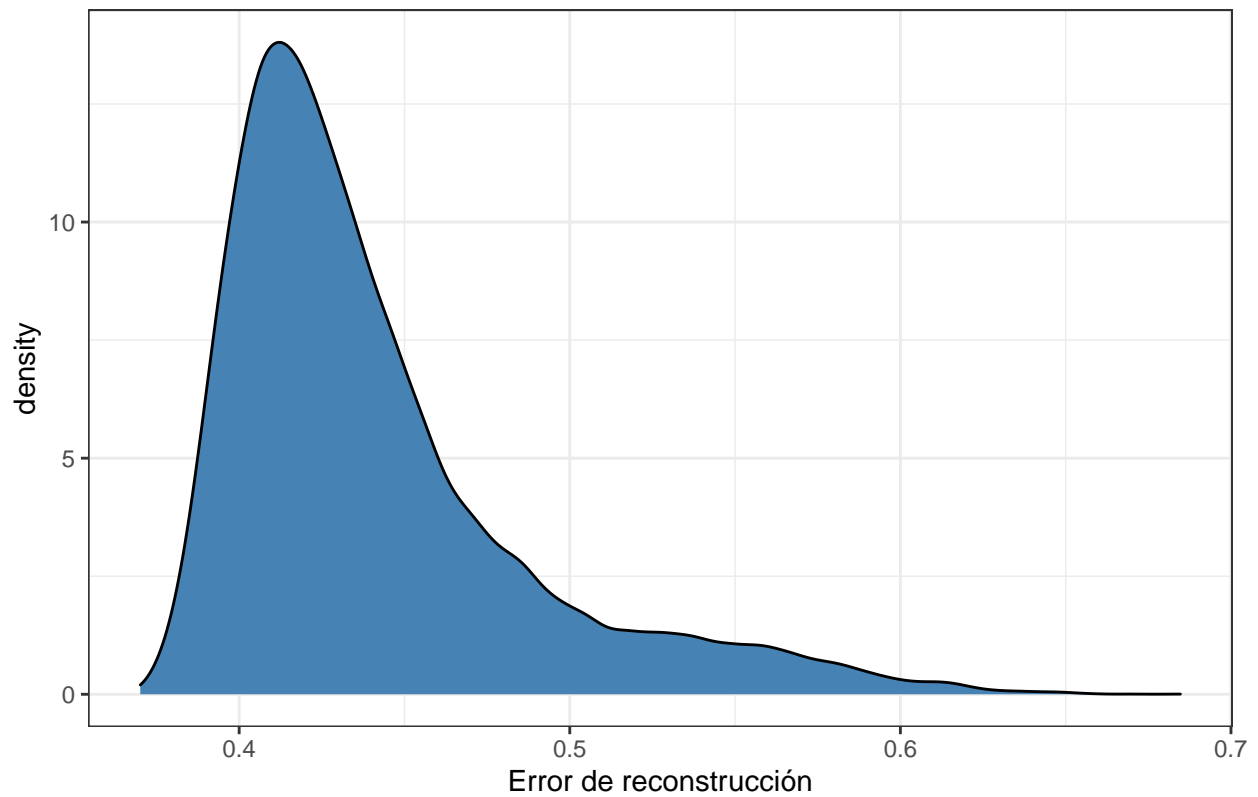


```
# Distribución del error de reconstrucción.  
tibble(df_temp) %>%  
  ggplot(aes(x = dbscan_score)) +  
  geom_density(fill = "steelblue") +  
  labs(title = "Density based Outlier Detection",  
        x = "Error de reconstrucción") +  
  theme_bw()
```

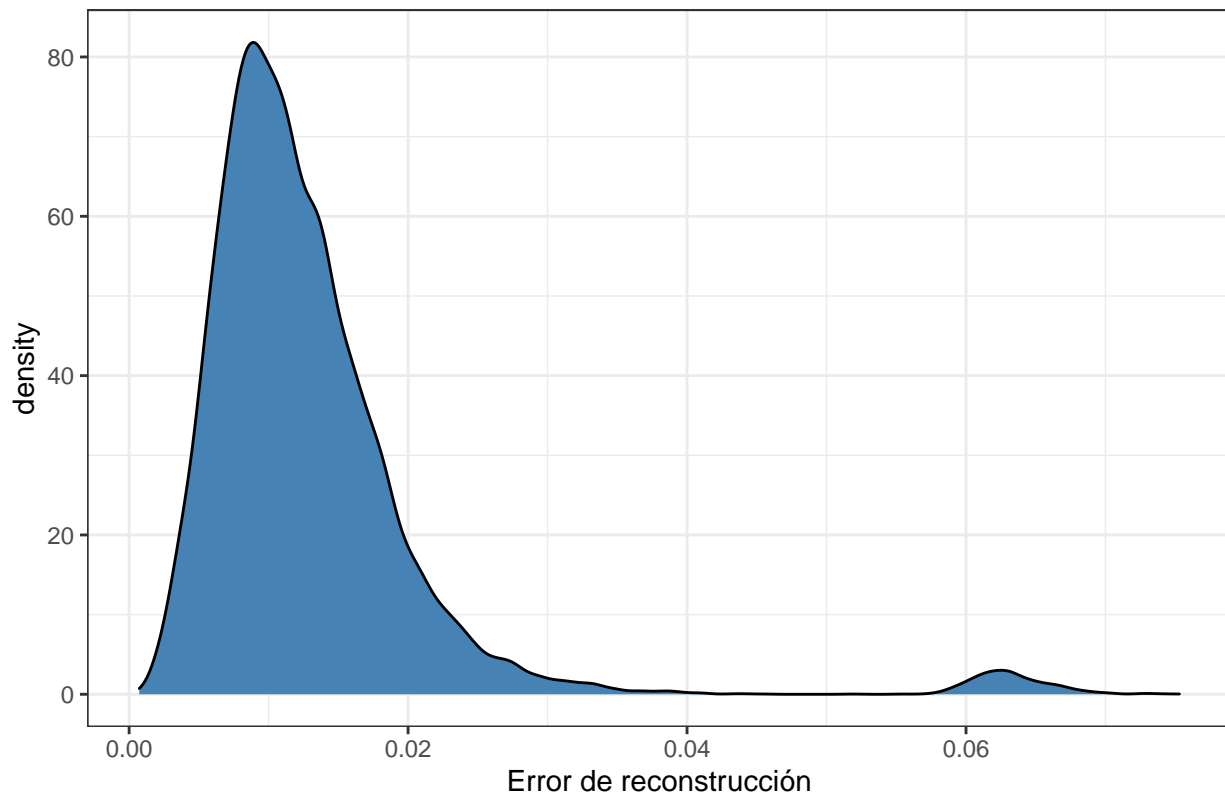
```
# Distribución del error de reconstrucción.  
tibble(df_temp) %>%  
  ggplot(aes(x = isolationforest_score)) +  
  geom_density(fill = "steelblue") +  
  labs(title = "Bagging of trees based Outlier Detection",  
        x = "Error de reconstrucción") +  
  theme_bw()
```

Bagging of trees based Outlier Detection



```
# Distribución del error de reconstrucción.
tibble(df_temp) %>%
  ggplot(aes(x = anomaly_score_autoencoder)) +
  geom_density(fill = "steelblue") +
  labs(title = "Deep Learning based Outlier Detection",
       x = "Error de reconstrucción") +
  theme_bw()
```

Deep Learning based Outlier Detection



4. Análisis de los datos

4.1. Selección de los grupos de datos que se quieren analizar/comparar (planificación de los análisis a aplicar).

Nuestro análisis se centrará en las máquinas cuando estas fallaron ($RUL = 0$). Además vamos a examinar dos grupos distintos, uno son las máquinas que fallaron con menos de 200 vuelos y el otro grupo son las máquinas con 200 o más vuelos.

```
df <- trainset2 %>%  
  ungroup() %>%  
  filter(RUL == 0) %>%  
  mutate(  
    label = if_else(time_in_cycles < 200, "bad", "good")  
  ) %>%  
  dplyr::select(-RUL, -time_in_cycles, -unit_number)  
  
table(df$label)
```

```
##  
## bad good  
## 52 48
```

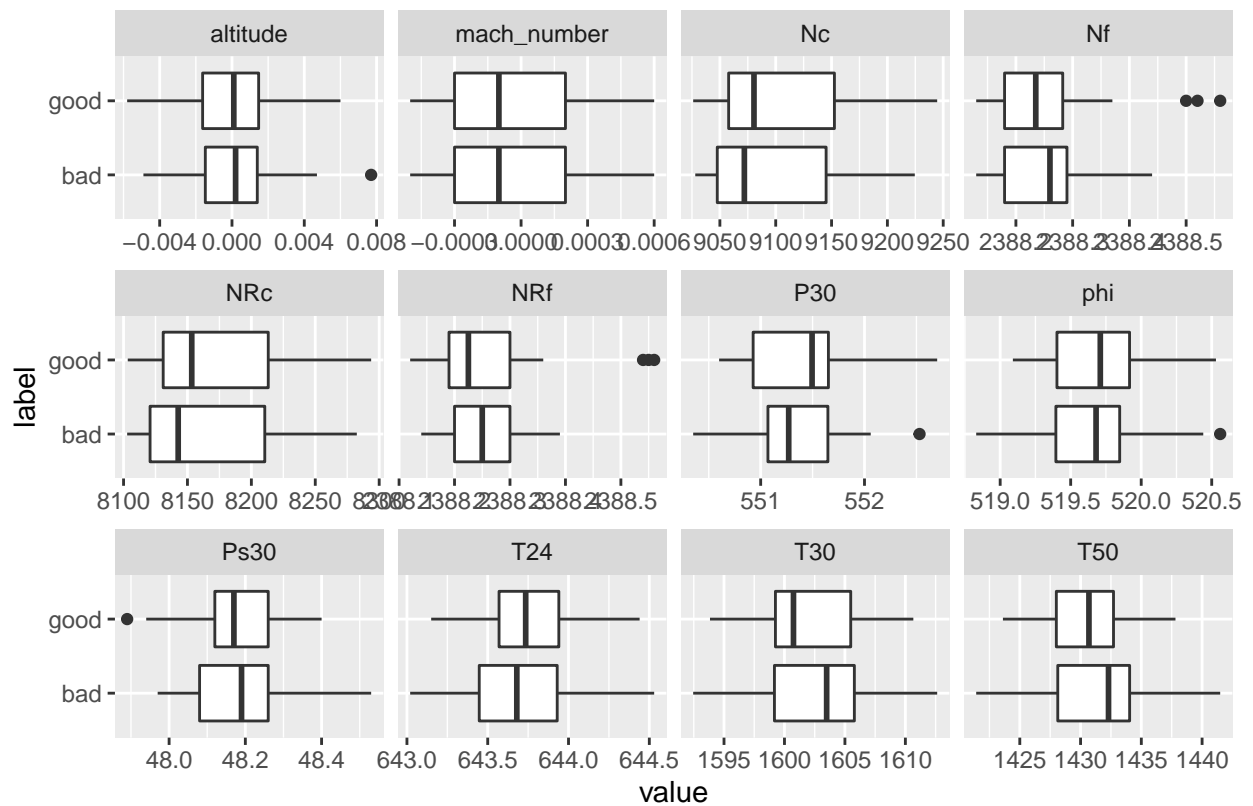
```
temp <- DataProfiling(df)
```

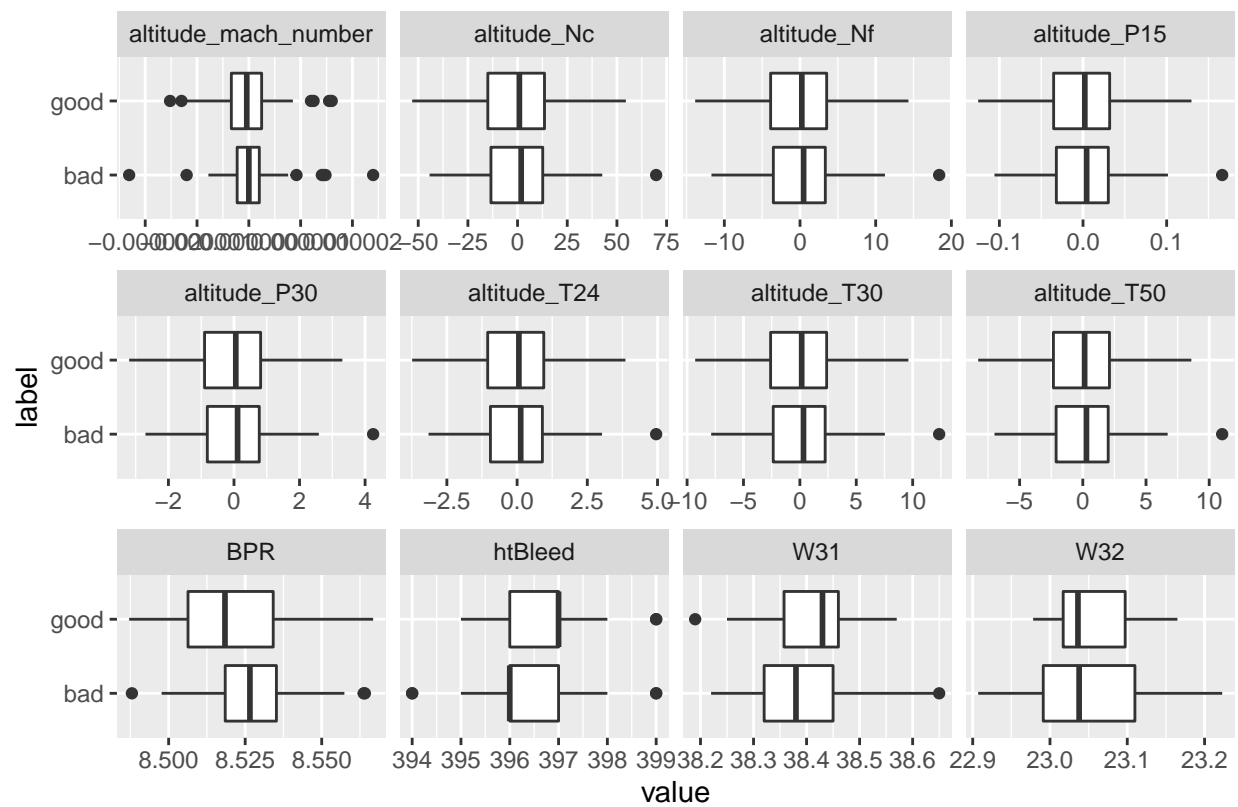
```
unique_columns <- temp$names[temp$n_unique == 1 & !is.na(temp$n_unique)]
print(unique_columns)
```

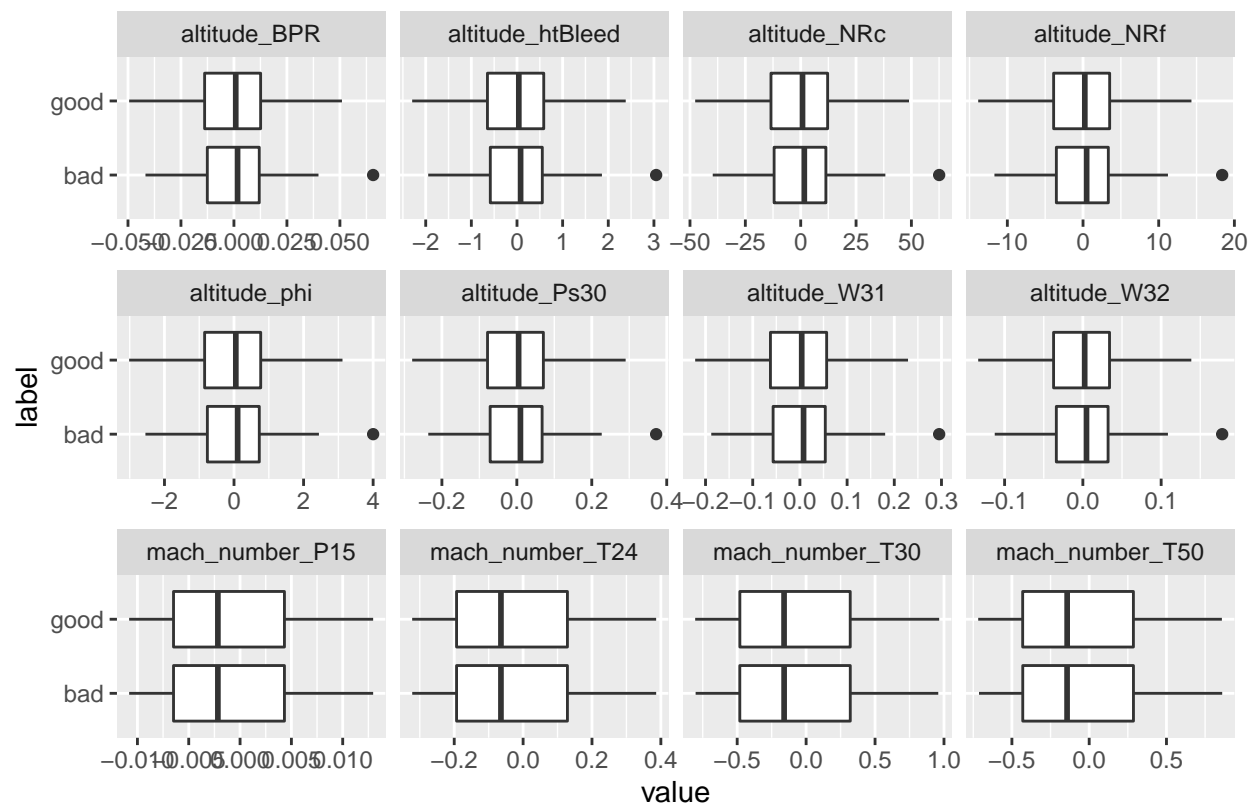
```
## [1] "P15"          "lag1_P15"     "lag2_P15"     "lag3_P15"     "lag4_P15"
## [6] "lag5_P15"     "rollsd3_P15"  "rolliqr3_P15" "rollsd5_P15"  "rolliqr5_P15"
## [11] "cummax_P15"
```

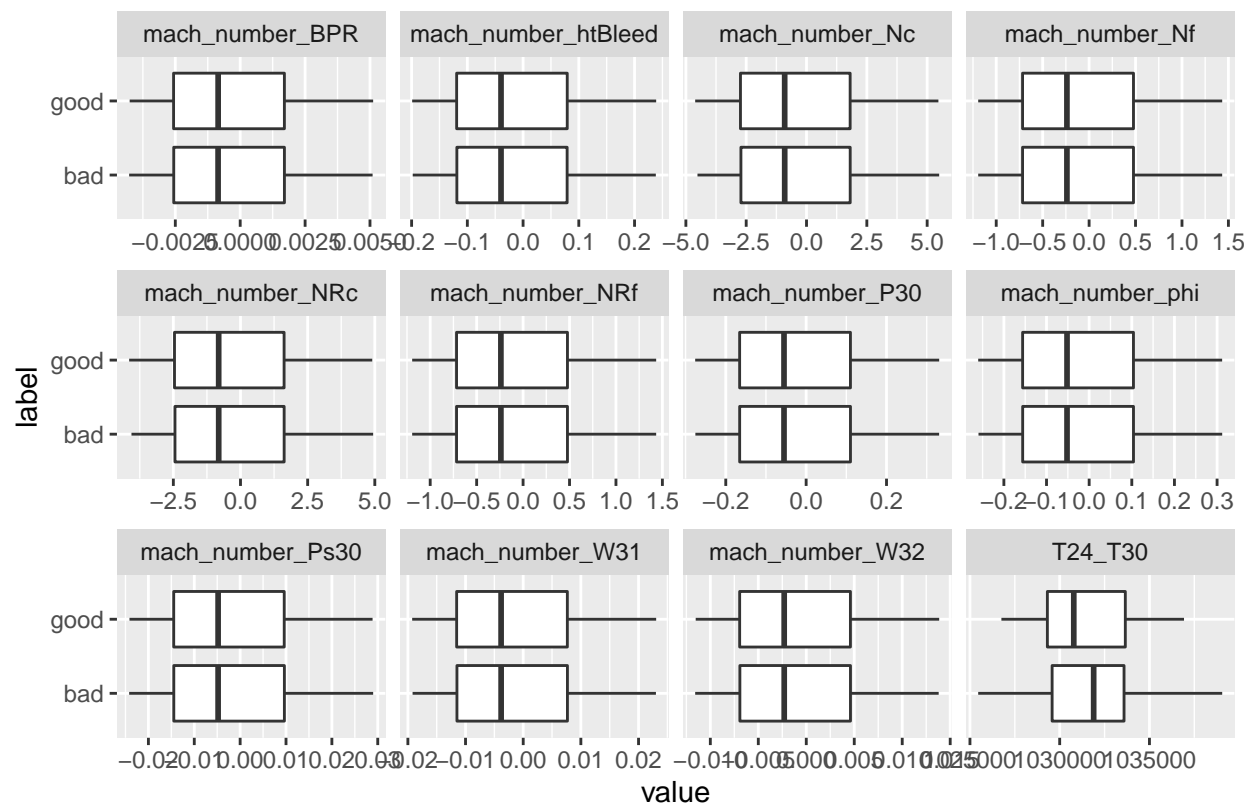
```
df <- df %>% dplyr::select(-one_of(unique_columns))
```

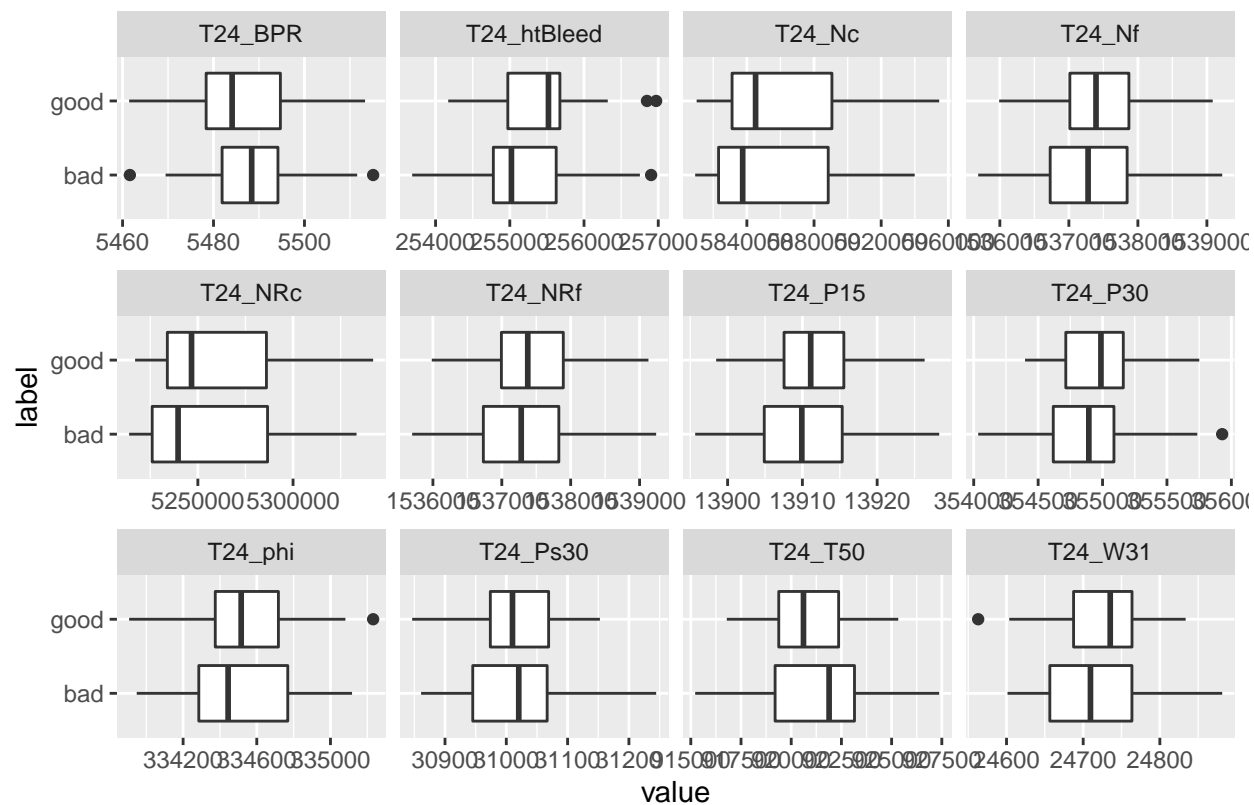
```
plot_boxplot(df, by = "label")
```

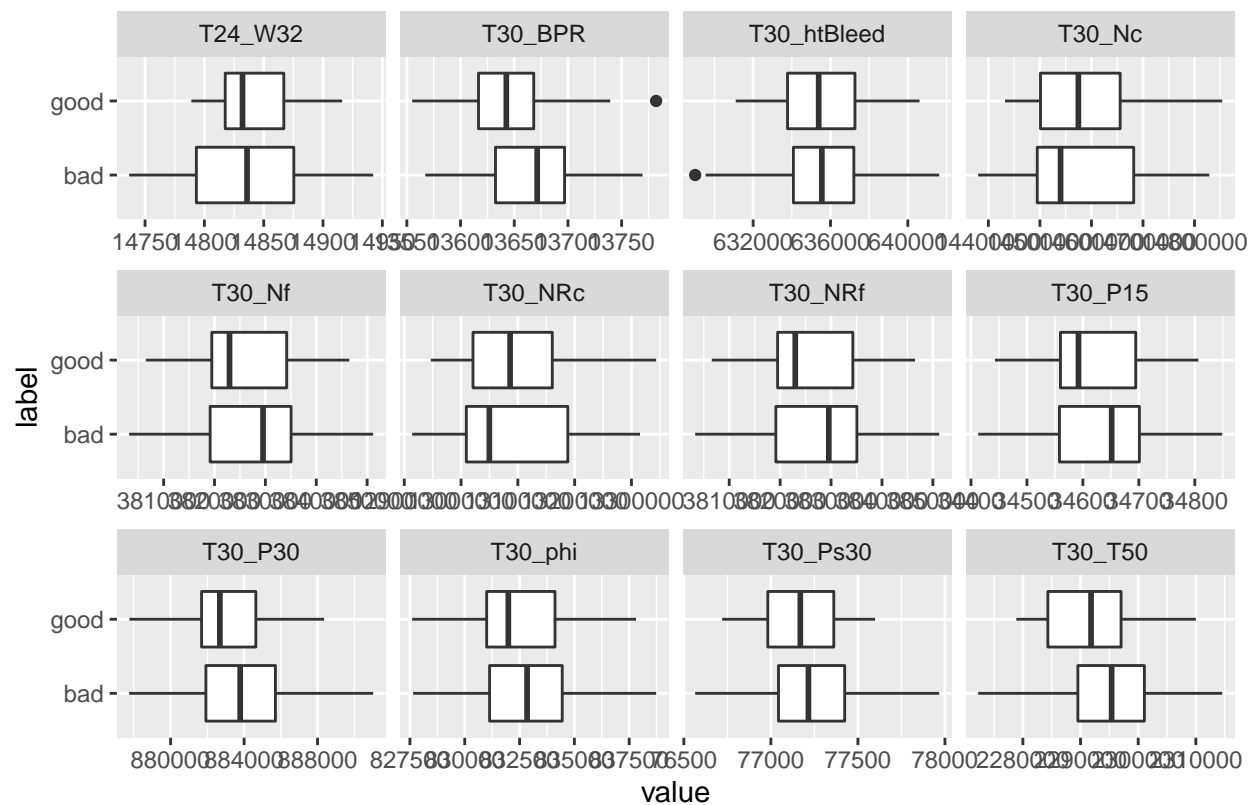


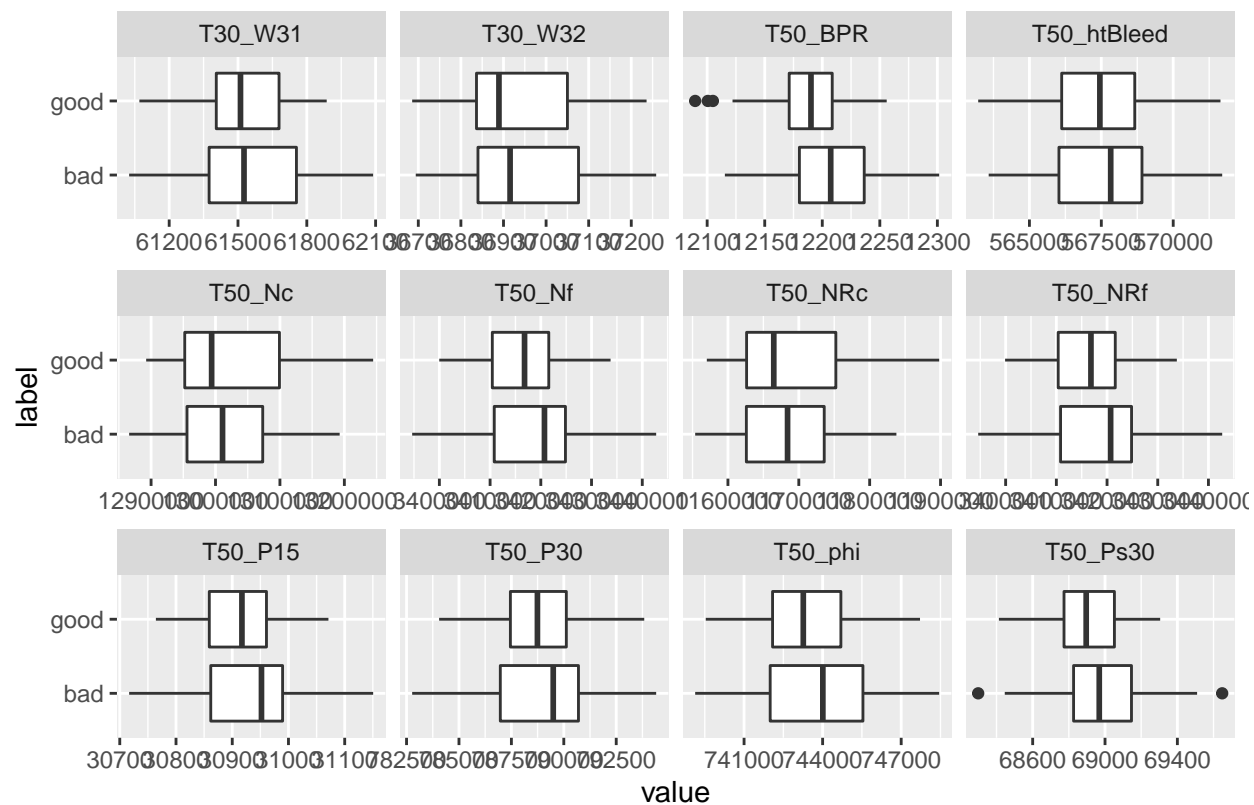


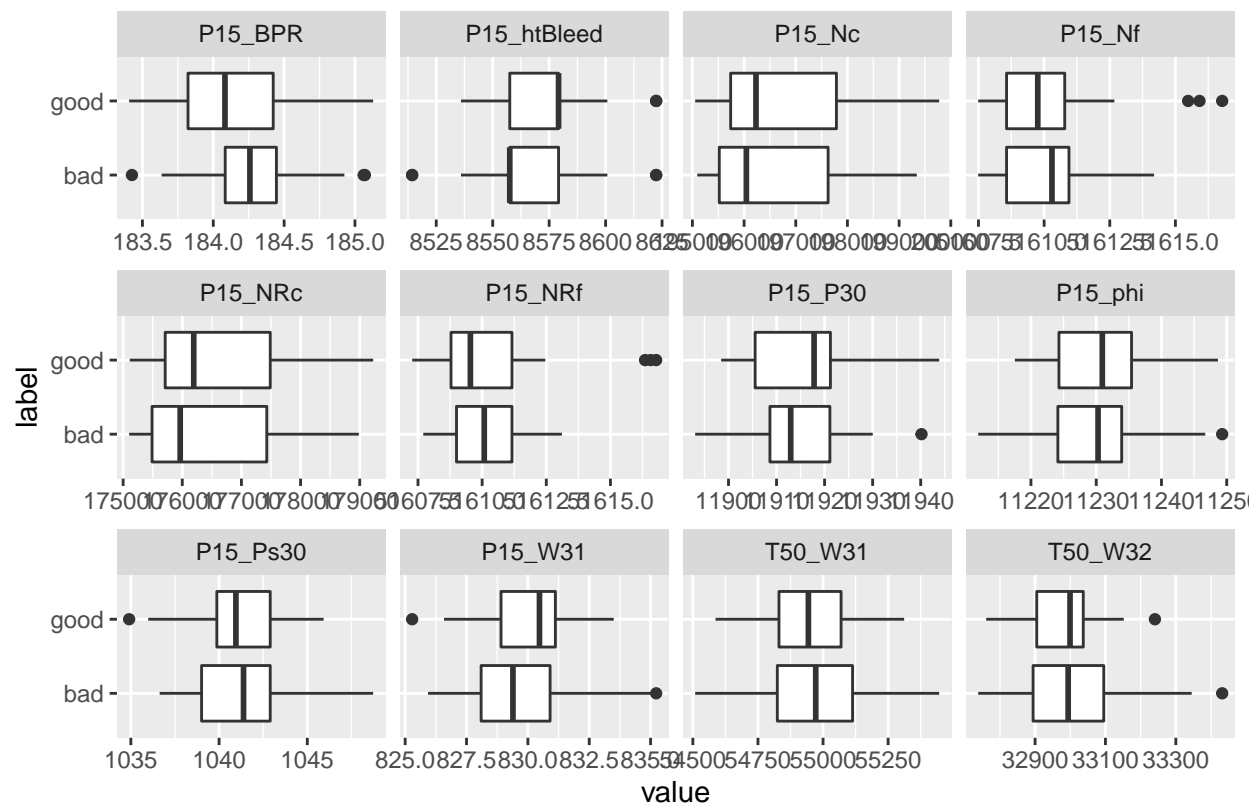


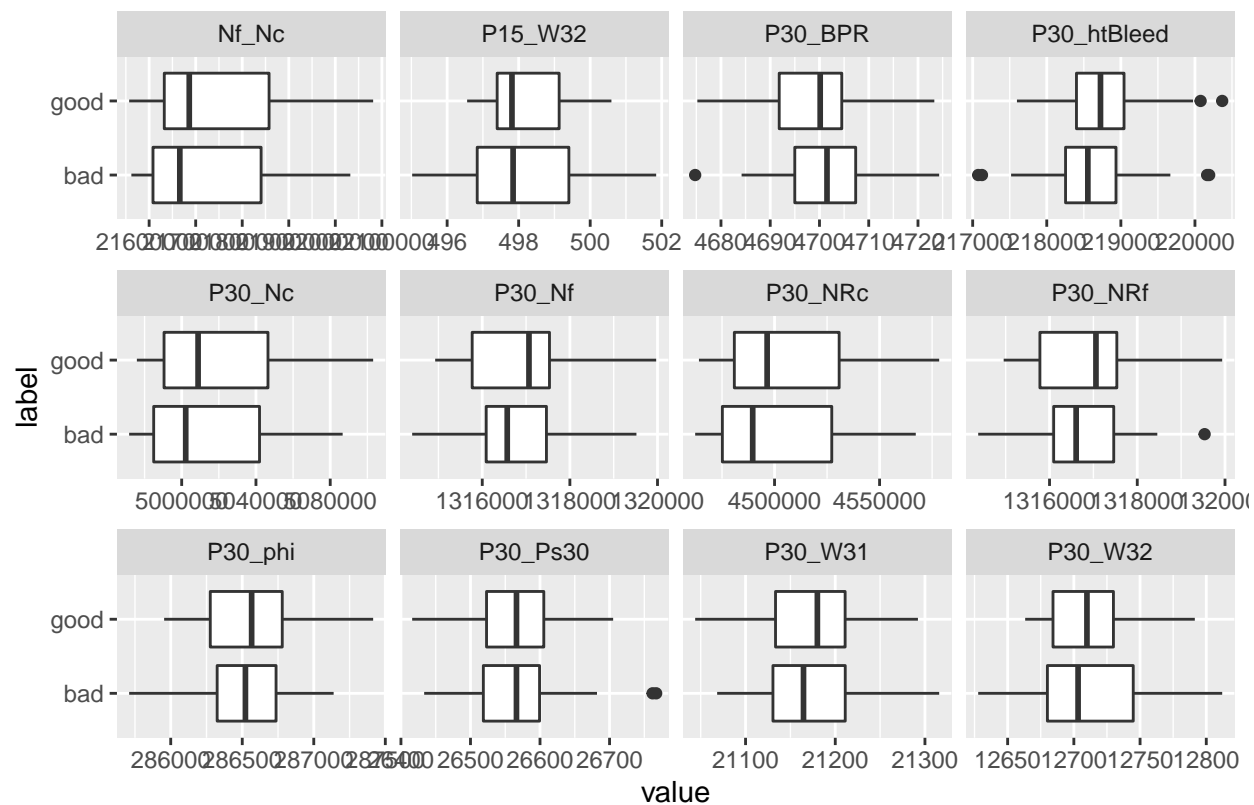


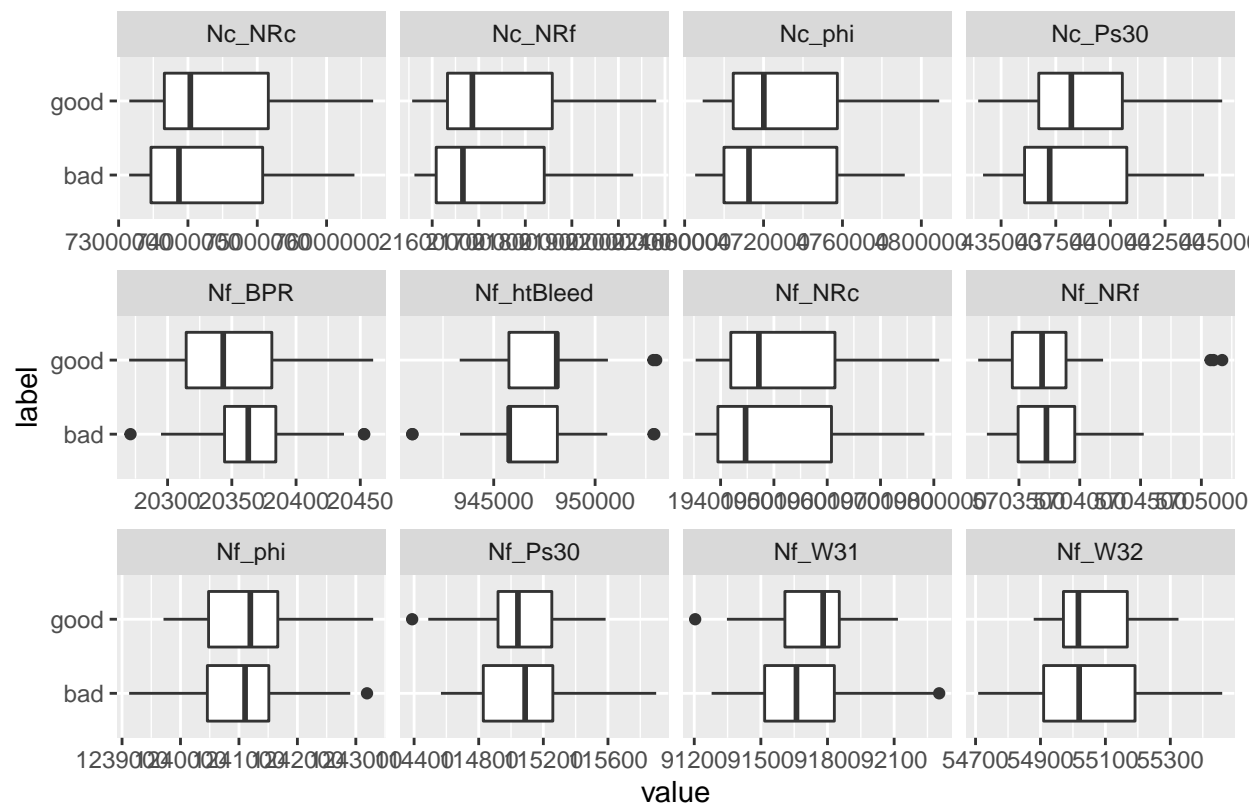


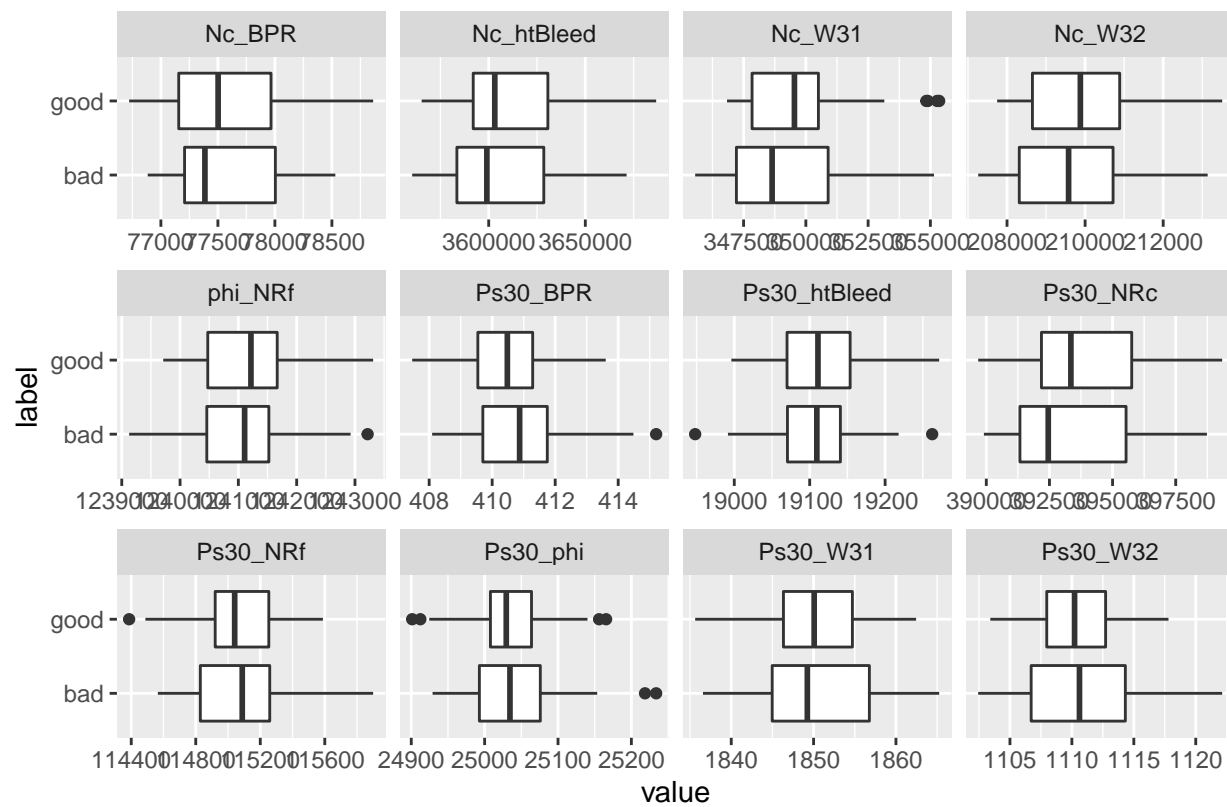


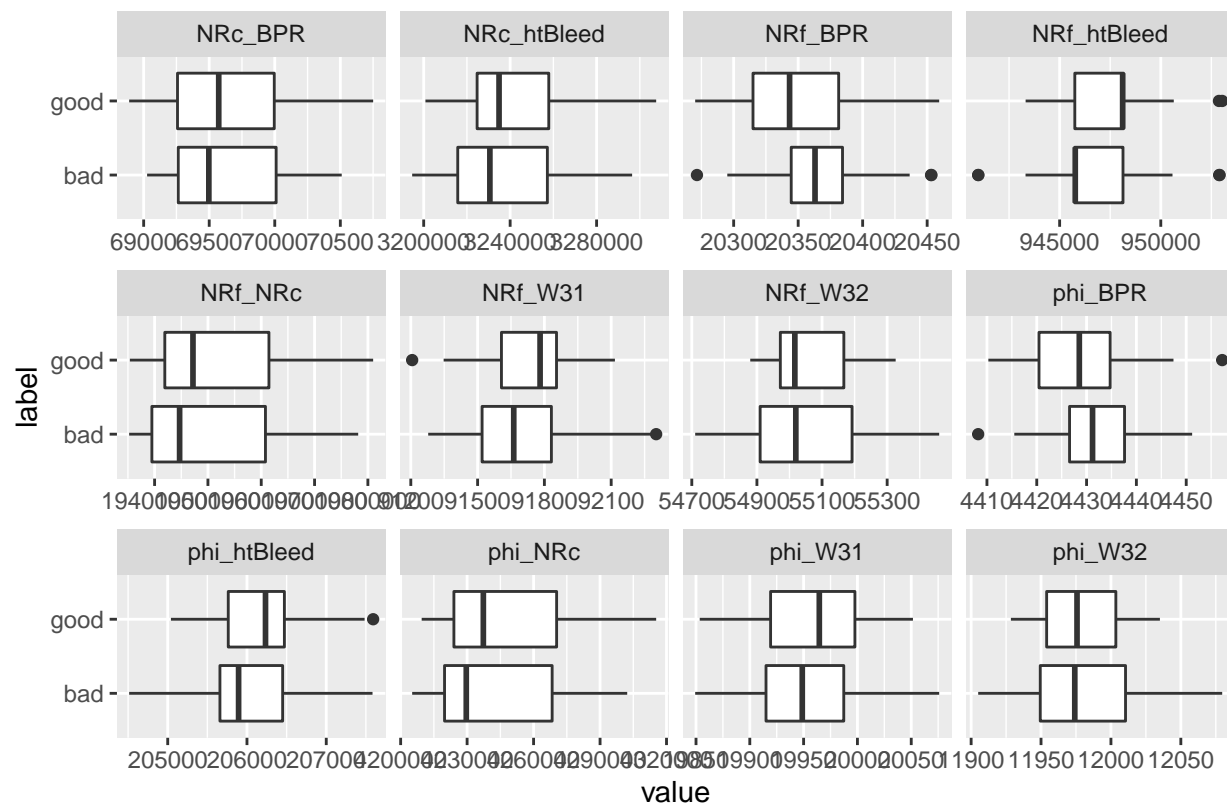


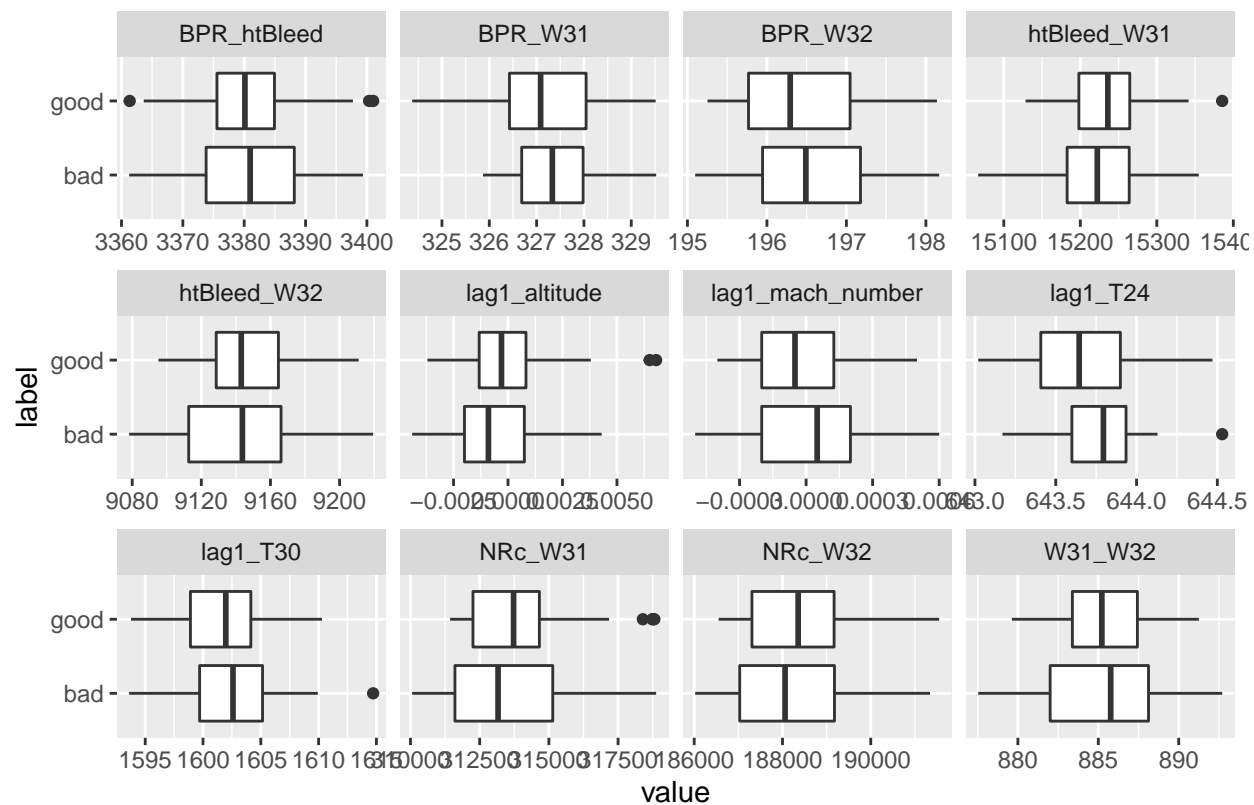


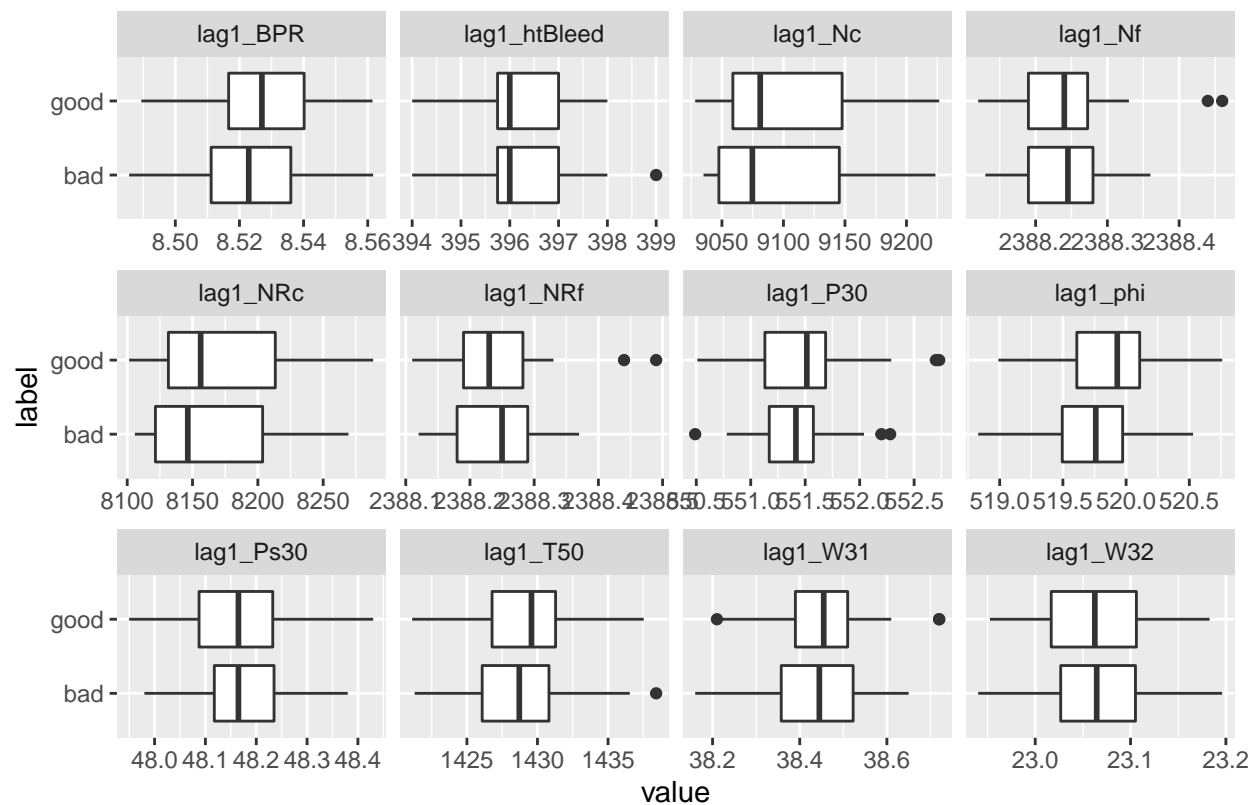


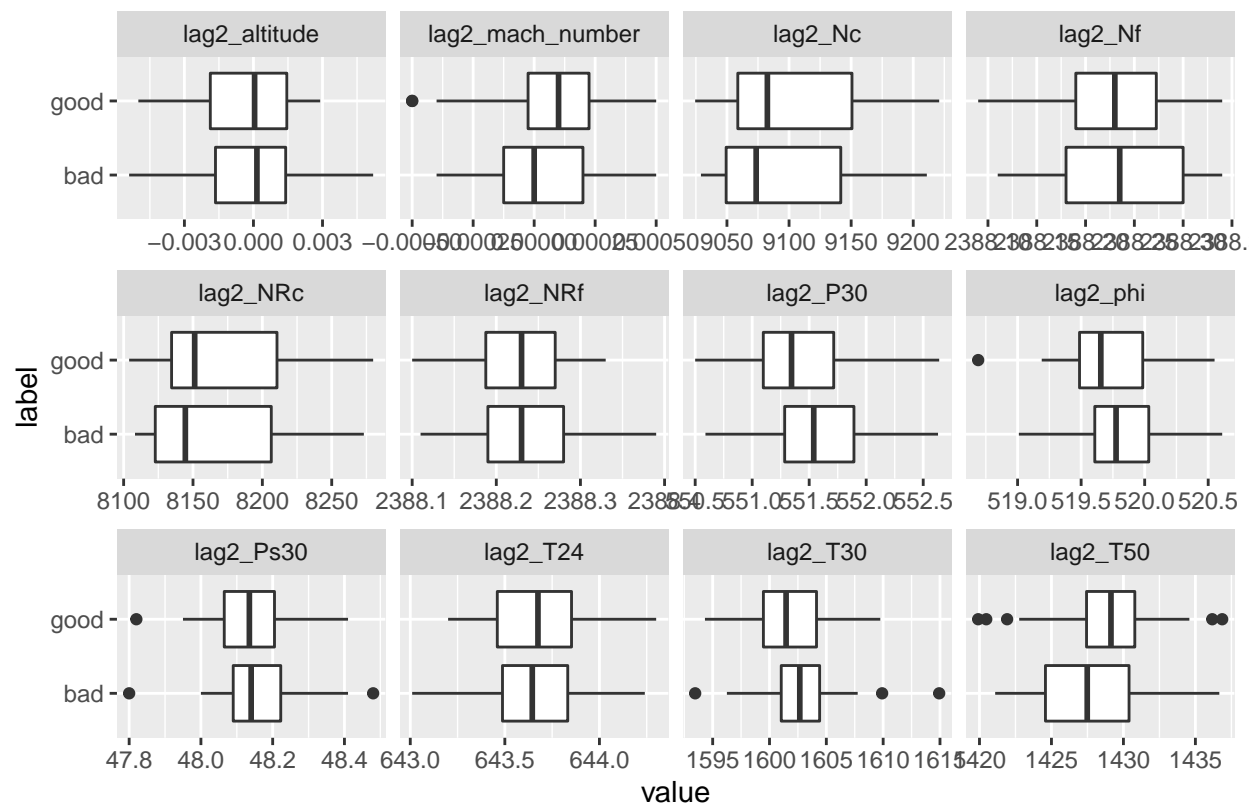


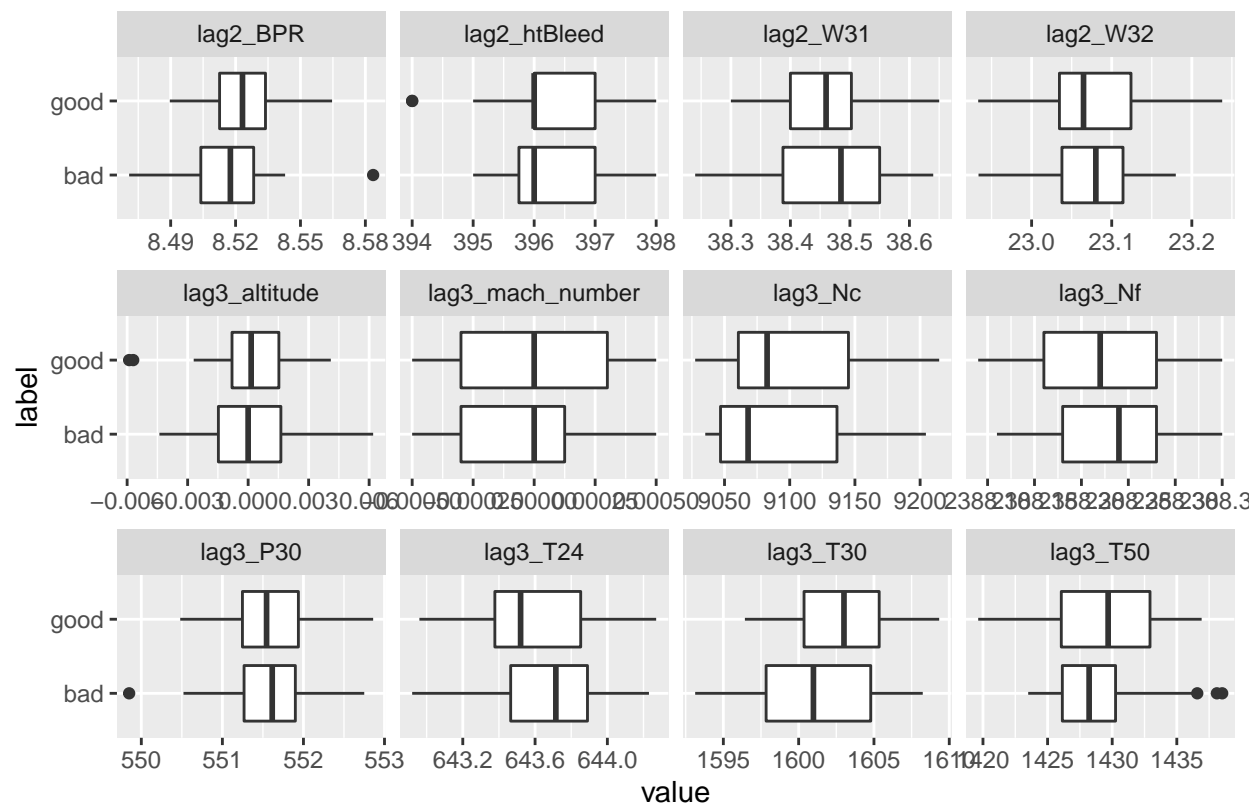


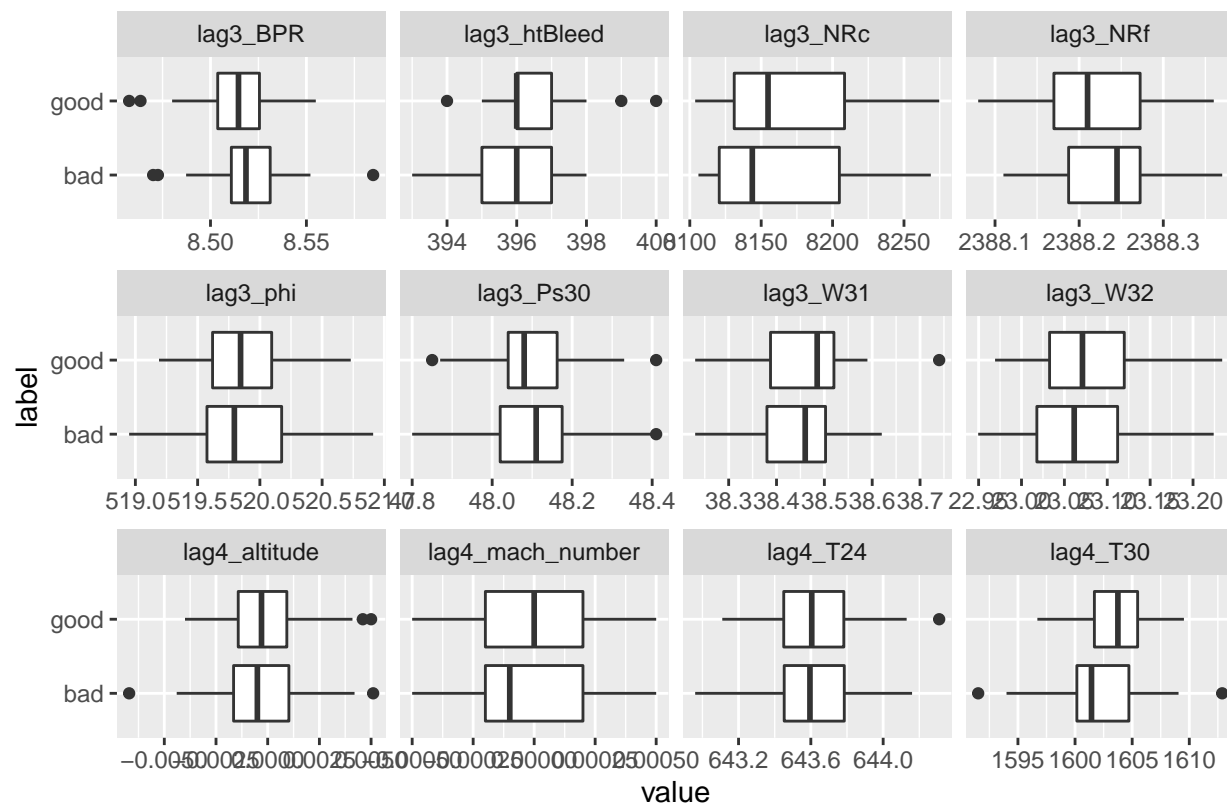


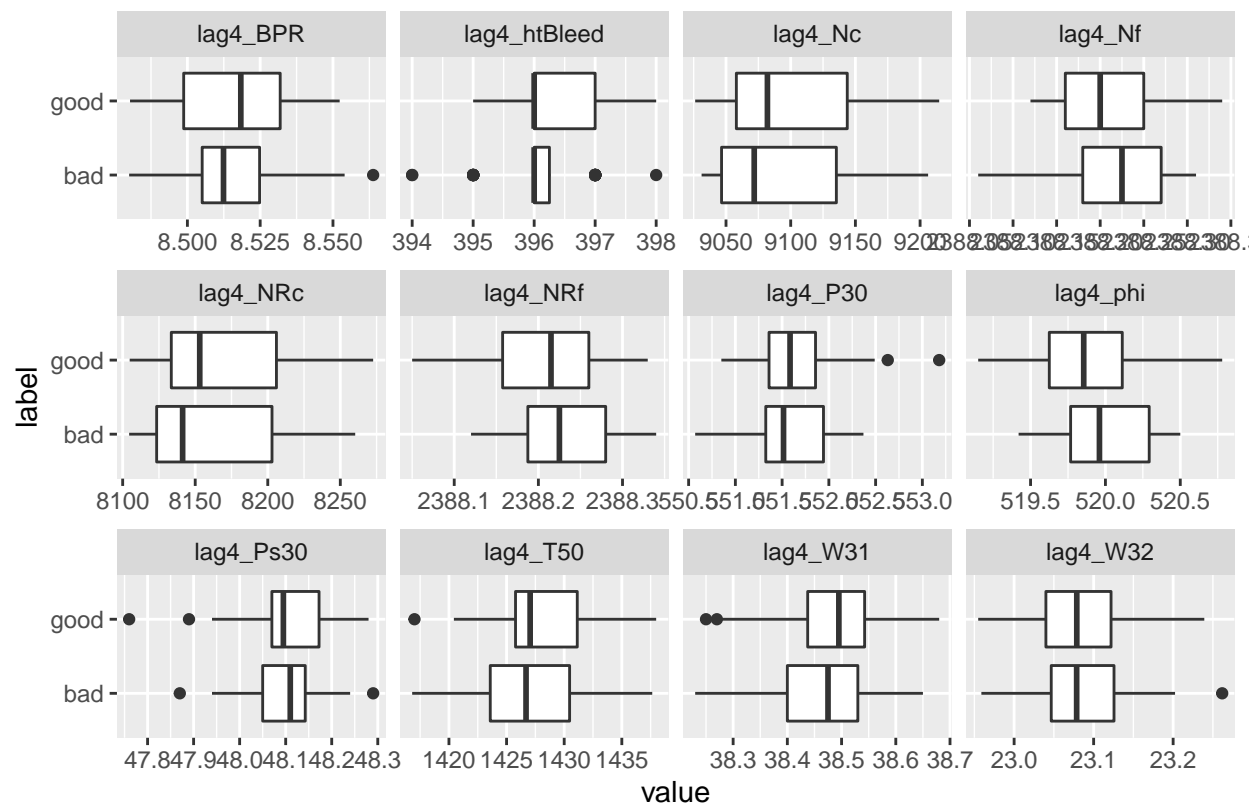


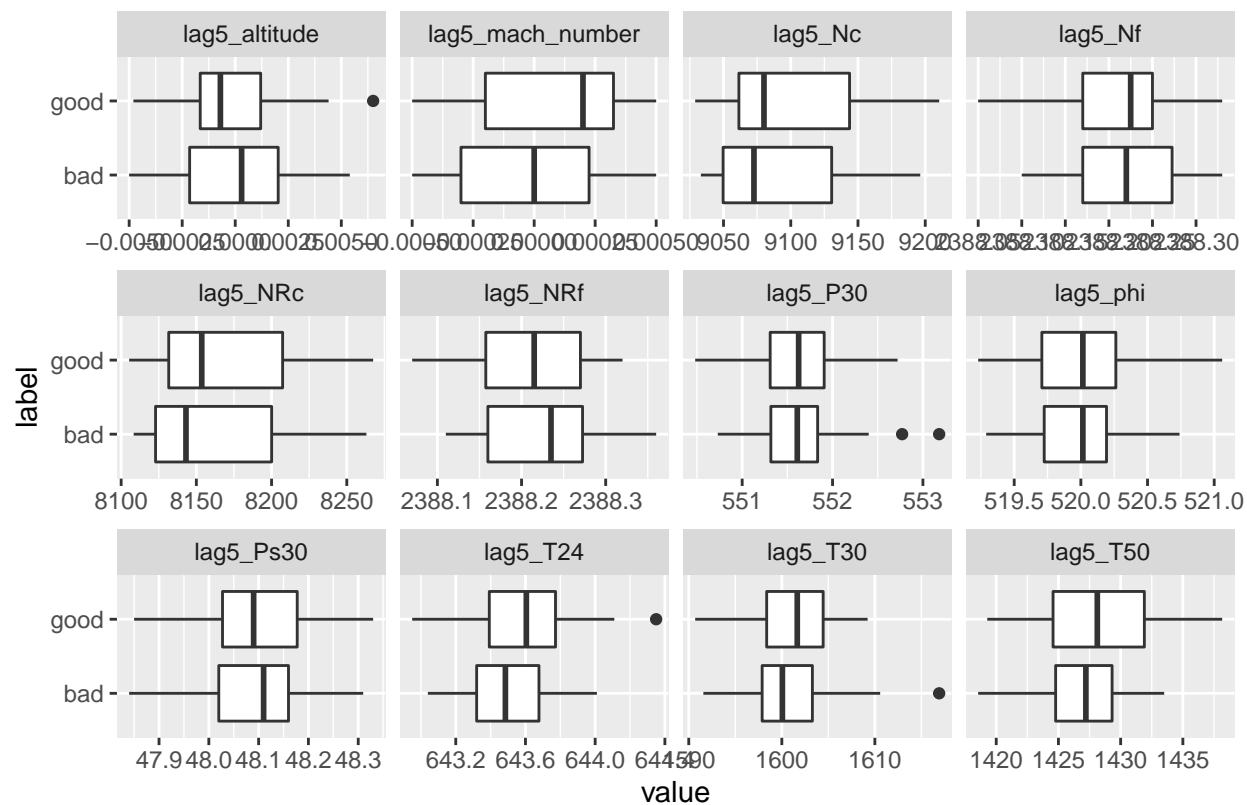


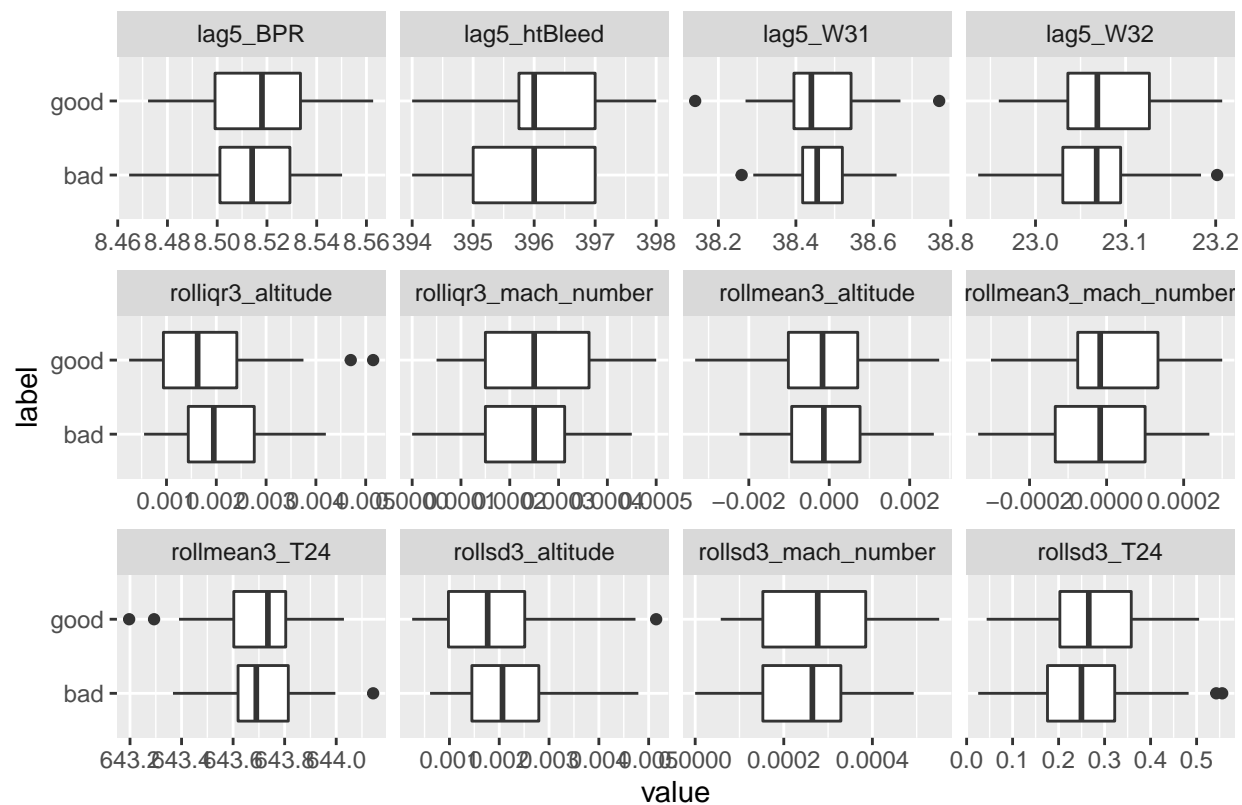


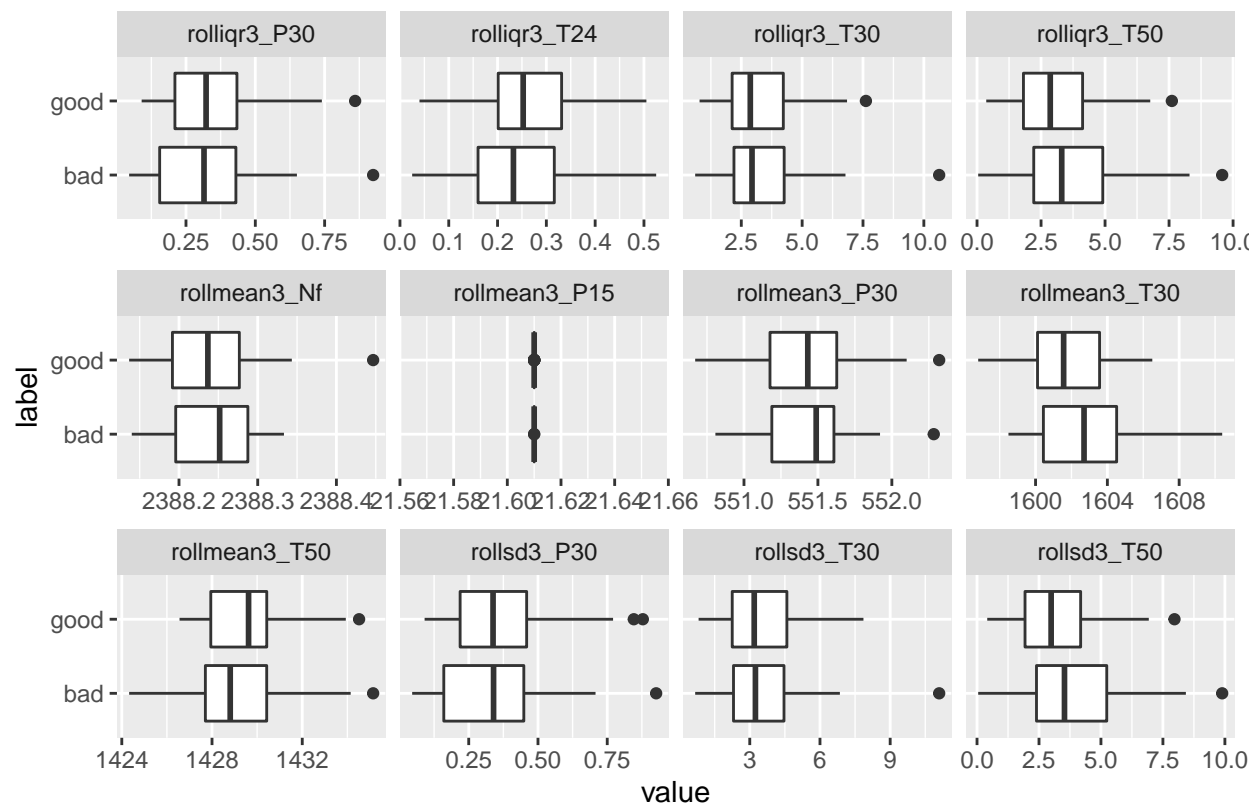


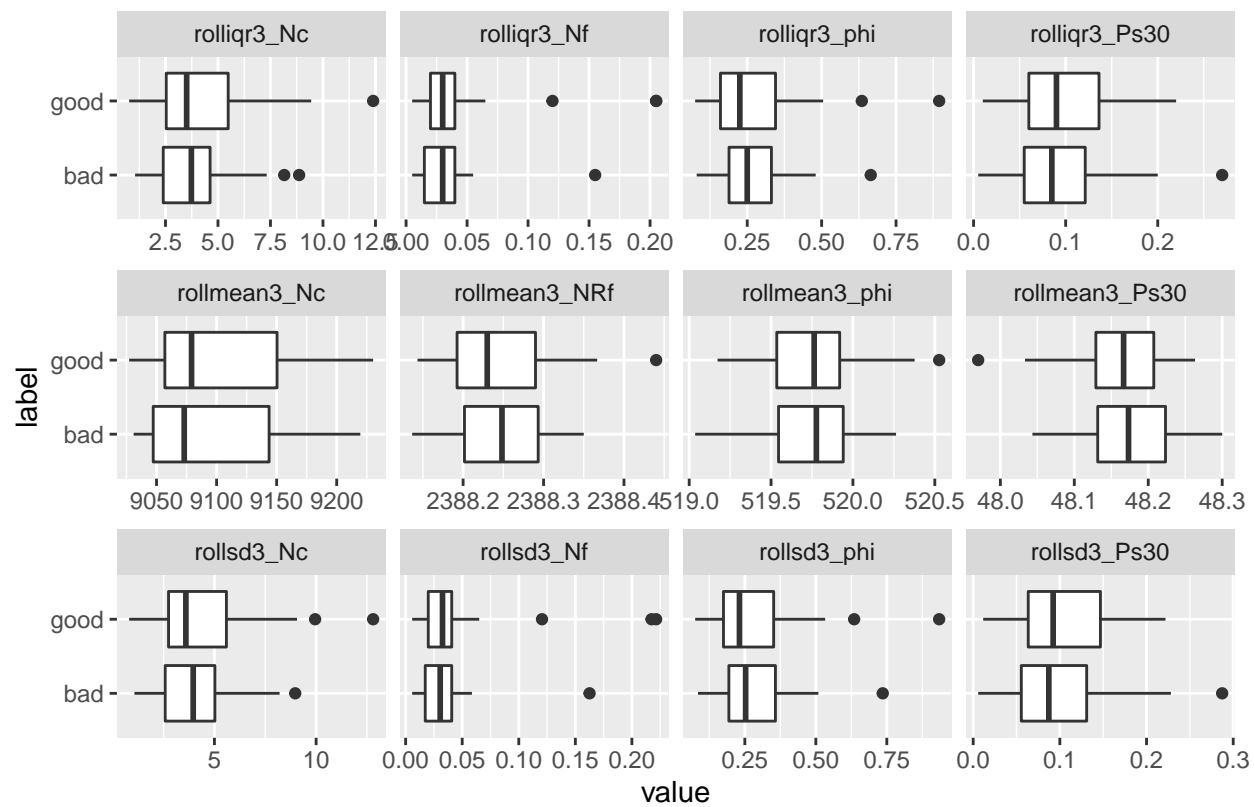


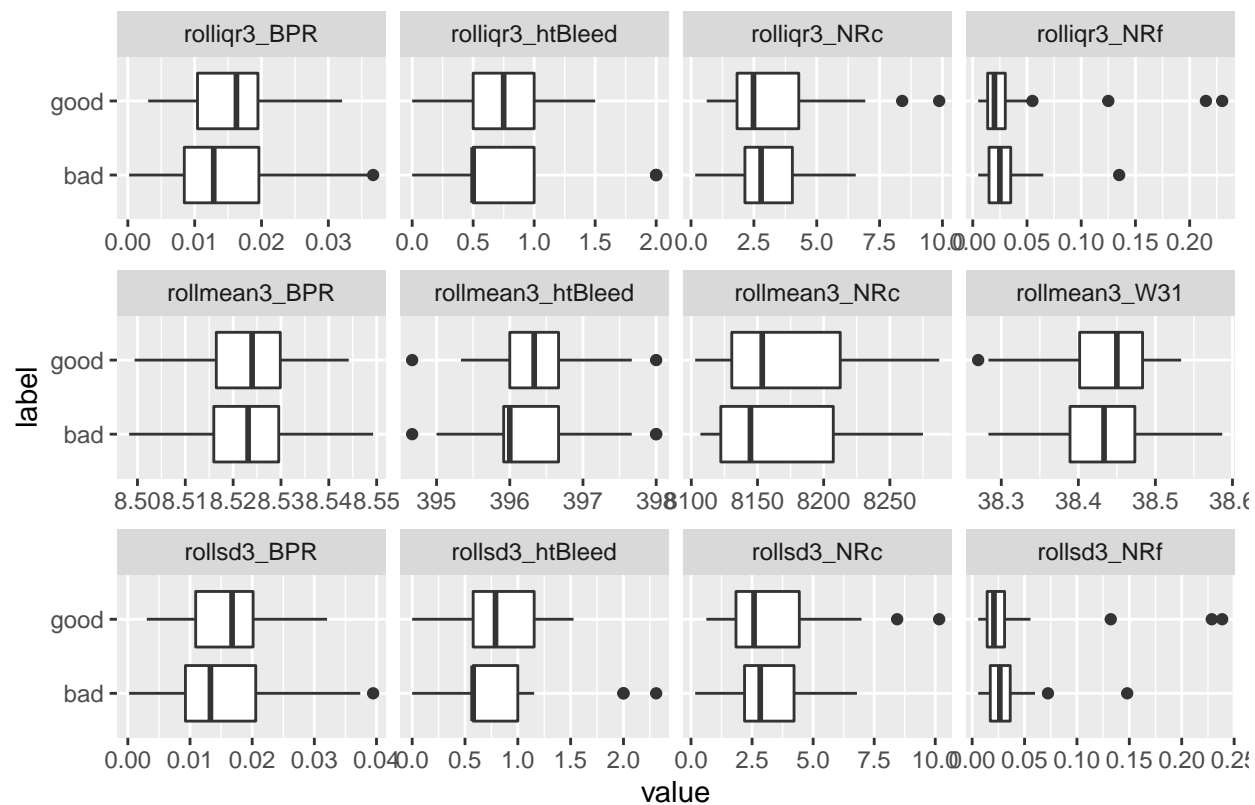


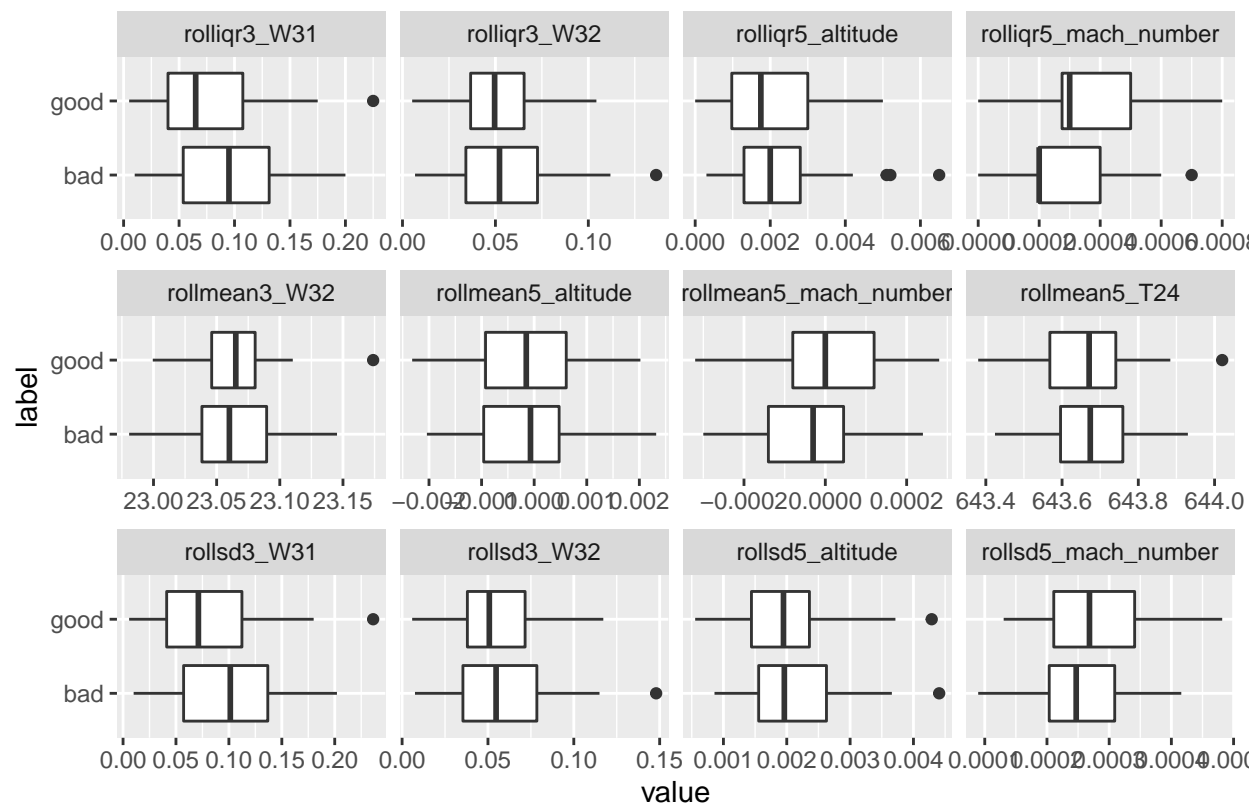


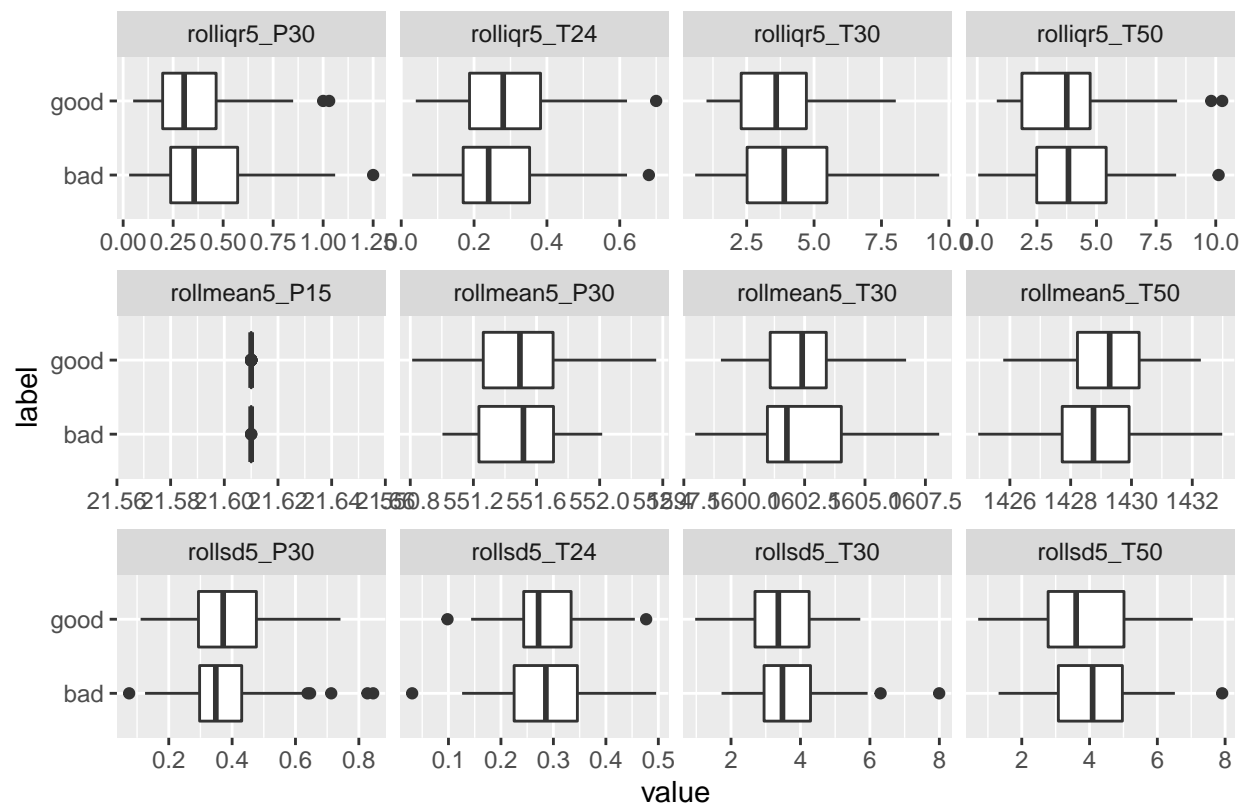


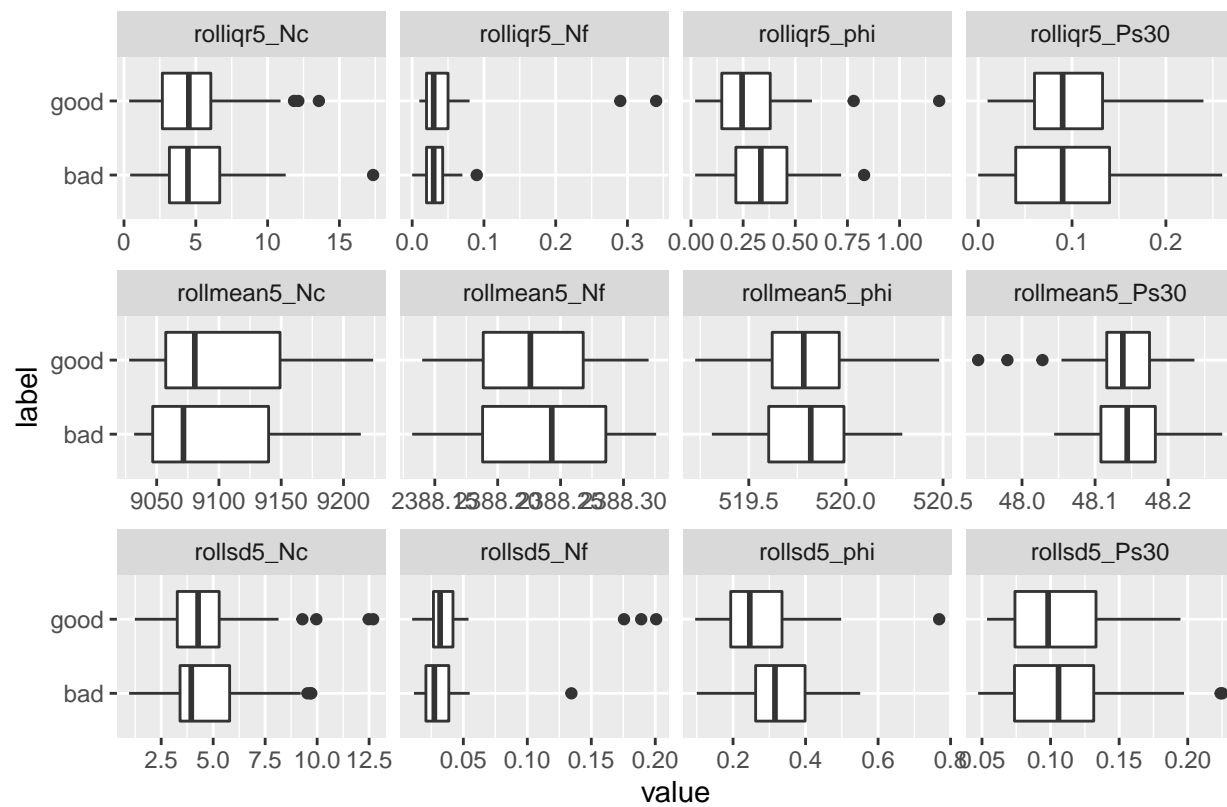


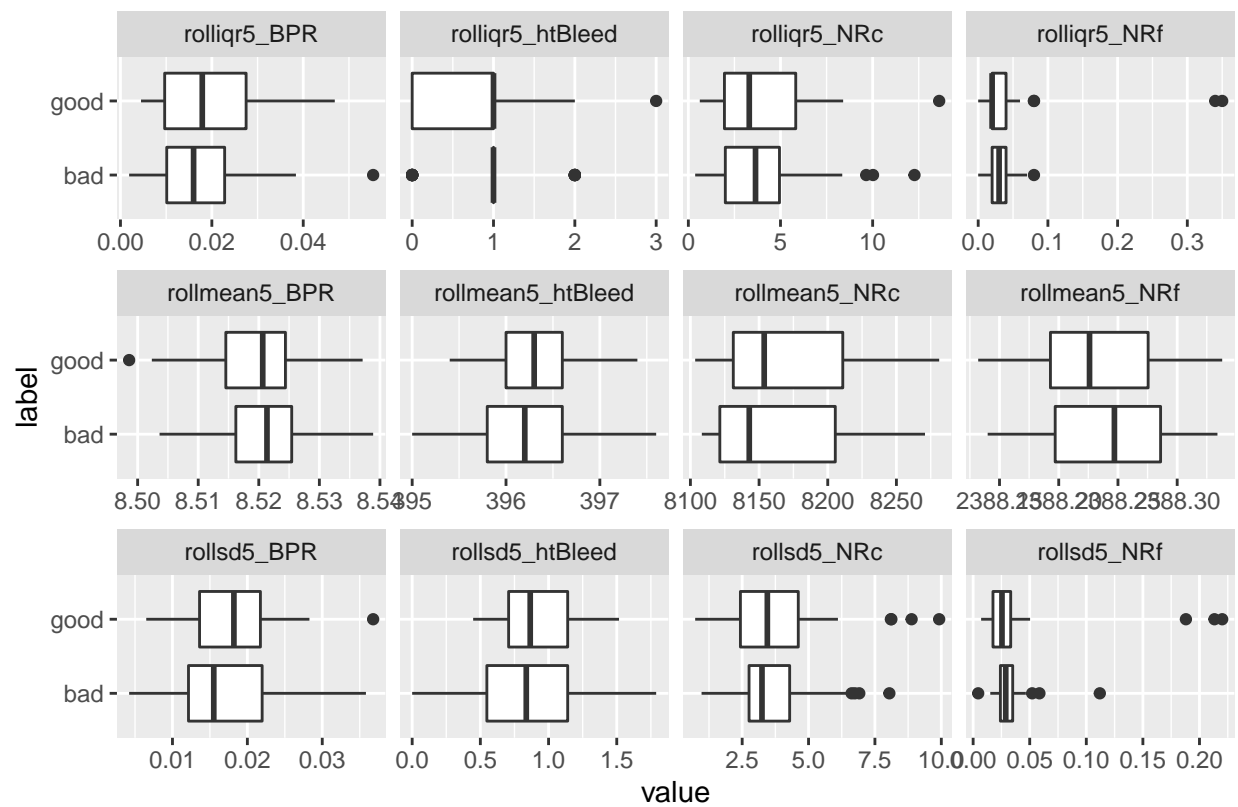


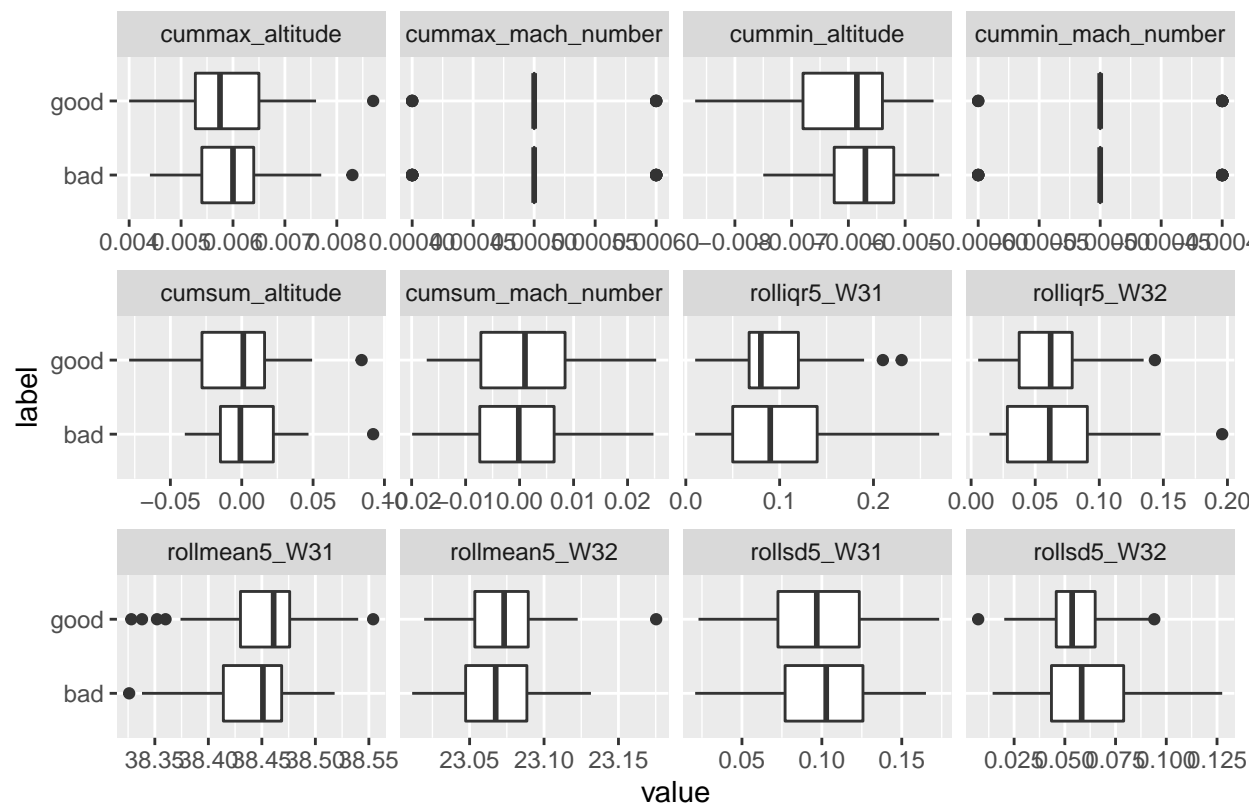


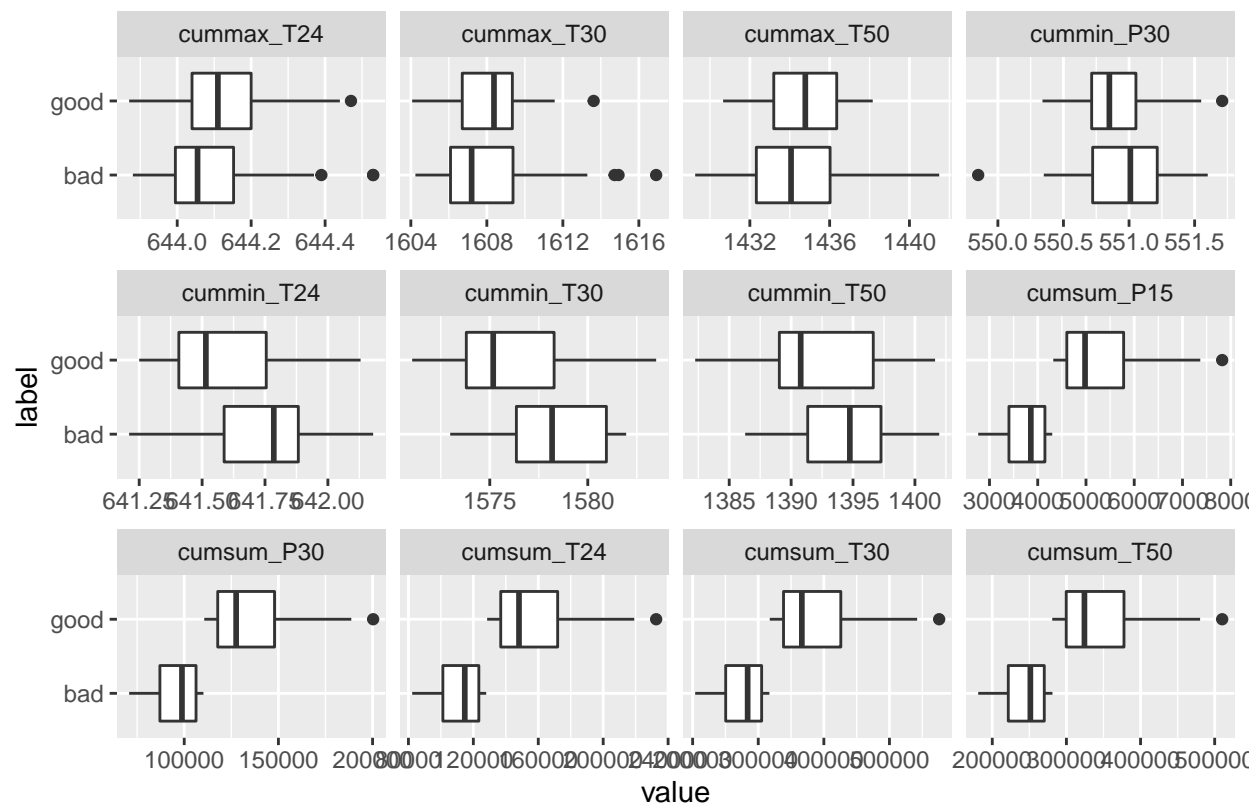


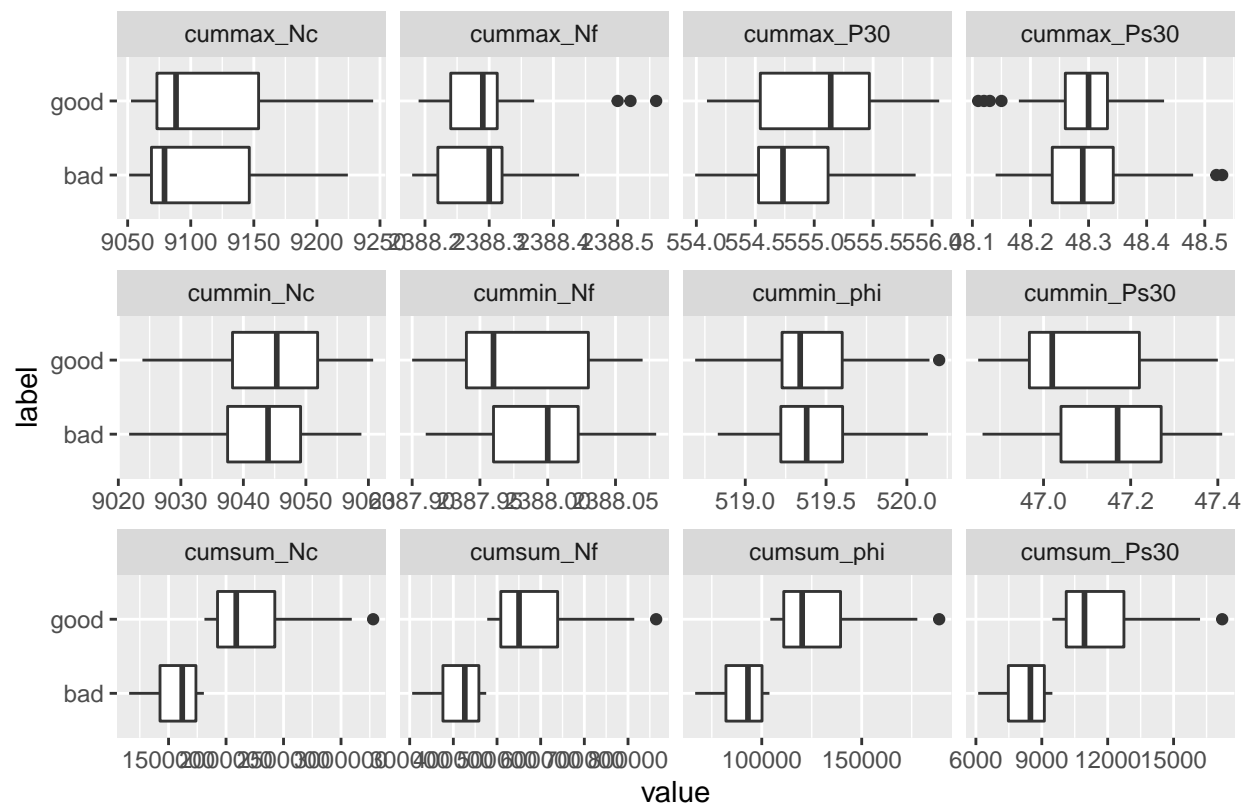


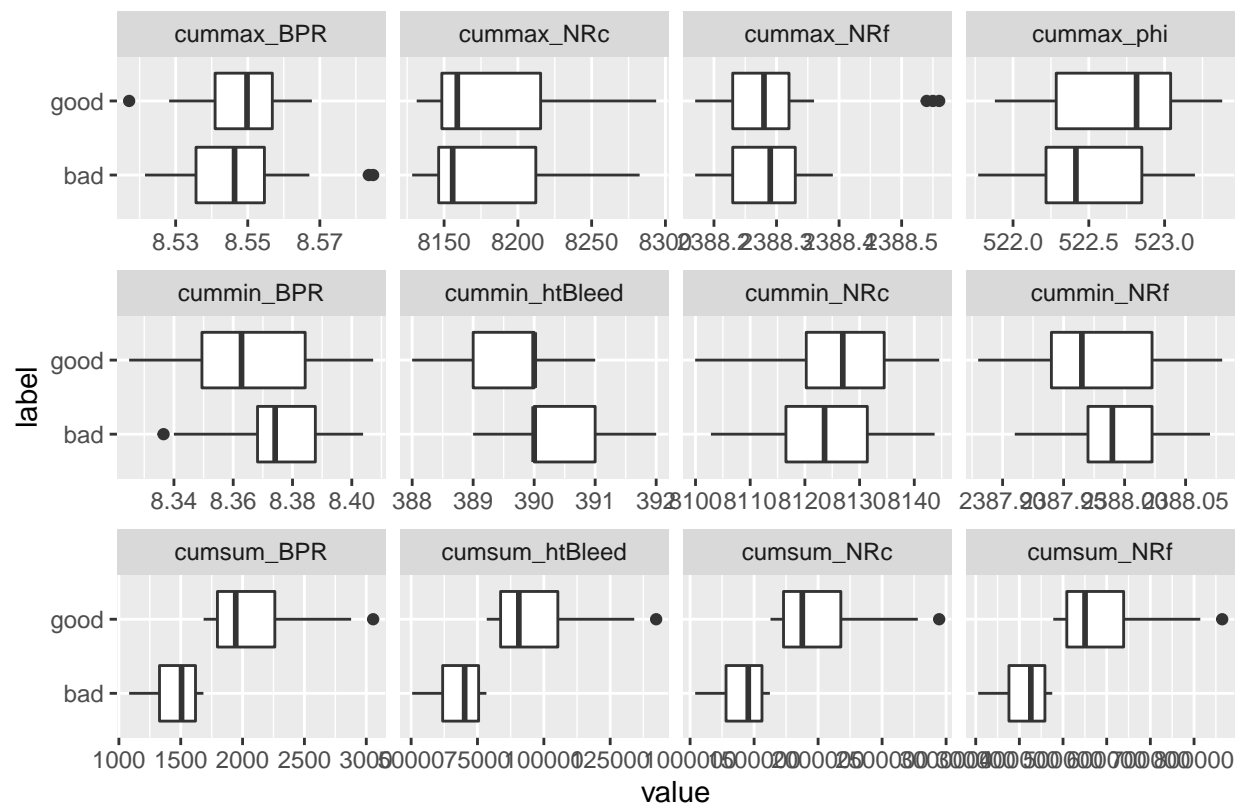


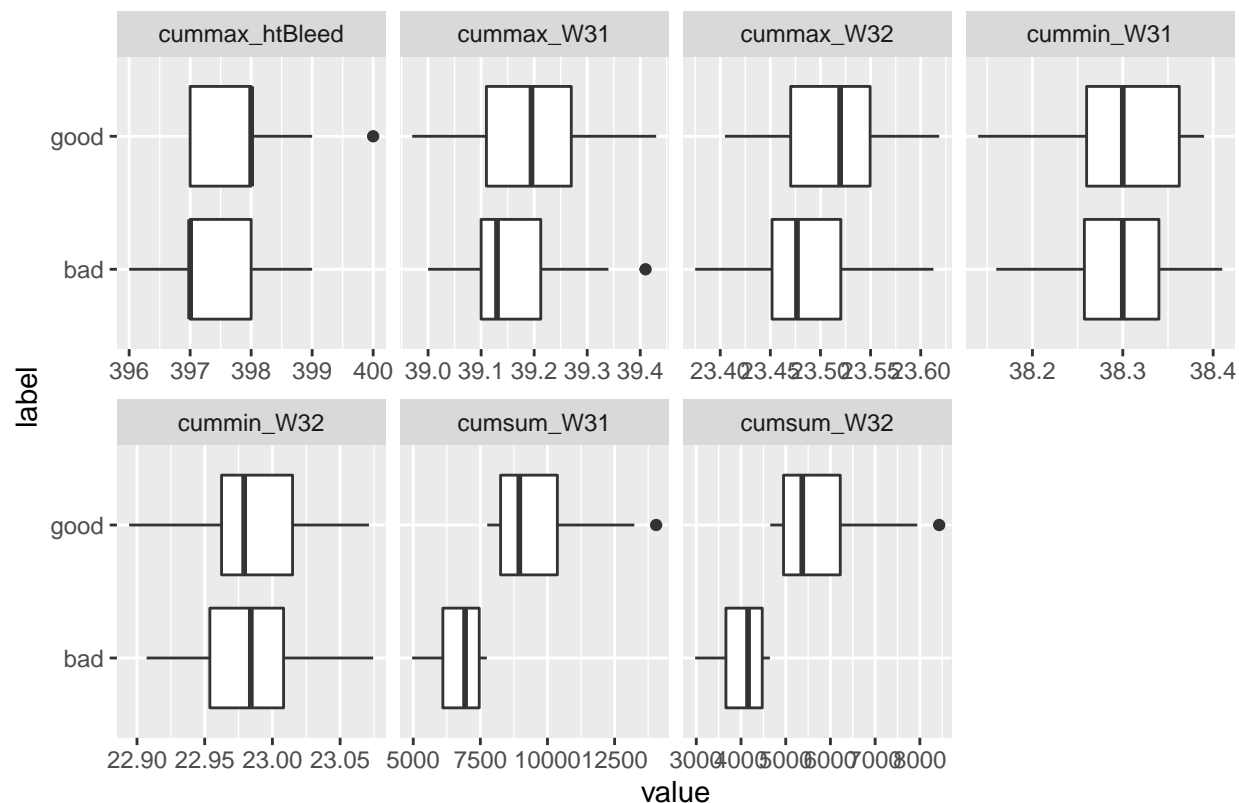












Page 32

4.2 Comprobación de la normalidad y homogeneidad de la varianza.

Primero, vamos a examinar si los datos siguen una distribución normal,

```
normalityChecker <- function(df){
  df$label = NULL
  result <- data.frame(variables = names(df), p_value = NA, result = NA)

  for (i in 1:nrow(result)) {
    variable <- df[,i]
    prueba <- shapiro.test(pull(variable))
    result$p_value[i] <- prueba$p.value
    result$result[i] <- if_else(prueba$p.value >= 0.05, "normally distributed", "not normally distributed")
  }

  return(result)
}

pruebas <- normalityChecker(df)

table(pruebas$result)

##
##      normally distributed not normally distributed
```

Ahora vamos a examinar la homocedasticidad.

Para las variables que siguen una distribución normal usaremos el **test de Levene** y el **test de Fligner-Killeen** para las columnas que no cumplen con la condición de normalidad.

```
homoscedasticityChecker <- function(df_in, normalitychecker_table){
  group = as.factor(df_in$label)
  df_in$label = NULL
  names(normalitychecker_table) <- c("variables", "normal_p_value", "normal_result")
  normalitychecker_table$homoscedasticity_p_value = NA
  normalitychecker_table$homoscedasticity_result = NA

  for (i in 1:nrow(normalitychecker_table)) {
    y <- pull(df_in[,i])
    temp_df <- data.frame(y = y, group = group)

    if(normalitychecker_table[i,3] == "normally distributed"){
      result_test <- leveneTest(y = y, group = group, data = temp_df)
      p_value = result_test$`Pr(>F)`[1]
    } else {
      result_test <- fligner.test(x = y, g = group, data = temp_df)
      p_value = result_test$p.value
    }
    normalitychecker_table$homoscedasticity_p_value[i] = p_value
    normalitychecker_table$homoscedasticity_result[i] = if_else(p_value < 0.05, "statistically different", "not statistically different")
  }
  return(normalitychecker_table)
}

pruebas2 <- homoscedasticityChecker(df = df, normalitychecker_table = pruebas)

table(pruebas2$homoscedasticity_result)
```

```
##
## not statistically different variances      statistically different variances
##                                     343                                     37
```

```
table(pruebas2$normal_result, pruebas2$homoscedasticity_result)
```

```
##
##                                not statistically different variances
## normally distributed                                179
## not normally distributed                                164
##
##                                statistically different variances
## normally distributed                                14
## not normally distributed                                23
```

4.3 Aplicación de pruebas estadísticas para comparar los grupos de datos.

Queremos saber si el promedio de nuestra variables para el grupo de máquinas buenas (m_A) es significativamente diferente al de las máquinas malas (m_B).

En este caso, tenemos dos grupos de muestras no relacionados (es decir, independientes o no apareados). Por lo tanto, es posible utilizar una prueba t independiente para evaluar si las medias son diferentes.

```
significantDifferenceChecker <- function(df_in, normalitychecker_table){

  group = as.factor(df_in$label)
  df_in$label = NULL
  normalitychecker_table$mean_difference_p_value = NA
  normalitychecker_table$mean_difference_result = NA

  for (i in 1:nrow(normalitychecker_table)) {
    y <- pull(df_in[,i])
    temp_df <- data.frame(y = y, group = group)

    if(normalitychecker_table[i,3] == "normally distributed" & normalitychecker_table[i,5] == "statistical") {
      result_test <- t.test(x = temp_df$y[group == "good"], y = temp_df$y[group == "bad"])
    } else {
      result_test <- wilcox.test(x = temp_df$y[group == "good"], y = temp_df$y[group == "bad"])
    }

    normalitychecker_table$mean_difference_p_value[i] = result_test$p.value
    normalitychecker_table$mean_difference_result[i] = if_else(result_test$p.value < 0.05, "statistically significant", "not statistically significant")
  }

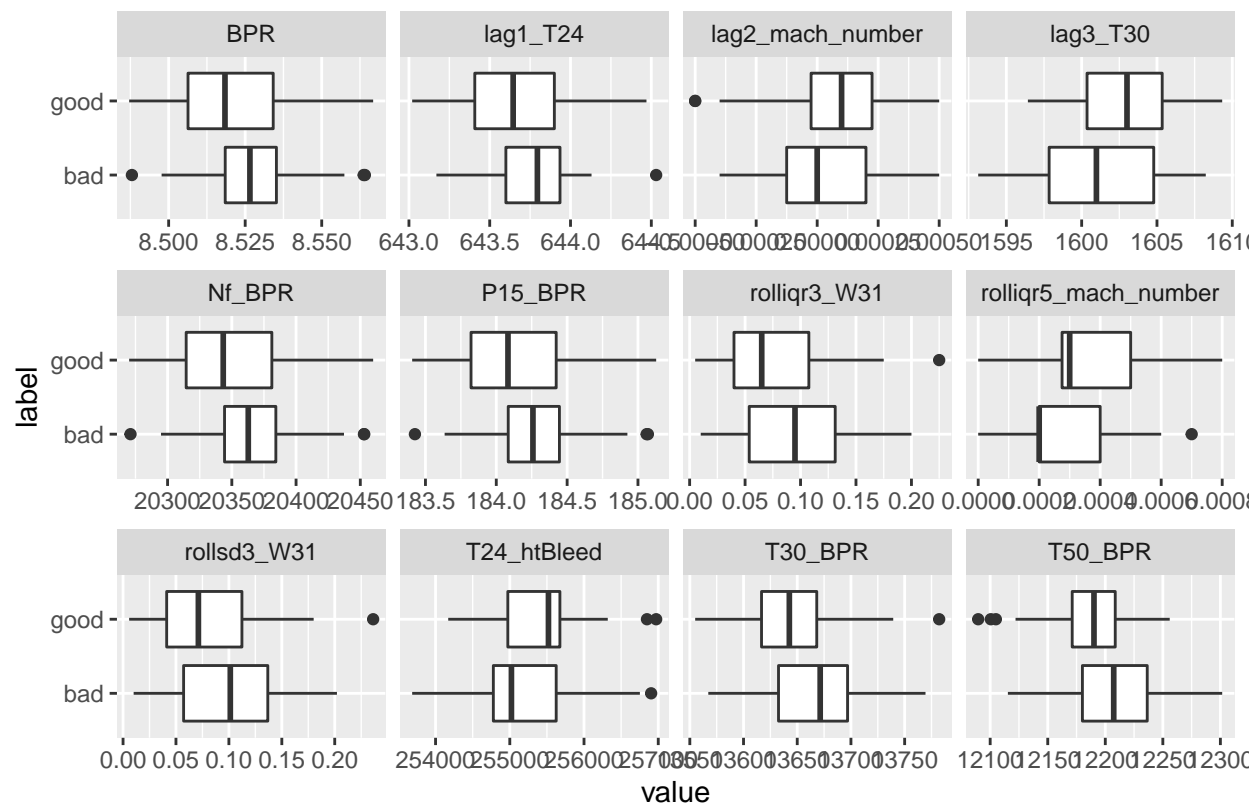
  return(normalitychecker_table)
}

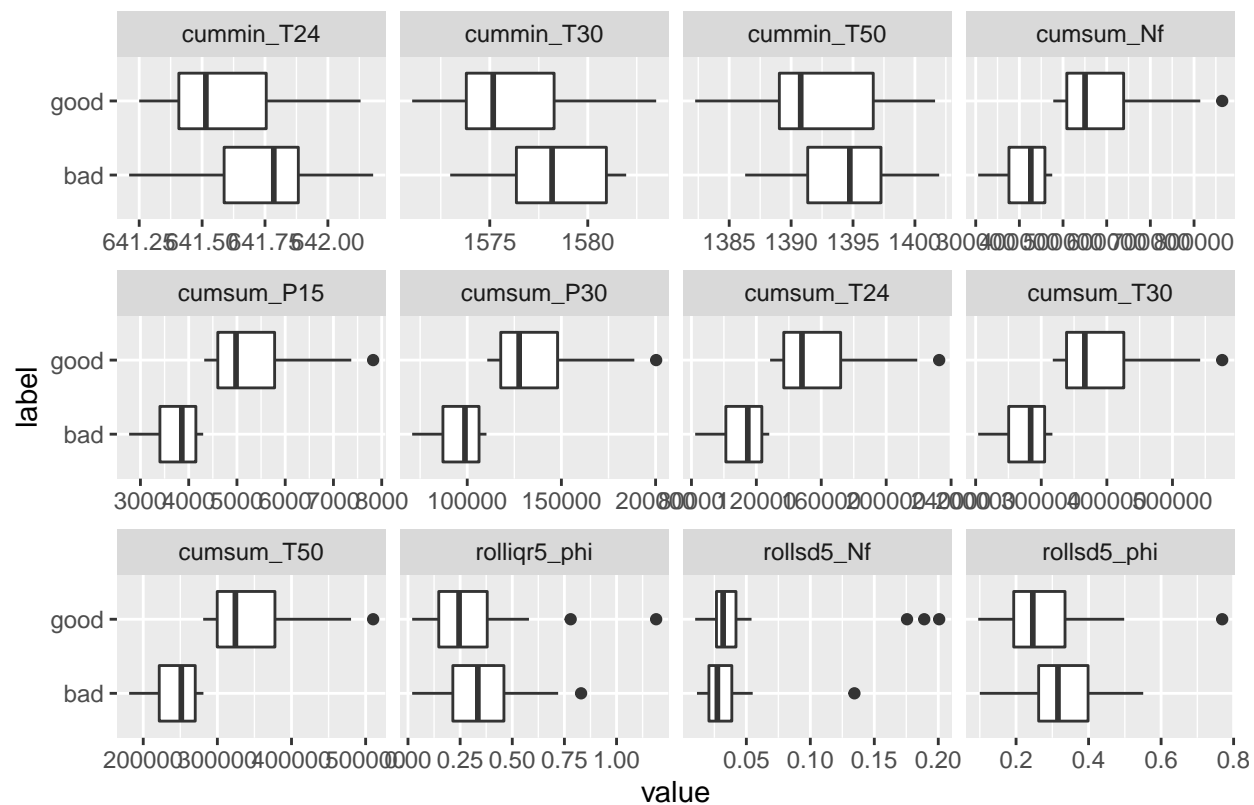
pruebas3 <- significantDifferenceChecker(df_in = df, normalitychecker_table = pruebas2)

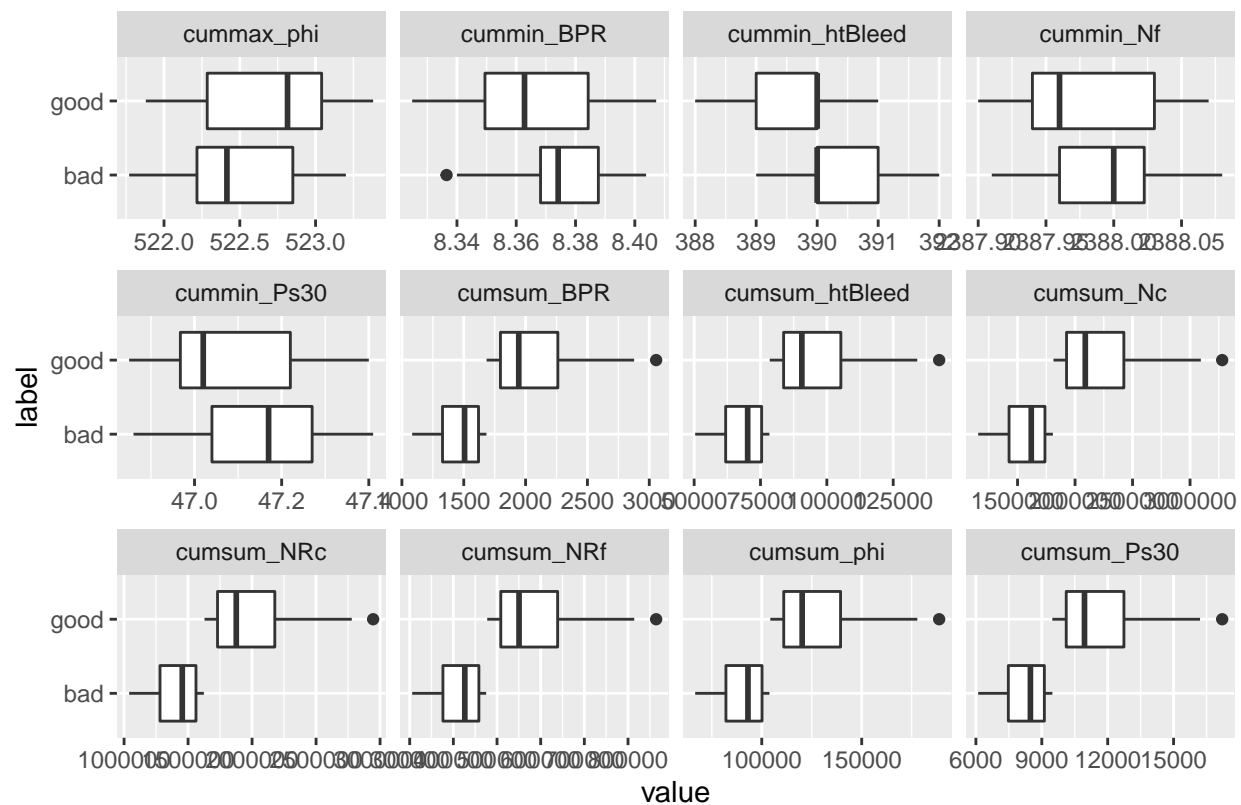
table(pruebas3$mean_difference_result)

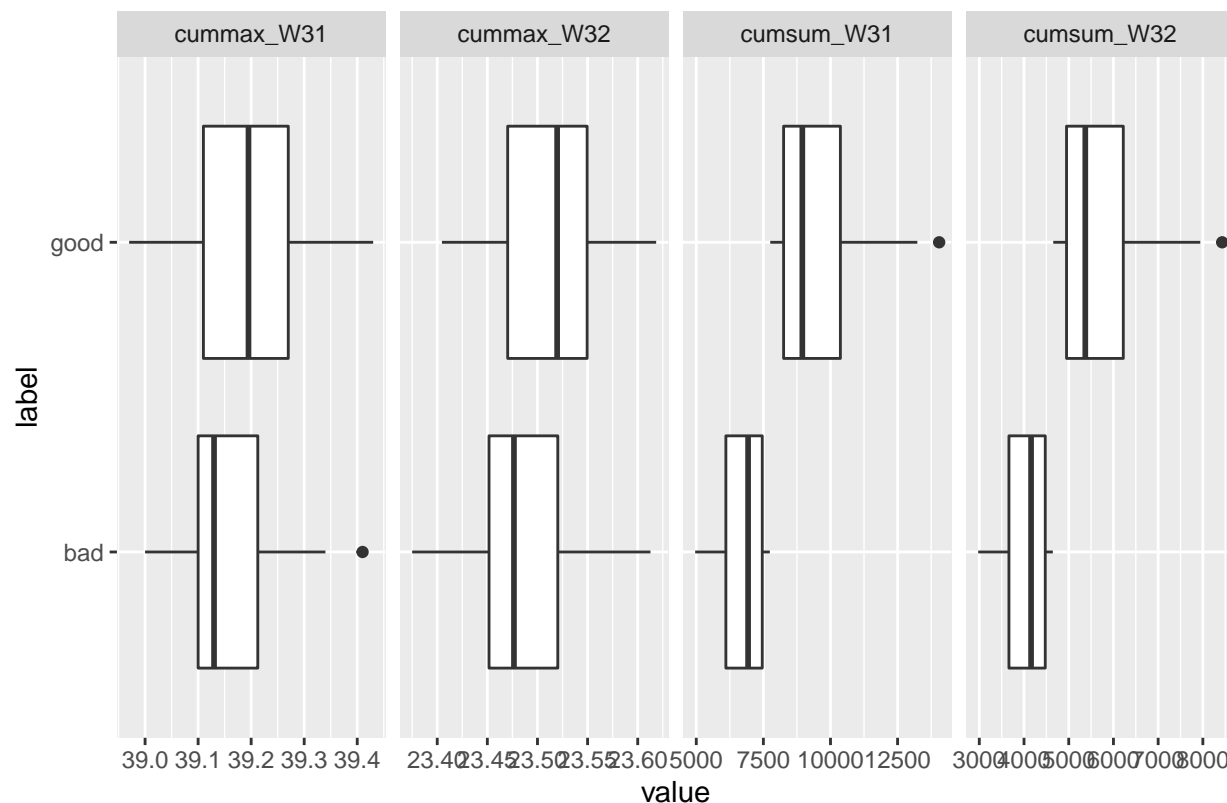
##
## not statistically different means      statistically different means
##                                339                                41

columns_significant <- pruebas3 %>% filter(mean_difference_result == "statistically different means") %>%
  plot_boxplot(df[,c(columns_significant,"label")], by = "label")
```









Page 4

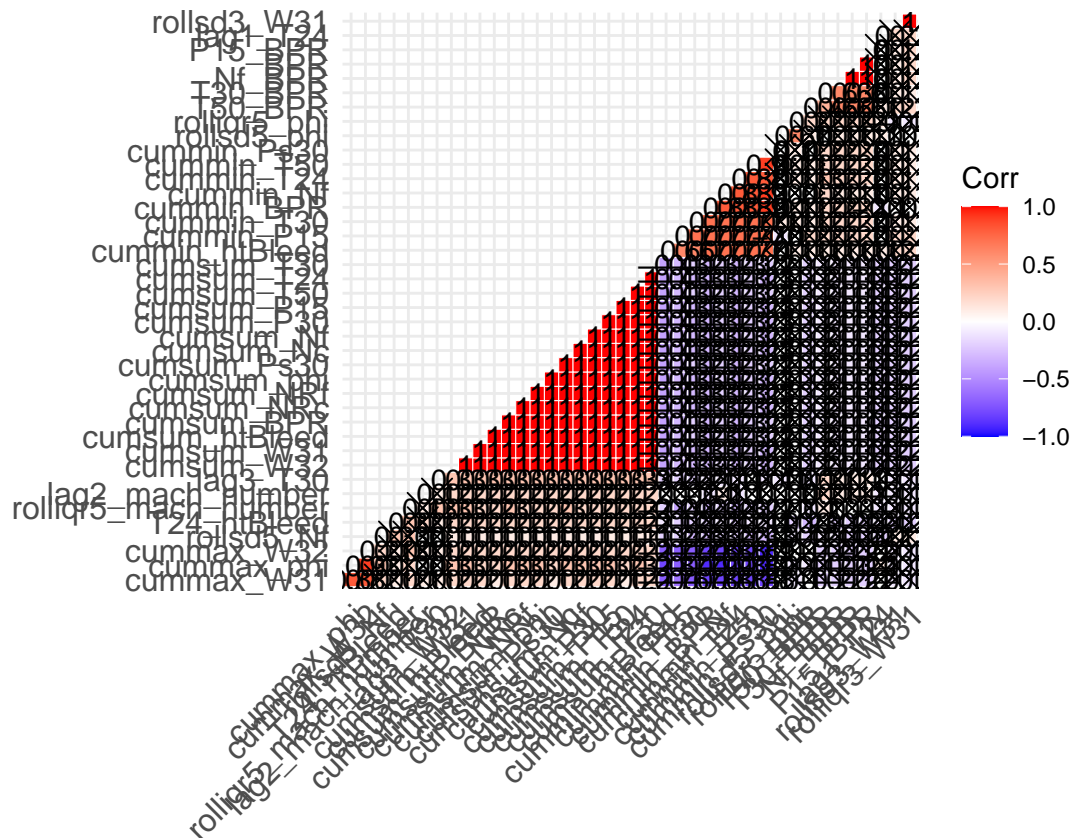
Ahora vamos a estudiar el problema de multicolinealidad o de cuando dos variables son altamente correlacionadas, al punto de ser una combinación lineal. Además, examinaremos si la correlación es significativamente diferente de cero. Usaremos la prueba de Spearman porque gran parte de las variables no son normalmente distribuidas.

```
df_temp <- df[,columns_significant]

# Compute a correlation matrix
corr <- round(cor(df_temp), 1)
# Compute a matrix of correlation p-values
p.mat <- cor_pmat(df_temp, method="spearman")

# Visualize the lower triangle of the correlation matrix
# Barring the no significant coefficient
corr.plot <- ggcorrplot(
  corr, hc.order = TRUE, type = "lower", lab = TRUE, outline.col = "white",
  p.mat = p.mat
)

corr.plot
```



El resultado indica que hay muchas correlaciones que no son estadísticamente diferente de cero y, por otro lado, algunas variables son combinación lineales. Este última condición sería un problema en un modelo de regresión lineal.

```
df_temp <- df[,c(columns_significant[!grepl('cumsum', columns_significant)], "label")]
df_temp$label <- as.factor(df_temp$label)

model <- glm(label ~., data = df_temp, family = binomial) %>%
  stepAIC(trace = FALSE)

# Summarize the final selected model
summary(model)
```

```
##
## Call:
## glm(formula = label ~ T30_BPR + lag2_mach_number + lag3_T30 +
##       rollsd3_W31 + rolliqr3_W31 + rollsd5_Nf + rollsd5_phi + cummin_T30 +
##       cummax_phi + cummin_htBleed, family = binomial, data = df_temp)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.9656  -0.7588  -0.1976   0.6786   2.3426
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  2785.571308 1038.209952   2.683  0.00730 **
```

```
## T30_BPR          -0.017296      0.006311  -2.741  0.00613 **
## lag2_mach_number 3426.358515 1208.903800   2.834  0.00459 **
## lag3_T30         0.114932      0.082391   1.395  0.16303
## rollsd3_W31      -126.389551   65.188692  -1.939  0.05252 .
## rolliqr3_W31     125.561192   67.570453   1.858  0.06314 .
## rollsd5_Nf       20.661018   12.321427   1.677  0.09357 .
## rollsd5_phi      -5.280365    2.381735  -2.217  0.02662 *
## cummin_T30       -0.315528    0.146278  -2.157  0.03100 *
## cummax_phi       -3.105803    1.304948  -2.380  0.01731 *
## cummin_htBleed   -1.567459    0.629747  -2.489  0.01281 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 138.469  on 99  degrees of freedom
## Residual deviance:  92.621  on 89  degrees of freedom
## AIC: 114.62
##
## Number of Fisher Scoring iterations: 5
```

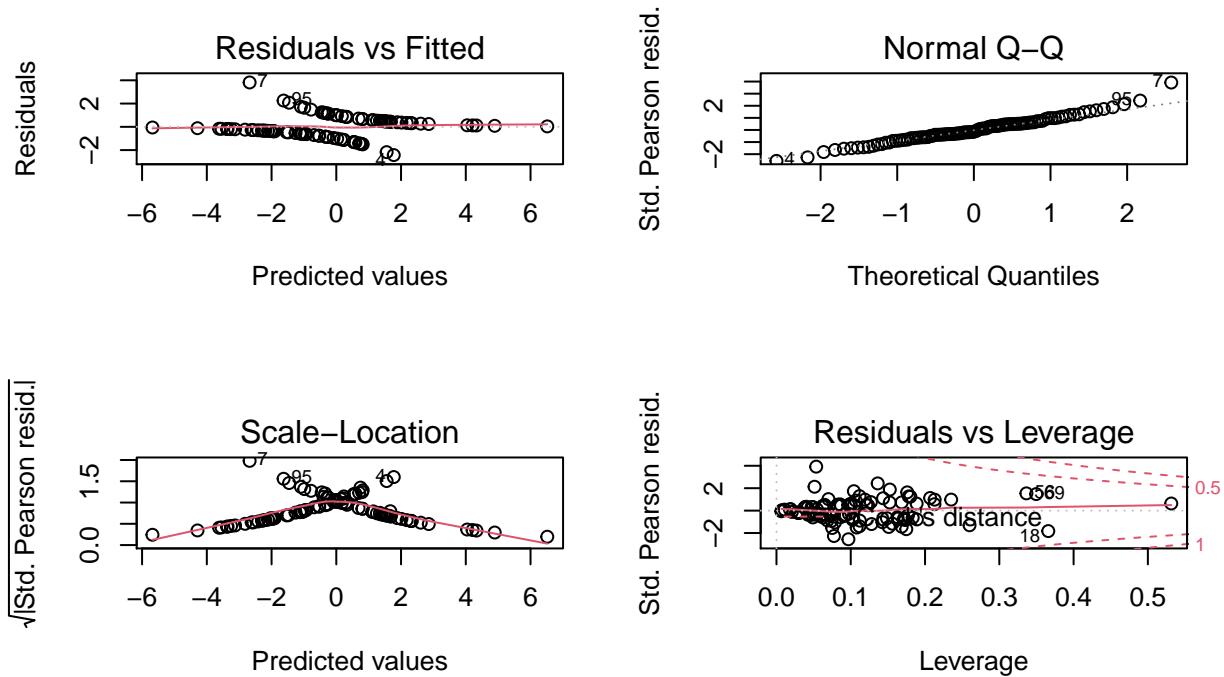
Hemos usado una selección de variable paso a paso basada en el criterio de información de akaike. Un modelo con un AIC bajo se caracteriza por una baja complejidad (minimiza p) y un buen ajuste. Las variables más importantes para explicar si una máquina tiene buen o malo rendimiento son:

- T30_BPR
- lag2_mach_number
- lag3_T30
- rollsd3_W31
- rolliqr3_W31
- rollsd5_Nf
- rollsd5_phi
- cummin_T30
- cummax_phi
- cummin_htBleed

5. Representación de los resultados a partir de tablas y gráficas.

```
oldpar <- par(oma = c(0, 0, 3, 0), mfrow = c(2, 2))
plot(model)
```

del ~ T30_BPR + lag2_mach_number + lag3_T30 + rollsd3_W31 + rolli



Vemos el intervalo de confianza de los parametros.

```
confint(model, level = 0.95)
```

```
##              2.5 %      97.5 %
## (Intercept)  839.08652419 4960.568982273
## T30_BPR      -0.03058919 -0.005583626
## lag2_mach_number 1176.20980745 5977.450565982
## lag3_T30      -0.04353586  0.283654081
## rollsd3_W31   -261.98460973 -3.264000799
## rolliqr3_W31  -2.31786819  265.908040671
## rollsd5_Nf     -0.22533057  51.641153057
## rollsd5_phi    -10.24300577 -0.756712201
## cummin_T30     -0.61570085 -0.036416478
## cummax_phi     -5.82876638 -0.650232964
## cummin_htBleed -2.88617946 -0.386898867
```

```
pred <- as.factor(round(predict(model, newdata = df_temp, type="response"),0))

observaciones <- as.factor(if_else(df$label=="bad", 0, 1))

matriz <- confusionMatrix(observaciones, pred)
matriz
```

```
## Confusion Matrix and Statistics
```

```

##
##           Reference
## Prediction  0  1
##           0 39 13
##           1 13 35
##
##           Accuracy : 0.74
##           95% CI : (0.6427, 0.8226)
##           No Information Rate : 0.52
##           P-Value [Acc > NIR] : 0.000005478
##
##           Kappa : 0.4792
##
## Mcnemar's Test P-Value : 1
##
##           Sensitivity : 0.7500
##           Specificity : 0.7292
##           Pos Pred Value : 0.7500
##           Neg Pred Value : 0.7292
##           Prevalence : 0.5200
##           Detection Rate : 0.3900
##           Detection Prevalence : 0.5200
##           Balanced Accuracy : 0.7396
##
##           'Positive' Class : 0
##

```

6. Resolución del problema. A partir de los resultados obtenidos, ¿cuáles son las

conclusiones? ¿Los resultados permiten responder al problema?

Nuestra conclusión es que la información contenida en las variables: lag2_mach_number lag3_T30
rollsd3_W31
rolliq3_W31
rollsd5_Nf
rollsd5_phi
cummin_T30
cummax_phi
cummin_htBleed

Consigue explicar un 75% sobre cuando una máquina tendrá una vida larga o corta. Nuestro sample size es de 100 máquinas, nuestro modelo es lineal y predecimos sobre el mismo dato usado en el entrenamiento, lo que puede causar overfitting (parcialmente ignorado porque regresión logística es bastante limitada en generar overfitting).

Este resultado no responde realmente el problema sobre los factores que determinan una vida larga en las máquinas, pero sirve para indicar factores importantes que pueden ser explorados con nuevas variables, modelos más potentes y diferentes “problem framings”.

Futuramente continuaremos esta investigación con los otros datasets y tres “problem framings” distintos:

- (1) Como un problema de clasificación para predecir el riesgo de la máquina fallar en su siguiente ciclo;
- (2) Un análisis de supervivencia;

- (3) Un modelos de regresión para predecir el **Remaining Useful Life**