

# models

Tao

2023-01-20

```
library("dplyr")

##
## Attaching package: 'dplyr'
## The following objects are masked from 'package:stats':
##
##   filter, lag
## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union
#install.packages('corrplot')
library(corrplot)

## corrplot 0.92 loaded
library(RColorBrewer)
# install.packages("gbm")
library("gbm")

## Loaded gbm 2.1.8
# install.packages("caret")
library("caret")

## Loading required package: ggplot2
## Loading required package: lattice
#install.packages("pdp")
library("pdp")          # model visualization
library("ggplot2")      # model visualization
#install.packages("lime")
library("lime")         # model visualization

##
## Attaching package: 'lime'
## The following object is masked from 'package:dplyr':
##
##   explain
library("pROC")

## Type 'citation("pROC")' for a citation.
##
```

```

## Attaching package: 'pROC'

## The following objects are masked from 'package:stats':
##
##      cov, smooth, var
#install.packages("e1071", repos="http://R-Forge.R-project.org")
library("e1071")
library( "MASS" )      #      used to generate correlated variables

##
## Attaching package: 'MASS'

## The following object is masked from 'package:dplyr':
##
##      select
library("sp")
library("Hmisc")      #      used for graphing se bars

## Loading required package: survival

##
## Attaching package: 'survival'

## The following object is masked from 'package:caret':
##
##      cluster

## Loading required package: Formula

##
## Attaching package: 'Hmisc'

## The following object is masked from 'package:e1071':
##
##      impute

## The following objects are masked from 'package:dplyr':
##
##      src, summarize

## The following objects are masked from 'package:base':
##
##      format.pval, units
#install.packages("randomForest")
require("randomForest")

## Loading required package: randomForest

## randomForest 4.7-1.1

## Type rfNews() to see new features/changes/bug fixes.

##
## Attaching package: 'randomForest'

## The following object is masked from 'package:ggplot2':
##
##      margin

```

```
## The following object is masked from 'package:dplyr':
##
##      combine

#install.packages("e1071")
library(e1071)
library(caret)
library("ModelMetrics")

##
## Attaching package: 'ModelMetrics'

## The following object is masked from 'package:pROC':
##
##      auc

## The following objects are masked from 'package:caret':
##
##      confusionMatrix, precision, recall, sensitivity, specificity

## The following object is masked from 'package:base':
##
##      kappa

library("foreign")
#install.packages("rfUtilities")
library("rfUtilities")

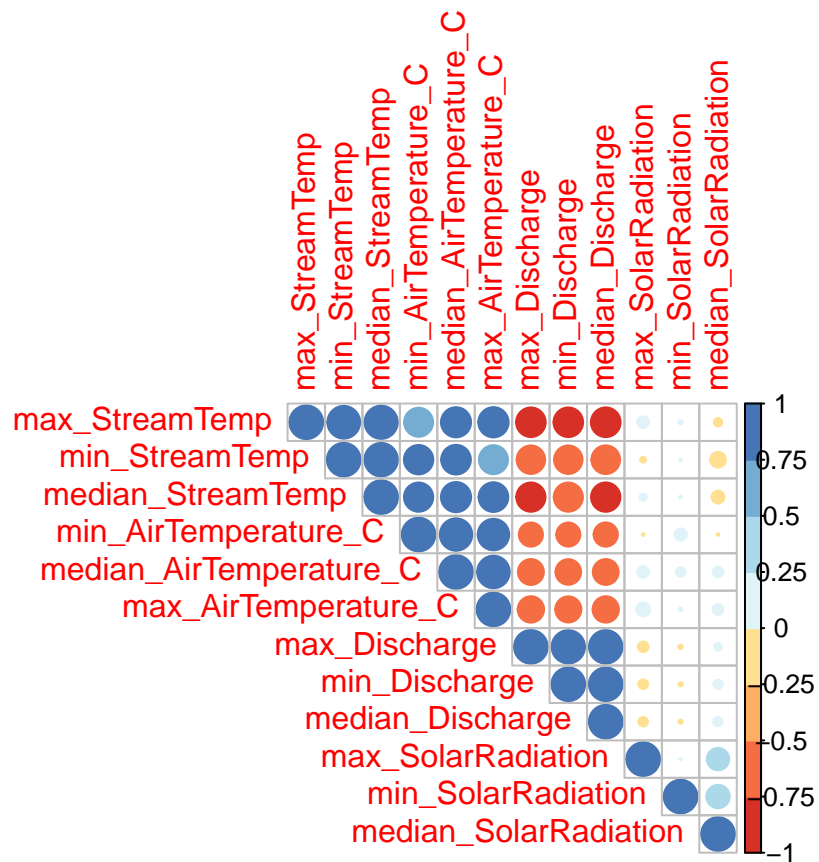
##
## Attaching package: 'rfUtilities'

## The following object is masked from 'package:ModelMetrics':
##
##      logLoss
```

## Load data

```
#load("daily_df_summer.Rdata")
load("LowerWeather_daily_df_summer.Rdata")

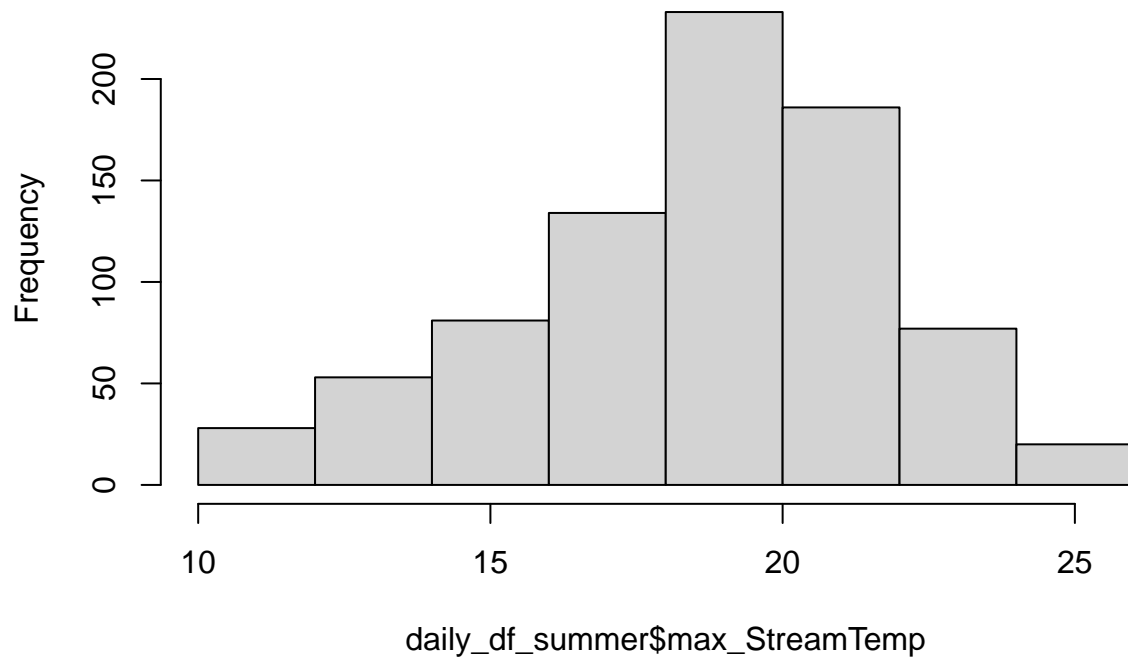
M <-cor(daily_df_summer[,c(2:13)])
corrplot(M, type="upper", order="hclust",
         col=brewer.pal(n=8, name="RdYlBu"))
```



```
#stream T, Air T, DISCHARGE
```

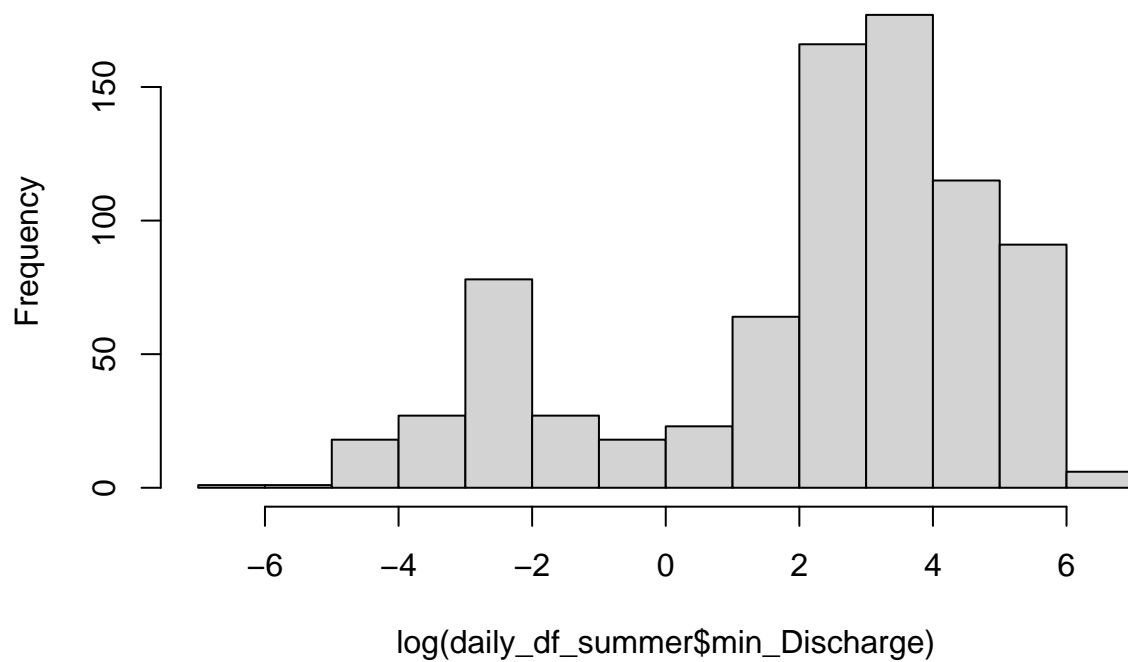
```
hist(daily_df_summer$max_StreamTemp)
```

**Histogram of daily\_df\_summer\$max\_StreamTemp**



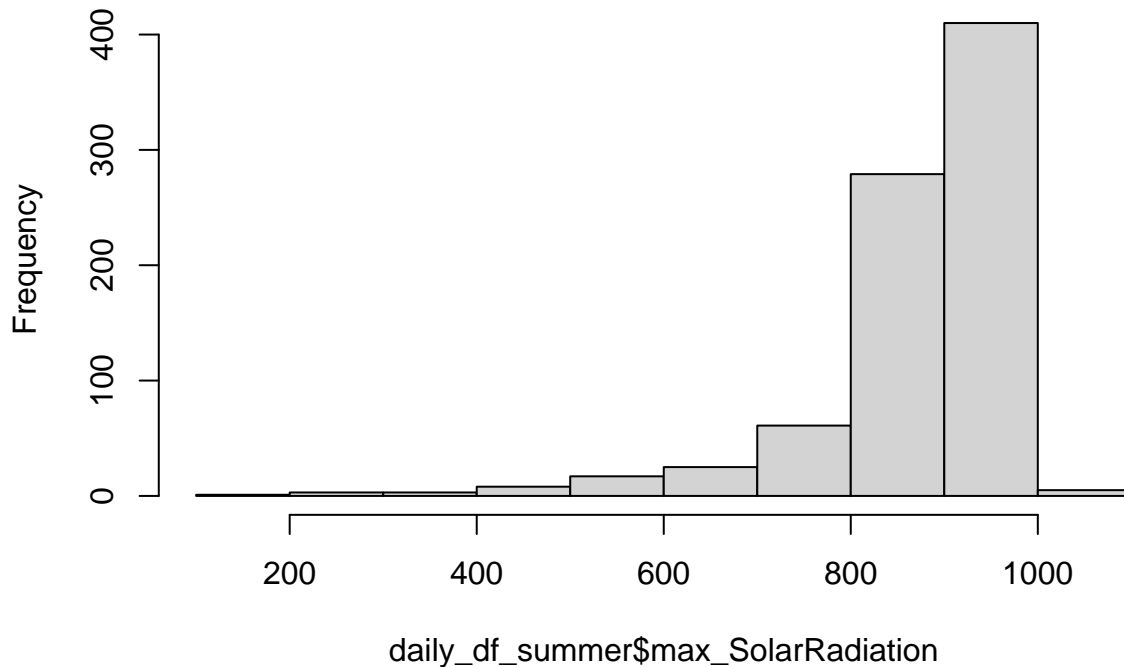
```
hist(log(daily_df_summer$min_Discharge))
```

**Histogram of log(daily\_df\_summer\$min\_Discharge)**



```
hist(daily_df_summer$max_SolarRadiation)
```

## Histogram of daily\_df\_summer\$max\_SolarRadiation



```
summary(lm(daily_df_summer$max_StreamTemp~ daily_df_summer$min_Discharge + daily_df_summer$max_AirTemper
```

```
##
## Call:
## lm(formula = daily_df_summer$max_StreamTemp ~ daily_df_summer$min_Discharge +
##     daily_df_summer$max_AirTemperature_C + daily_df_summer$max_SolarRadiation)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.1533 -1.0511 -0.2339  0.8206  4.9264
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    10.1207674   0.5394640   18.761  <2e-16 ***
## daily_df_summer$min_Discharge -0.0164048   0.0007577  -21.651  <2e-16 ***
## daily_df_summer$max_AirTemperature_C  0.3295152   0.0127151   25.915  <2e-16 ***
## daily_df_summer$max_SolarRadiation -0.0001829   0.0004863   -0.376    0.707
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.476 on 808 degrees of freedom
## Multiple R-squared:  0.7803, Adjusted R-squared:  0.7795
## F-statistic: 956.8 on 3 and 808 DF,  p-value: < 2.2e-16
```

```
summary(lm(daily_df_summer$max_StreamTemp~ log(daily_df_summer$min_Discharge) + daily_df_summer$max_Air
```

```
##
## Call:
## lm(formula = daily_df_summer$max_StreamTemp ~ log(daily_df_summer$min_Discharge) +
##     daily_df_summer$max_AirTemperature_C + daily_df_summer$max_SolarRadiation)
```

```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.0969 -0.9133  0.0992  0.9018  4.6423
##
## Coefficients:
##                                Estimate Std. Error t value Pr(>|t|)
## (Intercept)                   7.5512207   0.4785185  15.780 < 2e-16 ***
## log(daily_df_summer$min_Discharge) -0.4891809   0.0208759 -23.433 < 2e-16 ***
## daily_df_summer$max_AirTemperature_C 0.3660805   0.0112962  32.407 < 2e-16 ***
## daily_df_summer$max_SolarRadiation  0.0017420   0.0004806   3.625 0.000307 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.431 on 808 degrees of freedom
## Multiple R-squared:  0.7933, Adjusted R-squared:  0.7926
## F-statistic: 1034 on 3 and 808 DF, p-value: < 2.2e-16
summary(lm(daily_df_summer$max_StreamTemp~ log(daily_df_summer$min_Discharge) + daily_df_summer$max_Air'

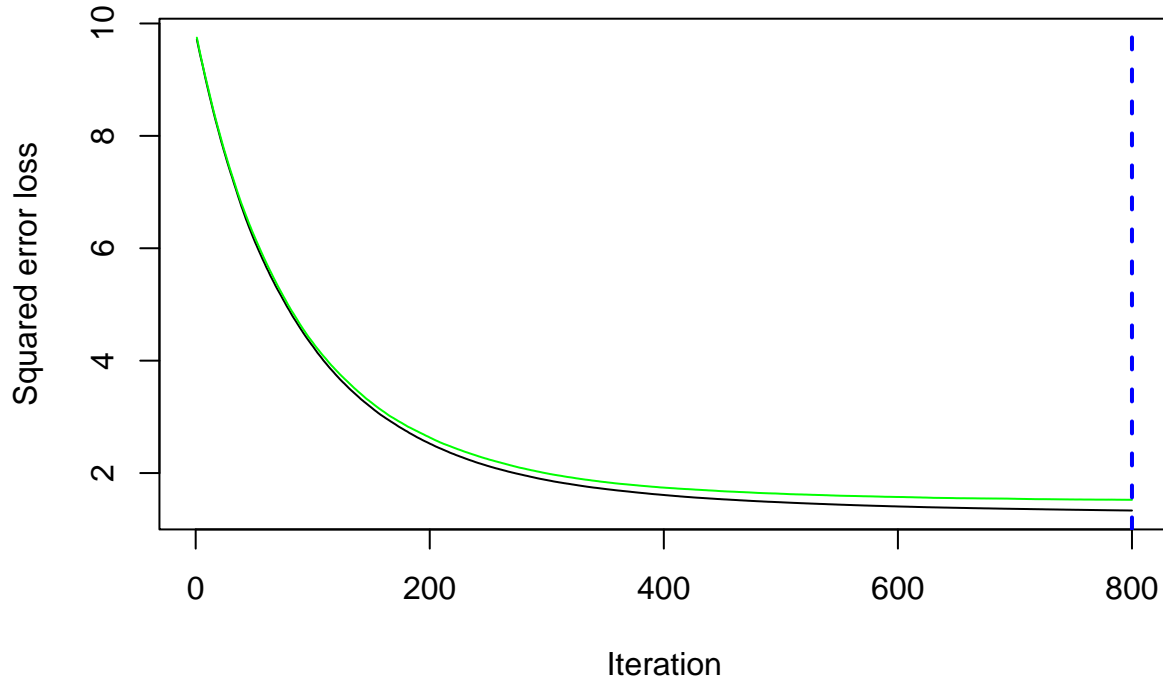
##
## Call:
## lm(formula = daily_df_summer$max_StreamTemp ~ log(daily_df_summer$min_Discharge) +
##     daily_df_summer$max_AirTemperature_C)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.0945 -0.9297  0.0771  0.9220  4.3511
##
## Coefficients:
##                                Estimate Std. Error t value Pr(>|t|)
## (Intercept)                   8.75027   0.34835   25.12 <2e-16 ***
## log(daily_df_summer$min_Discharge) -0.47460   0.02064  -23.00 <2e-16 ***
## daily_df_summer$max_AirTemperature_C 0.37578   0.01106   33.99 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.442 on 809 degrees of freedom
## Multiple R-squared:  0.79, Adjusted R-squared:  0.7895
## F-statistic: 1521 on 2 and 809 DF, p-value: < 2.2e-16
# set seed for generating random data.
set.seed(0)
# createDataPartition() function from the caret package to split the original dataset into a training a
variables<-c("max_StreamTemp","min_Discharge","max_AirTemperature_C", "max_SolarRadiation")
parts = createDataPartition( daily_df_summer$max_StreamTemp , p = .8, list = F)
train = daily_df_summer[parts, variables ]
test = daily_df_summer[-parts, variables ]
# feature and target array
test_x = test[, -1]
test_y = test[, 1]

model_gbm = gbm(train$max_StreamTemp ~.,
                 data = train,
                 distribution = "gaussian",
```

```
cv.folds = 10,
shrinkage = .01,
n.minobsinnode = 10,
n.trees = 800)
```

```
# model performance
```

```
perf_gbm1 = gbm.perf( model_gbm, method = "cv")
```



```
print(model_gbm)
```

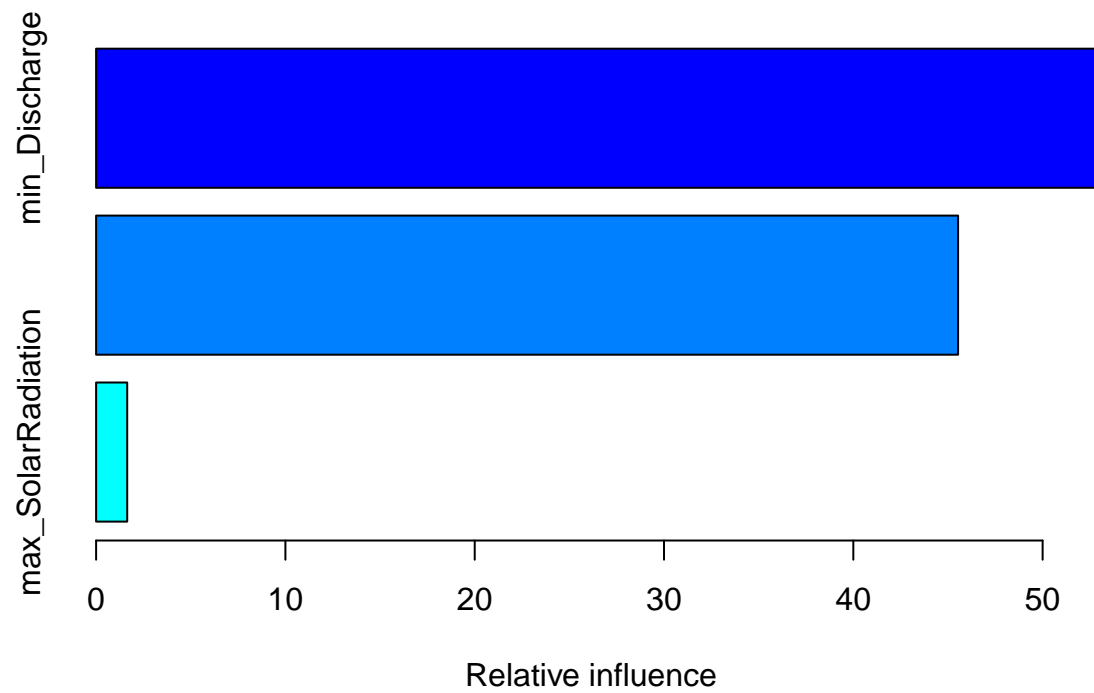
```
## gbm(formula = train$max_StreamTemp ~ ., distribution = "gaussian",
##      data = train, n.trees = 800, n.minobsinnode = 10, shrinkage = 0.01,
##      cv.folds = 10)
## A gradient boosted model with gaussian loss function.
## 800 iterations were performed.
## The best cross-validation iteration was 800.
## There were 3 predictors of which 3 had non-zero influence.
```

```
summary(model_gbm)
```

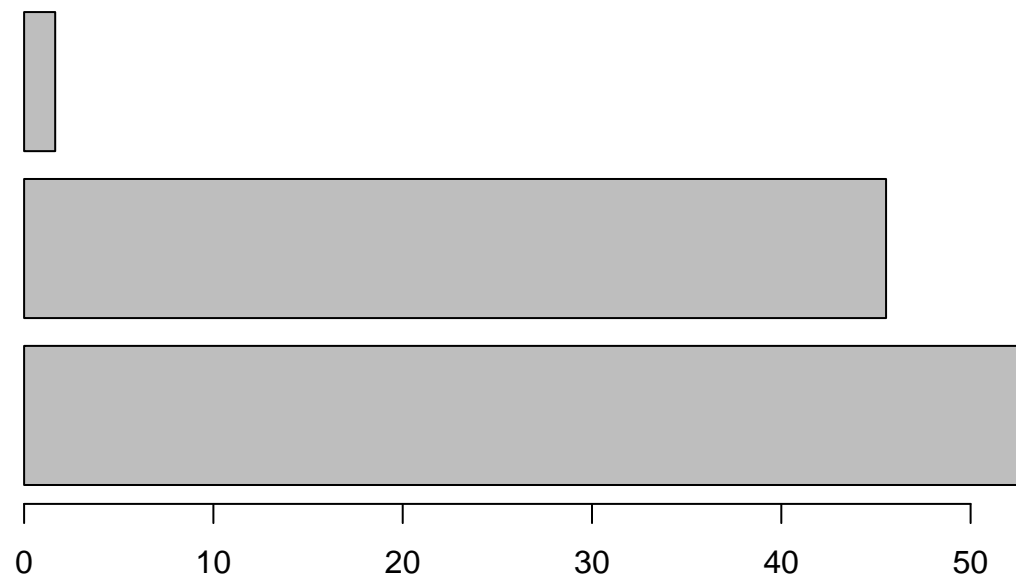
```
##                var    rel.inf
## min_Discharge      min_Discharge 52.824348
## max_AirTemperature_C max_AirTemperature_C 45.533356
## max_SolarRadiation    max_SolarRadiation  1.642296
```

```
rinf<-summary(model_gbm)
```

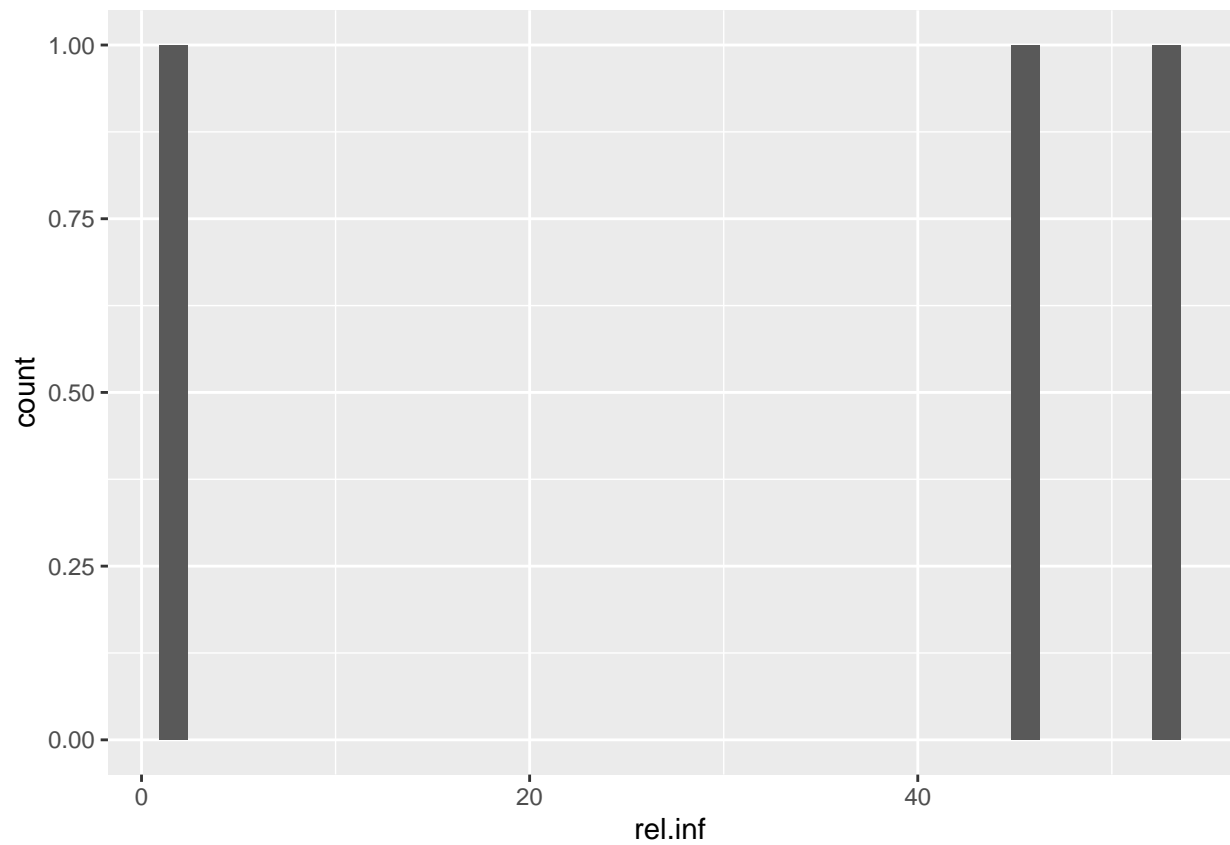




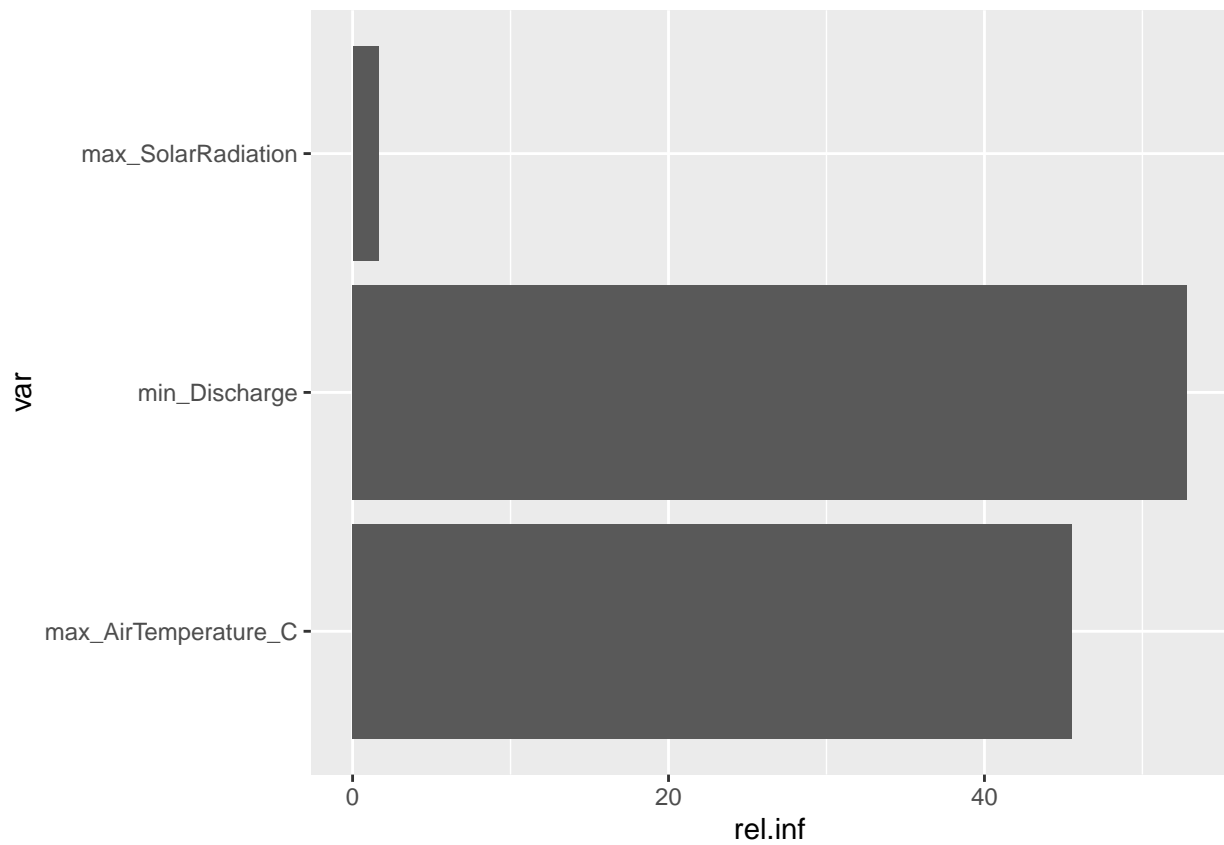
```
barplot( rinf$rel.inf , horiz = TRUE, las = 1)
```



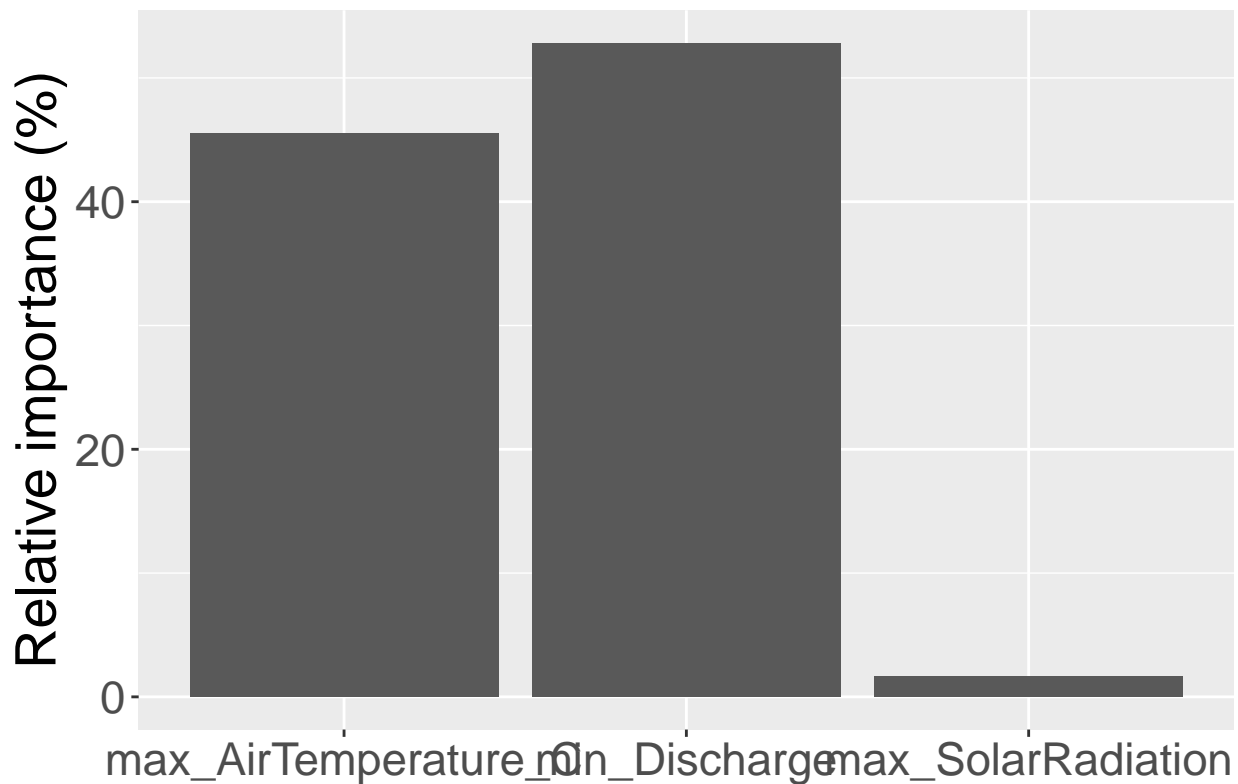
```
ggplot(rinf, aes(rel.inf)) + geom_bar()
```



```
rinf$var<- factor(rinf$var, levels=c( "max_AirTemperature_C" ,"min_Discharge" , "max_SolarRadiat" ))
ggplot( rinf, aes( var , rel.inf ))+ geom_col()+
coord_flip()
```



```
ggplot( rinf )+ geom_bar( aes( x=var, y= rel.inf), stat = "summary")+ scale_x_discrete(labels= c( "r", "i", "n", "f" ))
## No summary function supplied, defaulting to `mean_se()`
```



```
pred_y = predict.gbm(model_gbm, test_x)
```

```
## Using 800 trees...
```

```
residuals = test_y$max_StreamTemp - pred_y
summary(test_y$max_StreamTemp)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  10.50   16.70   19.10   18.63   20.70   26.00
```

```
xlim=c(5,30)
```

```
RMSE = sqrt(mean(residuals^2))
```

```
cat('The root mean square error of the test data is ', round(RMSE,3), '\n')
```

```
## The root mean square error of the test data is  1.271
```

```
y_test_mean = mean( test_y$max_StreamTemp )
```

```
# Calculate total sum of squares
```

```
tss = sum(( test_y$max_StreamTemp - y_test_mean)^2 )
```

```
# Calculate residual sum of squares
```

```
rss = sum(residuals^2)
```

```
# Calculate R-squared
```

```
rsq = 1 - (rss/tss)
```

```
cat('The R-square of the test data is ', round(rsq,3), '\n')
```

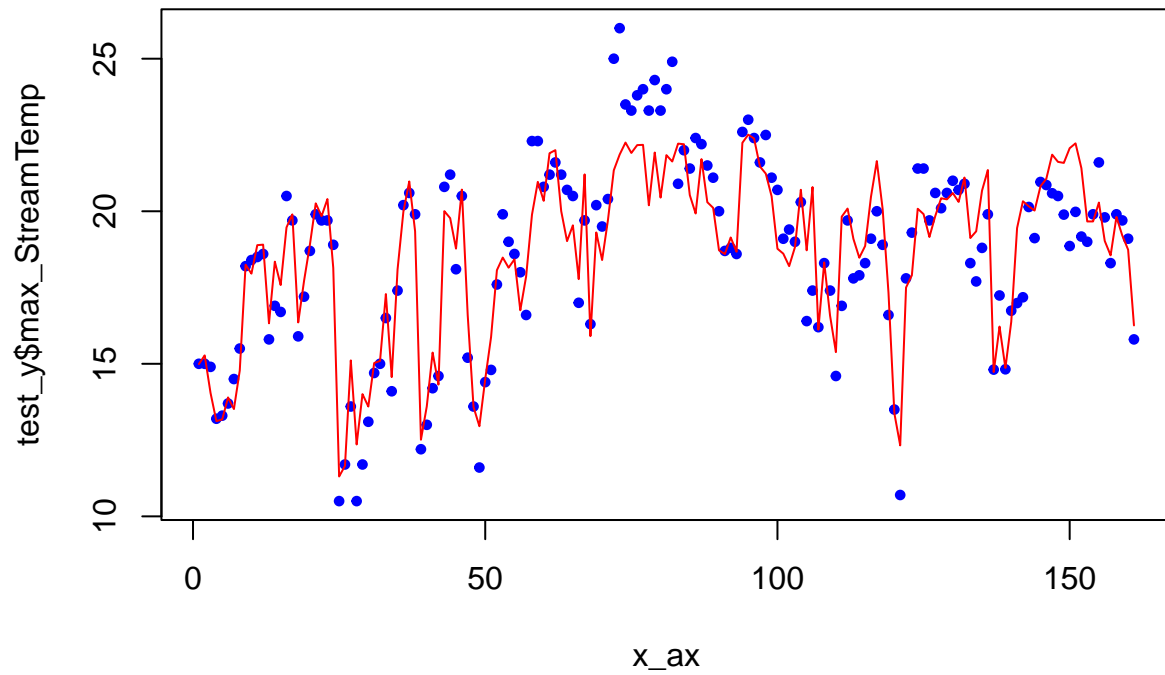
```
## The R-square of the test data is  0.839
```

```
# visualize the model, actual and predicted data
```

```
x_ax = 1:length(pred_y)
```

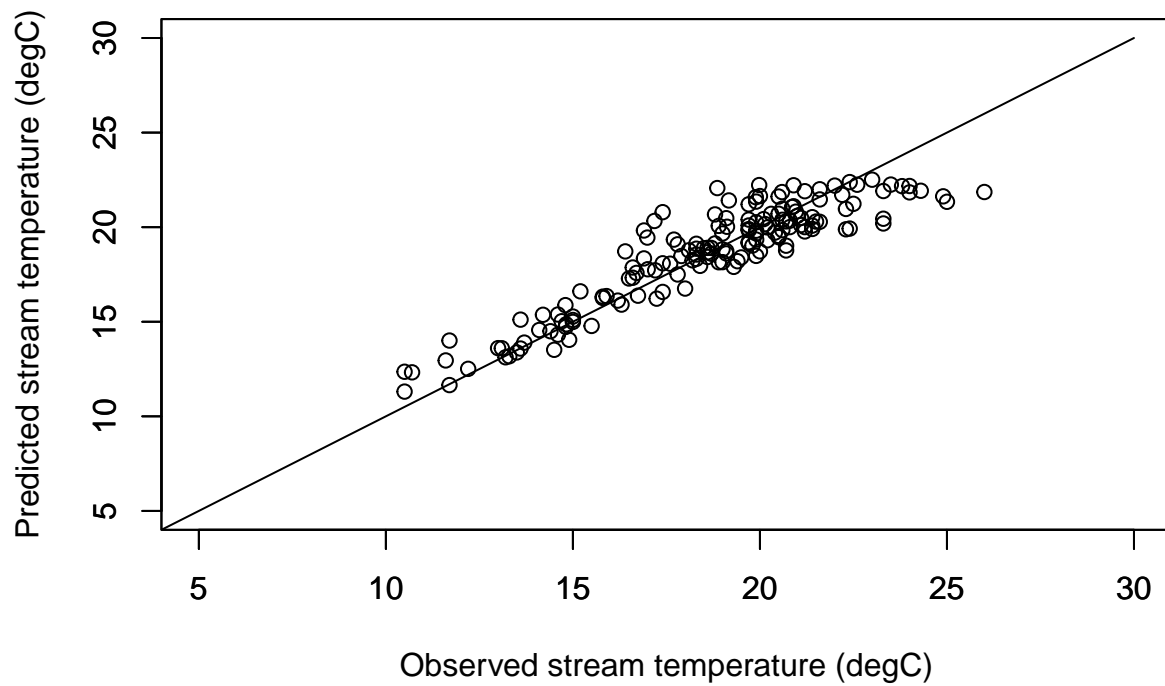
```
plot(x_ax, test_y$max_StreamTemp , col="blue", pch=20, cex=.9)
```

```
lines(x_ax, pred_y, col="red", pch=20, cex=.9)
```

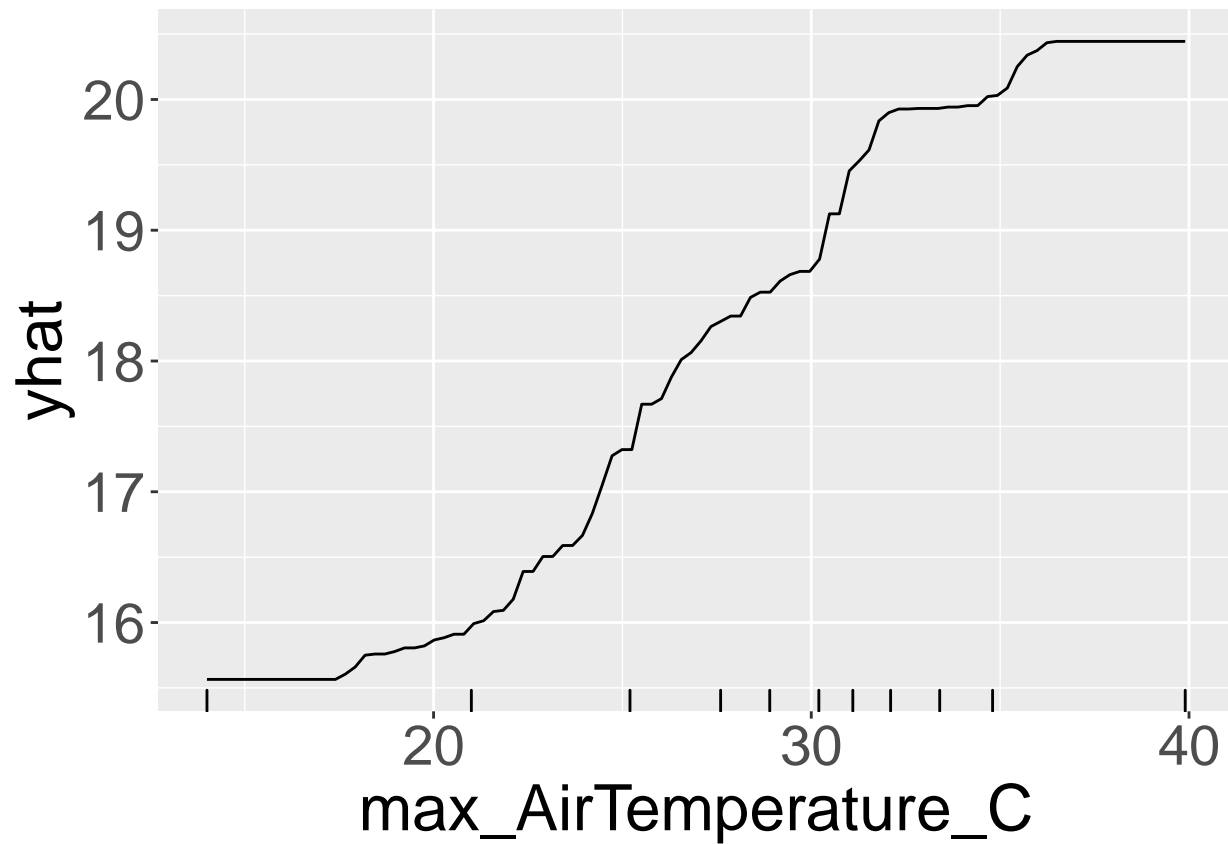


```
plot( test_y$max_StreamTemp, pred_y,xlim= xlim ,ylim= xlim, xlab="Observed stream temperature (degC)",
par(new=T)
x=seq(1,30)
plot(x,x,type="l",xlim= xlim ,ylim= xlim,xlab="",ylab=""))
```

## GBM



```
model_gbm %>%
  partial(pred.var = "max_AirTemperature_C" , n.trees = model_gbm$n.trees, grid.resolution = 100) %>%
  autoplot(rug = TRUE, train = train)+theme(axis.text=element_text(size=21),
    axis.title=element_text(size=24))
```



```
#, "min_Discharge"

model_gbm %>%
  partial(pred.var = "min_Discharge" , n.trees = model_gbm$n.trees, grid.resolution = 100) %>%
  autoplot(rug = TRUE, train = train)+theme(axis.text=element_text(size=21),
    axis.title=element_text(size=24))
```

