```
1. Multiple y format (LR)
 X_train <- read.csv("ModelData/X_train2.csv")[,-1]; Y_train <- read.csv("ModelData/Y_train2.csv")[,-1]</pre>
 X_val <- read.csv("ModelData/X_val2.csv")[,-1]; Y_val <- read.csv("ModelData/Y_val2.csv")[,-1]</pre>
 X_test <- read.csv("ModelData/X_test2.csv")[,-1]; Y_test <- read.csv("ModelData/Y_test2.csv")[,-1]</pre>
 colnames(Y_train) <- c("Y1","Y2","Y3","Y4","Y5","Y6","Y7")
 Data_train <- cbind(Y_train,X_train); #head(Data_train); dim(Data_train)</pre>
 colnames(Y_val) <- c("Y1","Y2","Y3","Y4","Y5","Y6","Y7")</pre>
 Data_val <- cbind(Y_val, X_val); #head(Data_val); dim(Data_val)</pre>
 colnames(Y_test) <- c("Y1","Y2","Y3","Y4","Y5","Y6","Y7")
 Data_test <- cbind(Y_test, X_test); #head(Data_test); dim(Data_test)</pre>
1.1. Linear regression without interactions
 m1 <- lm(cbind(Data_train$Y1,Data_train$Y2,Data_train$Y3,Data_train$Y4,Data_train$Y5,Data_train$Y6,Data_train$Y7)
 ~ . , Data_train[,-c(1:7)])
 #summary(m1)
 #cat("MSE:", mean(m1$residuals^2))
 #cat("RMSE:", sqrt(mean(m1$residuals^2)))
 y_pred <- predict(m1, Data_val)</pre>
 y_real <- Y_val</pre>
 y_pred <- rbind(y_pred[,1],y_pred[,2],y_pred[,3],y_pred[,4],y_pred[,5],y_pred[,6],y_pred[,7])</pre>
 y_real <- rbind(y_real[,1],y_real[,2],y_real[,3],y_real[,4],y_real[,5],y_real[,6],y_real[,7])</pre>
 cat("MSE:", mean((y_real-y_pred)^2))
 ## MSE: 0.1117158
 cat("RMSE:", sqrt(mean((y_real-y_pred)^2)))
 ## RMSE: 0.3342391
 y_pred <- predict(m1, Data_test)</pre>
 y_real <- Y_test</pre>
 y_pred <- rbind(y_pred[,1],y_pred[,2],y_pred[,3],y_pred[,4],y_pred[,5],y_pred[,6],y_pred[,7])</pre>
 y_real <- rbind(y_real[,1],y_real[,2],y_real[,3],y_real[,4],y_real[,5],y_real[,6],y_real[,7])</pre>
 cat("MSE:", mean((y_real-y_pred)^2))
 cat("RMSE:", sqrt(mean((y_real-y_pred)^2)))
 ## RMSE: 0.3081202
2. One y format, "all var & station" interactions (LR, LASSO, Ridge, Decision Tree)
 Train <- read.csv("ModelData/Train2 long.csv"); #head(Train); dim(Train)</pre>
 Val <- read.csv("ModelData/Val2_long.csv"); #head(Val); dim(Val)</pre>
 Test <- read.csv("ModelData/Test2_long.csv"); #head(Test); dim(Test)</pre>
```

# categorical var for(i in c(4:6)){ Train[,i] <- as.character(Train[,i])</pre> Val[,i] <- as.character(Val[,i])</pre> Test[,i] <- as.character(Test[,i])</pre> 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

summarise(median=median(original\_mrt\_flow),iqr=IQR(original\_mrt\_flow))

## Warning in predict.lm(lr\_model, Val): prediction from a rank-deficient fit may

## Warning in predict.lm(lr\_model, Test[, -c(30:31)]): prediction from a

## rank-deficient fit may be misleading

## Original RMSE: 494.7286

library(glmnet)

## RMSE: 0.3192557

## Loading required package: Matrix

cat("RMSE:", sqrt(mean((y\_real-y\_pred)^2)))

y\_pred <- predict(lasso\_model, x)</pre>

y\_real <- Test\$mrt\_flow</pre>

## Original MAE: 337.1932

2.3. Ridge with interactions

#y <- rbind(Train, Val)\$mrt\_flow</pre>

y <- Train\$mrt flow

## RMSE: 0.3165786

y\_real <- Test\$mrt\_flow</pre>

library(glmnet)

y\_real <- Val\$mrt\_flow</pre>

## MSE: 0.09878929

## RMSE: 0.3143076

## Original RMSE: 486.4082

## Original MAE: 305.1055

library(rpart)

#summary(dt\_model)

## [1] "bike\_flow"

## MSE: 0.02402419

## RMSE: 0.1549974

## Original RMSE: 301.4344

## Original MAE: 175.7418

cat("RMSE:", sqrt(mean((y\_real-y\_pred)^2)))

## [7] "aqi"

"hour"

"relative\_humidity"

## [4] "previous\_mrt\_flow" "day\_in\_a\_week"

set.seed(1)

2.5. Regression tree (decision tree)

y\_pred <- predict(elastic\_model, x)</pre>

cat("MSE:", mean((y\_real-y\_pred)^2))

cat("RMSE:", sqrt(mean((y\_real-y\_pred)^2)))

set.seed(1)

y\_pred <- predict(ridge\_model, x)</pre>

scaling\_back <- group\_by(Test, mrt\_station)%>%

## be misleading

Test <- merge(Test,scaling back)</pre> Train  $\leftarrow$  Train[,-3]; Val  $\leftarrow$  Val[,-3]; Test  $\leftarrow$  Test[,-3] #summary(Train); summary(Val); summary(Test) 2.1. Linear regression with interactions lr\_model <- lm(mrt\_flow~ . \* mrt\_station, Train)</pre> #lr\_model <- lm(mrt\_flow ~ . \* mrt\_station, rbind(Train, Val))</pre> #summary(lr model) #cat("MSE:", mean(m1\$residuals^2)) #cat("RMSE:", sqrt(mean(m1\$residuals^2))) y\_pred <- predict(lr\_model, Val)</pre>

y\_real <- Val\$mrt\_flow</pre> cat("MSE:", mean((y\_real-y\_pred)^2)) ## MSE: 0.104414 cat("RMSE:", sqrt(mean((y\_real-y\_pred)^2))) ## RMSE: 0.3231316 y\_pred <- predict(lr\_model, Test[,-c(30:31)])</pre>

y\_real <- Test\$mrt\_flow</pre> cat("MSE:", mean((y\_real-y\_pred)^2)) ## MSE: 0.06850735 cat("RMSE:", sqrt(mean((y\_real-y\_pred)^2))) ## RMSE: 0.2617391

cat("Original RMSE:", sqrt(mean(( (y\_real\*Test\$iqr+Test\$median) - (y\_pred\*Test\$iqr+Test\$median) )^2)))

cat("Original MAE:", mean(abs((y\_real\*Test\$iqr+Test\$median) - (y\_pred\*Test\$iqr+Test\$median))))

## Original MAE: 314.6577 2.2. LASSO with interactions f <- as.formula(mrt\_flow ~ .\* mrt\_station) # using .\*. for all interactions y <- Train\$mrt\_flow x <- model.matrix(f, Train)[,-1] # using model.matrix to take advantage of f #y <- rbind(Train, Val)\$mrt\_flow</pre> #x <- model.matrix(f, rbind(Train, Val))[,-1] # using model.matrix to take advantage of f</pre>

## Loaded glmnet 4.1-7 lasso\_kfold <- cv.glmnet(x, y, alpha=0, nfolds=10)</pre> lasso\_best\_lambda <- lasso\_kfold\$lambda.min</pre> lasso\_model <- glmnet(x, y, alpha=0, lambda=lasso\_best\_lambda)</pre> x <- model.matrix(mrt\_flow ~.\*mrt\_station, rbind(Train, Val))[,-1][-(1:nrow(Train)),]</pre> y\_pred <- predict(lasso\_model, x)</pre> y\_real <- Val\$mrt\_flow</pre> cat("MSE:", mean((y\_real-y\_pred)^2)) ## MSE: 0.1019242

cat("MSE:", mean((y\_real-y\_pred)^2)) ## MSE: 0.07014739 cat("RMSE:", sqrt(mean((y\_real-y\_pred)^2))) ## RMSE: 0.2648535 cat("Original RMSE:", sqrt(mean(( (y\_real\*Test\$iqr+Test\$median) - (y\_pred\*Test\$iqr+Test\$median) )^2))) ## Original RMSE: 539.6949

x <- model.matrix(mrt\_flow ~.\*mrt\_station, rbind(Train,Test[,-c(30:31)]))[,-1][-(1:nrow(Train)),]

cat("Original MAE:", mean(abs((y\_real\*Test\$iqr+Test\$median) - (y\_pred\*Test\$iqr+Test\$median))))

f <- as.formula(mrt\_flow ~ .\* mrt\_station) # using .\*. for all interactions

x <- model.matrix(f, Train)[,-1] # using model.matrix to take advantage of f

#x <- model.matrix(f, rbind(Train, Val))[,-1] # using model.matrix to take advantage of f</pre>

 $x \leftarrow model.matrix(mrt_flow \sim .*mrt_station, rbind(Train, Test[, -c(30:31)]))[, -1][-(1:nrow(Train)),]$ 

#x <- model.matrix(mrt\_flow ~.\*mrt\_station, rbind(Train, Val, Test))[,-1][-(1:nrow(rbind(Train, Val))),]</pre>

#x <- model.matrix(mrt flow ~.\*mrt station, rbind(Train, Val, Test))[,-1][-(1:nrow(rbind(Train, Val))),]</pre>

library(glmnet) ridge\_kfold <- cv.glmnet(x, y, alpha=1, nfolds=10)</pre> ridge\_best\_lambda <- ridge\_kfold\$lambda.min</pre> ridge\_model <- glmnet(x, y, alpha=1, lambda=ridge\_best\_lambda)</pre> x <- model.matrix(mrt\_flow ~.\*mrt\_station, rbind(Train, Val))[,-1][-(1:nrow(Train)),]</pre> y\_pred <- predict(ridge\_model, x)</pre> y\_real <- Val\$mrt\_flow</pre> cat("MSE:", mean((y real-y pred)^2)) ## MSE: 0.100222 cat("RMSE:", sqrt(mean((y\_real-y\_pred)^2)))

```
cat("MSE:", mean((y_real-y_pred)^2))
 ## MSE: 0.06638792
 cat("RMSE:", sqrt(mean((y real-y pred)^2)))
 ## RMSE: 0.2576585
 cat("Original RMSE:", sqrt(mean(( (y_real*Test$iqr+Test$median) - (y_pred*Test$iqr+Test$median) )^2)))
 ## Original RMSE: 481.8453
 cat("Original MAE:", mean(abs((y_real*Test$iqr+Test$median) - (y_pred*Test$iqr+Test$median))))
 ## Original MAE: 303.8942
2.4. Elastic (combine LASSO and Ridge) with interactions
 f <- as.formula(mrt_flow ~ . * mrt_station) # using .*. for all interactions
 x <- model.matrix(f, Train)[,-1] # using model.matrix to take advantage of f
 y <- Train$mrt_flow
```

elastic\_model <- glmnet(x, y, alpha=0.05, lambda=0.01) # choose alpha and beta by the performance in val

x <- model.matrix(mrt\_flow ~.\*mrt\_station, rbind(Train, Val))[,-1][-(1:nrow(Train)),]</pre>

 $x \leftarrow model.matrix(mrt_flow \sim .*mrt_station, rbind(Train,Test[,-c(30:31)]))[,-1][-(1:nrow(Train)),]$ y\_pred <- predict(elastic\_model, x)</pre> y\_real <- Test\$mrt\_flow</pre> cat("MSE:", mean((y\_real-y\_pred)^2)) ## MSE: 0.06627837 cat("RMSE:", sqrt(mean((y\_real-y\_pred)^2))) ## RMSE: 0.2574459

cat("Original RMSE:", sqrt(mean(( (y\_real\*Test\$iqr+Test\$median) - (y\_pred\*Test\$iqr+Test\$median) )^2)))

cat("Original MAE:", mean(abs((y\_real\*Test\$iqr+Test\$median) - (y\_pred\*Test\$iqr+Test\$median))))

dt\_model <- rpart(mrt\_flow~., Train, cp=0.000004) # choose cp by the performance in val

#printcp(dt\_model); plotcp(dt\_model) dt\_model\$variable.importance bike\_flow hour previous\_mrt\_flow no 21736.49162 1768.69502 ## 21515.46961 12988.46256 ## day\_in\_a\_week no2 mrt\_station nox 1648.53914 1245.69957 1163.52528 ## 1704.84437 ## o3\_8hr air\_pressure month air\_temperature 716.73295 522.52496 341.15955 335.46342 ## relative\_humidity о3 co\_8hr CO ## 283.58198 256.31721 174.65703 124.11705 ## sunshine\_duration windspeed winddirec aqi ## 94.74418 93.77009 93.05271 85.13969 ## pm2.5\_avg precipitation so2\_avg pm10\_avg 43.77014 ## 48.81502 44.62035 42.35497 ## pm10 pm2.5 status so2 24.42436 41.24490 23.61094 19.64605 setdiff((dt\_model\$frame\$var), "<leaf>") # variables used

"mrt\_station"

"air\_pressure"

"month"

## [10] "no2" "pm2.5" "pm10" ## [13] "no" "pm10\_avg" "nox" ## [16] "o3 8hr" "winddirec" "sunshine\_duration" ## [19] "precipitation" "co" "air\_temperature" ## [22] "o3" "so2" "windspeed" ## [25] "co\_8hr" "so2 avg" "pm2.5\_avg" y\_pred <- predict(dt\_model, Val)</pre> y\_real <- Val\$mrt\_flow</pre> cat("MSE:", mean((y\_real-y\_pred)^2)) ## MSE: 0.05324942 cat("RMSE:", sqrt(mean((y\_real-y\_pred)^2))) ## RMSE: 0.2307584 y\_pred <- predict(dt\_model, Test[,-c(30:31)])</pre> y\_real <- Test\$mrt\_flow</pre> cat("MSE:", mean((y\_real-y\_pred)^2))

cat("Original RMSE:", sqrt(mean(( (y\_real\*Test\$iqr+Test\$median) - (y\_pred\*Test\$iqr+Test\$median) )^2)))

cat("Original MAE:", mean(abs((y\_real\*Test\$iqr+Test\$median) - (y\_pred\*Test\$iqr+Test\$median))))