

EXPLORATORY DATA ANALYSIS

```
ft_fire <- read.csv("forestfires.csv",header=T)
data(ft_fire)

# Warning in data(ft_fire): data set 'ft_fire' not found

str(ft_fire)

## 'data.frame': 517 obs. of 13 variables:
## $ X : int 1 7 8 8 8 8 7 7
## $ Y : int 5 4 4 6 6 6 6 5 ...
## $ month: chr "mar" "oct" "oct" "mar" ...
## $ day : chr "1" "2" "3" "1" ...
## $ FFCMC: num 88.2 90.6 98.6 91.7 89.3 92.3 91.5 91.92 5 ...
## $ DMC : num 28.2 35.4 43.7 33.3 51.3 ...
## $ DC : num 84.3 86.9 1 686.7 77.5 282.2 ...
## $ ISI : num 5.1 6.7 6.7 9 9.6 14.7 8.5 18.7 7.1 ...
## $ temp : num 8.2 12.64 8.2 11.1 22.2 9.2 12.9 12.2 8 ...
## $ RH : int 51 33 33 97 99 29 27 86 63 48 ...
## $ wind : num 6.7 6 9 1.3 4 1.8 5.4 3.1 2.2 5.4 ...
## $ rain : num 0 0 0 0 2 0 0 0 0 0 ...
## $ area : num 0 0 0 0 0 0 0 0 0 ...

summary(ft_fire)

##      X      Y      month      day      Length:517
## Min.   1.0000   Min.   12.0   Length:517   Class: character
## 1st Qu.:3.000   1st Qu.:14.0   Class: character
## Median:14.000   Median:14.0   Mode: character
## Mean   14.669   Mean   14.3
## 3rd Qu.:17.000   3rd Qu.:15.0
## Max.    19.000   Max.    19.0

##      FFCMC      DMC      DC      ISI
## Min.   88.200   Min.   28.2   Min.   84.3   Min.   5.1
## 1st Qu.:98.20   1st Qu.:35.4   1st Qu.:86.6   1st Qu.:43.7
## Median:101.00   Median:38.8   Median:86.2   Median: 5.4
## Mean   100.64   Mean   33.0   Mean   84.7   Mean   6.622
## 3rd Qu.:102.90   3rd Qu.:43.7   3rd Qu.:91.9   3rd Qu.:10.800
## Max.    106.00   Max.   51.3   Max.   886.6   Max.   18.700

##      temp      RH      wind      rain
## Min.   8.200   Min.   33.00   Min.   1.300   Min.   0.000000
## 1st Qu.:15.50   1st Qu.:33.00   1st Qu.:2.700   1st Qu.:0.000000
## Median:19.30   Median:42.00   Median:4.000   Median:0.000000
## Mean   18.000   Mean   42.00   Mean   4.012   Mean:0.02166
## 3rd Qu.:22.80   3rd Qu.:53.00   3rd Qu.:4.900   3rd Qu.:0.000000
## Max.    22.80   Max.   100.00   Max.   19.400   Max.   16.400000

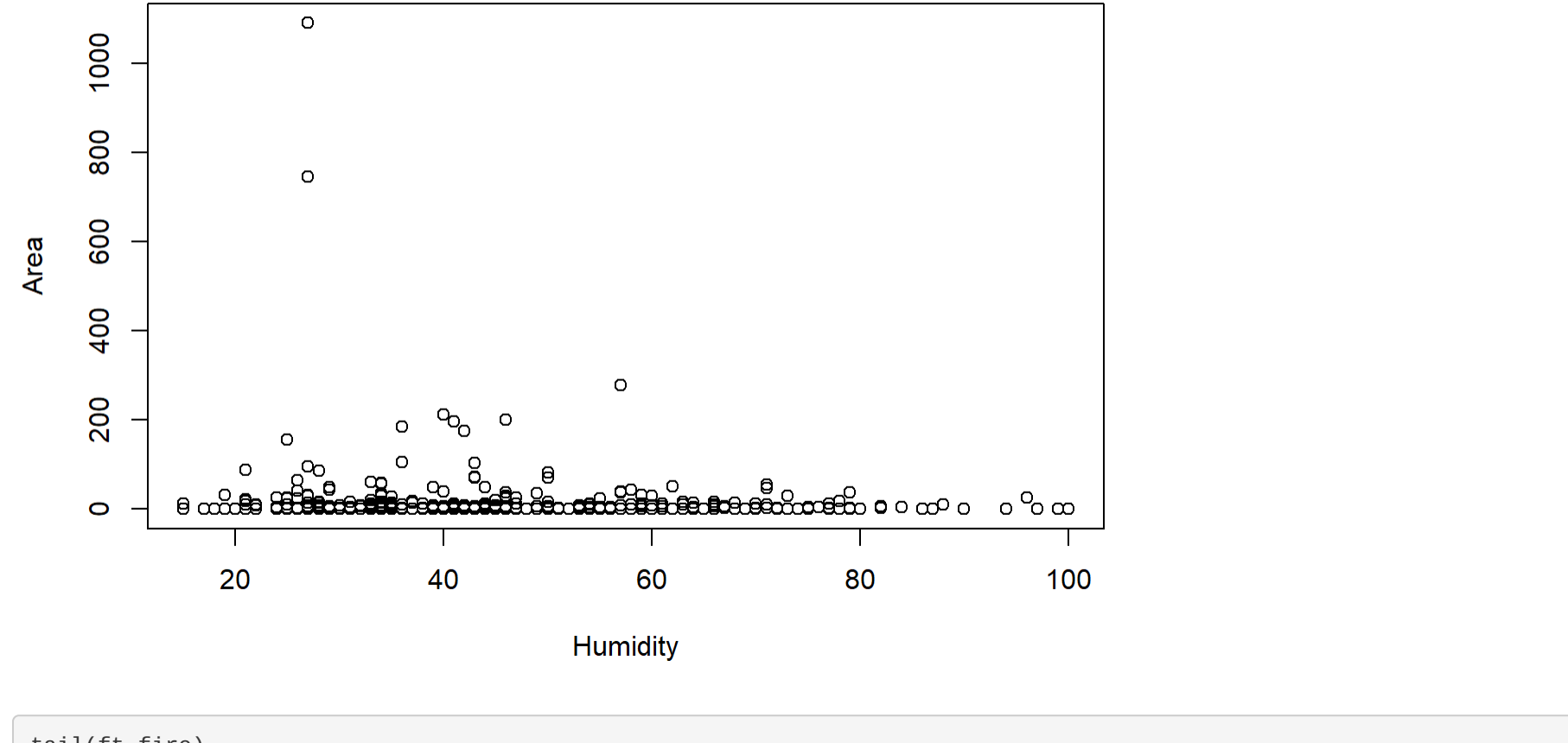
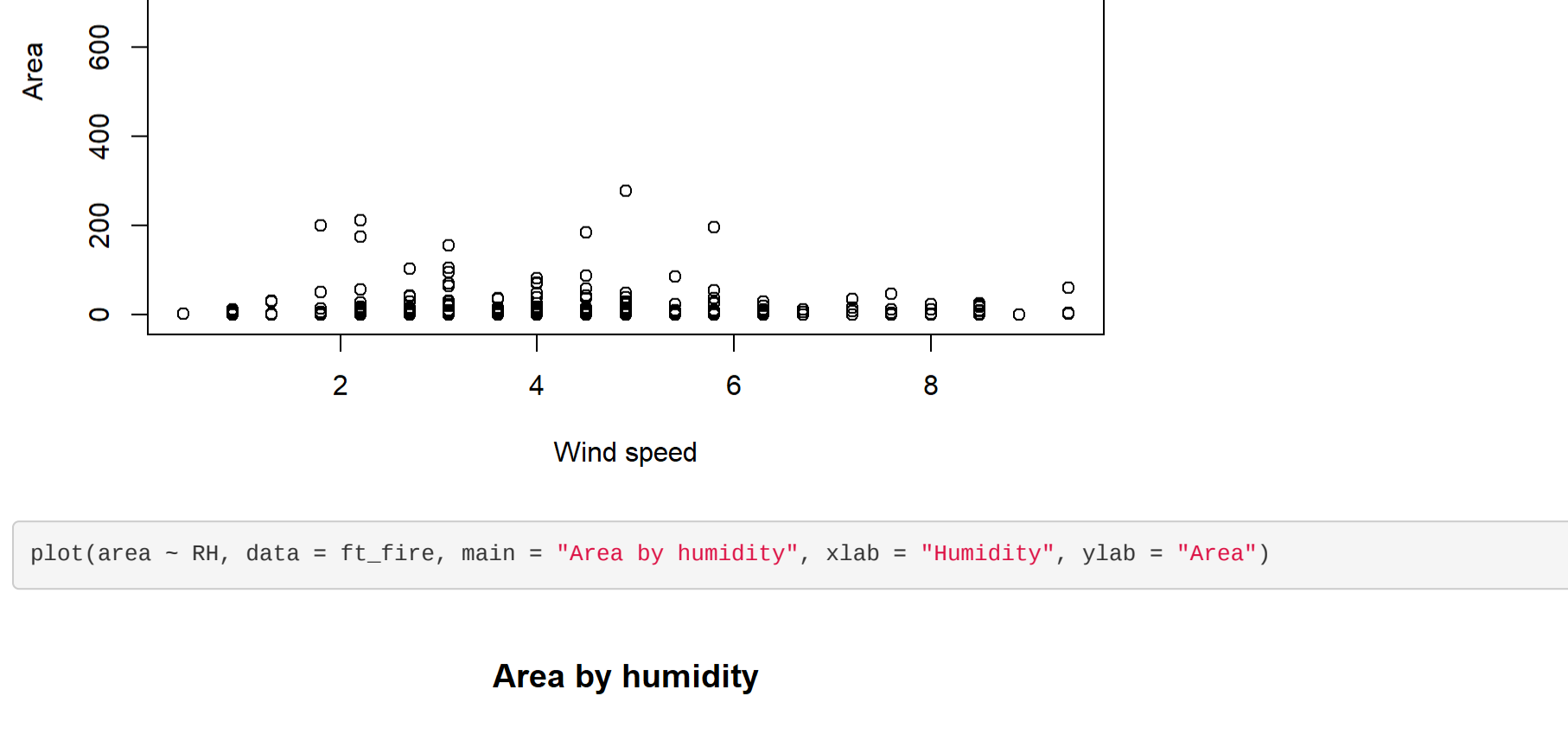
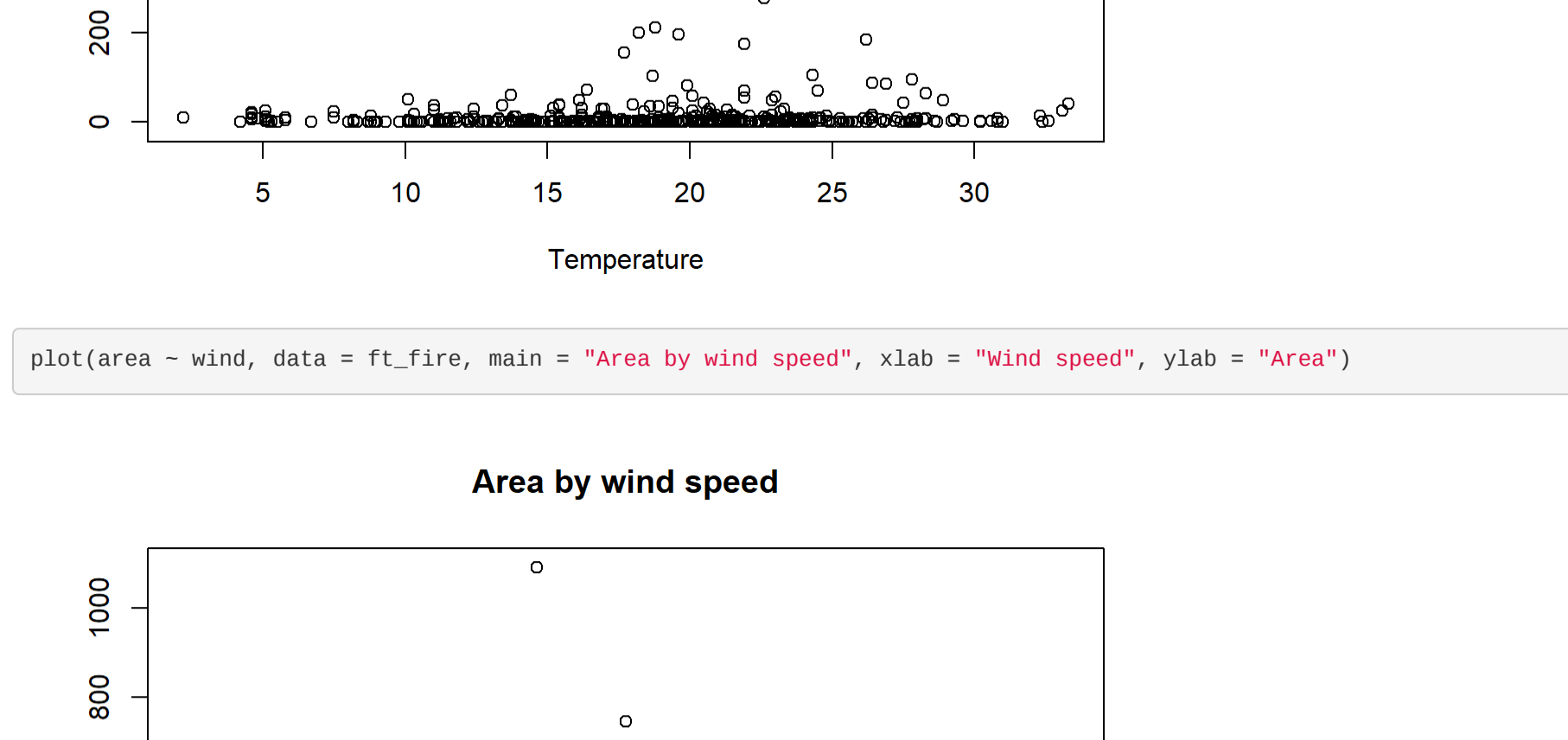
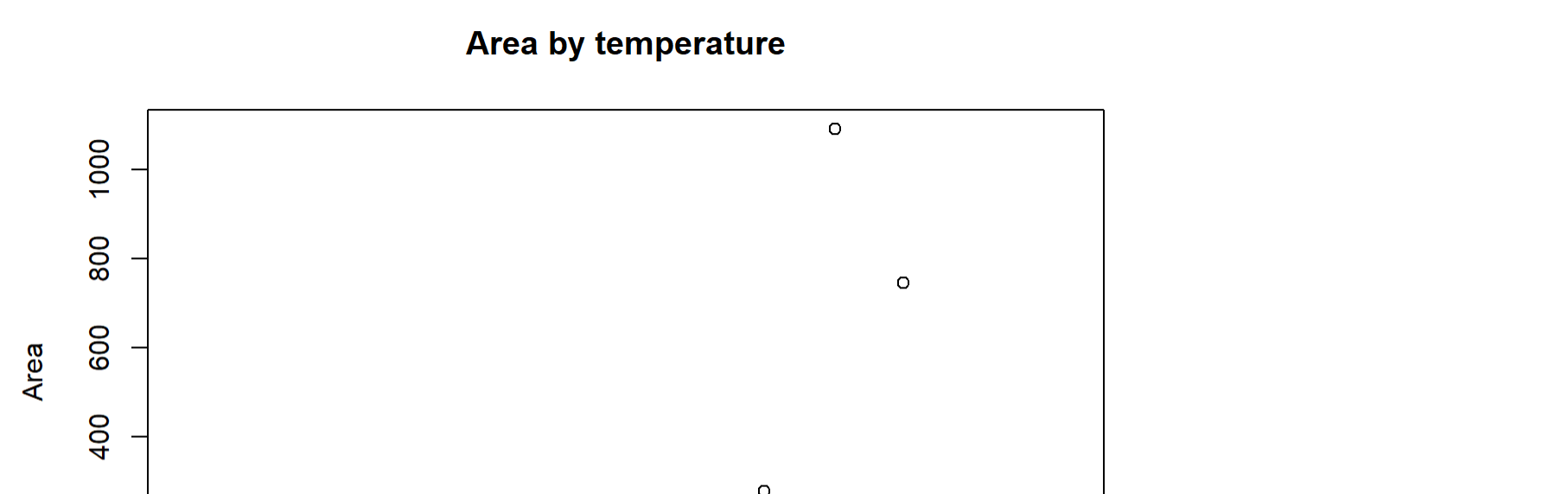
##      area
## Min.   : 0.00
## 1st Qu.: 0.00
## Median : 0.52
## Mean   : 12.85
## 3rd Qu.: 6.57
## Max.   :10990.84

head(ft_fire)

##      X Y month day FFCMC DMC DC ISI temp RH wind rain area
## 1 7 5 mar fri 86.2 26.2 94.3 5.1 8.2 51 6.7 0.0 0
## 2 7 4 oct tue 98.6 35.4 669.3 6.7 16.8 33 9.0 0.0 0
## 3 7 4 oct sat 99.6 43.7 686.9 9.7 16.8 33 1.3 0.0 0
## 4 8 6 mar fri 91.7 33.3 77.5 9.0 8.3 97 4.0 0.2 0
## 5 8 6 mar sun 89.3 33.3 102.9 9.6 11.4 99 1.0 0.0 0
## 6 8 6 aug sun 92.3 85.3 488.0 14.7 22.2 29 5.4 0.0 0

ft_fire$month <- as.numeric(factor(ft_fire$month))
ft_fire$day <- as.numeric(factor(ft_fire$day))
ggplot(data = ft_fire, aes(x = area)) +
  geom_histogram(fill = "blue", color = "black") +
  labs(title = "Histogram of Area", x = "Area", y = "Count") +
  theme_bw()

## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



```
tail(ft_fire)

##      X Y month day FFCMC DMC DC ISI temp RH wind rain area
## 512 8 6 2 4 81.6 56.7 655.6 1.9 27.8 35 2.7 0 0.00
## 513 4 3 2 4 81.6 56.7 655.6 1.9 27.8 32 2.7 0 6.44
## 514 2 4 2 4 81.6 56.7 655.6 1.9 21.0 71 0.0 0 54.29
## 515 7 4 2 4 81.6 56.7 655.6 1.9 21.2 70 6.7 0 11.15
## 516 1 4 2 3 94.4 146.0 614.7 11.3 25.6 42 4.0 0 0.00
## 517 6 3 10 0 79.9 3.0 396.7 1.1 11.8 35 4.5 0 0.00
```

DATA PREPROCESSING

```
set.seed(123)
Q1 <- quantile(ft_fire$area, 0.25)
Q3 <- quantile(ft_fire$area, 0.75)
log_area <- log(Q1 + (Q3 - Q1) * ft_fire$area)
scaled_predictors <- scale(ft_fire[, -14])

fire_nrm <- cbind(scaled_predictors, log_area = as.data.frame(ft_fire$log_area))

fire_nrm$log_area <- ft_fire$log_area

fire_nrm <- fire_nrm[, -which(names(fire_nrm) == "area")]
fire_nrm <- fire_nrm[, -which(names(fire_nrm) == "ft_fire$log_area")]
fire_nrm <- fire_nrm[, -which(names(fire_nrm) == "month")]
fire_nrm <- fire_nrm[, -which(names(fire_nrm) == "day")]

tail(fire_nrm)

##      X      Y      FFCMC      DMC      DC      ISI
## 511 0.5991227 0.5650757 0.07372876 0.8822143 0.8293628 -0.4194354
## 512 1.4687518 1.3743839 1.55159292 0.8398904 0.4806821 -1.5422418
## 513 -0.2685064 1.0535097 -1.55159292 -0.8398904 0.4806821 -1.5422418
## 514 1.8293713 -0.2442125 -1.55159292 -0.8398904 0.4806821 -1.5422418
## 515 -1.5695600 -0.2442125 0.66157852 0.9556083 0.2766938 0.4874468
## 516 0.5991227 -1.0535097 -1.91458859 -1.6790649 -1.7592911 -1.7149813
## 517 -0.1153148 1.0666083 0.8044892 -0.87812976 0.3576744
## 512 1.5644580 -0.5920088 -0.7148485 -0.87812976 0.0000000
## 513 1.5644580 -0.7703996 -0.7148485 -0.87812976 0.0000000
## 515 0.4896142 1.5581071 1.53562449 -0.87812976 2.4981519
## 516 1.1795101 -0.1619706 0.0179953 -0.87812976 0.0000000
## 517 -1.2351633 0.8377432 0.2983846 -0.87812976 0.0000000
```

LASSO REGRESSION

```
set.seed(123)
train_index <- createDataPartition(fire_nrm$log_area, p = 0.7, list = FALSE)
train <- fire_nrm[train_index, ]
test <- fire_nrm[-train_index, ]

xtrain <- train[, -13]
ytrain <- train$log_area
xtest <- test[, -13]
ytest <- test$log_area

lasso <- glmnet(xtrain, ytrain, alpha=1)
cv <- cv.glmnet(as.matrix(xtrain), ytrain, alpha=1)
bestlambda <- cv1lambda.min
print(paste0("Optimal lambda: ", bestlambda))

## [1] "Optimal lambda: 0.0567562491975646"

lasso_best <- glmnet(xtrain, ytrain, alpha=1, lambda=bestlambda)

plot(lasso, xvar = "lambda", label=T)
```



```
# Evaluate the model
ypred <- predict(lasso_best, as.matrix(xtest))

ypred <- exp(ypred) - 1
mse <- mean((exp(ytest) - 1) - ypred)^2
RMSE.lasso <- sqrt(mse)
MAD.lasso <- mean(abs((exp(ytest)-1) - ypred))

print(paste0("Optimal lambda: ", bestlambda))

## [1] "Optimal lambda: 0.0567562491975646"

print(paste0("RMSE: ", RMSE.lasso))

## [1] "RMSE: 3.62129748725139"

print(paste0("MSE: ", mse))

## [1] "MSE: 13.1137954911732"

print(paste0("MAD: ", MAD.lasso))

## [1] "MAD: 2.28599996826835"
```

RIDGE REGRESSION

```
set.seed(123)
train_index <- createDataPartition(fire_nrm$log_area, p = 0.7, list = FALSE)
train_data <- fire_nrm[train_index, ]
test_data <- fire_nrm[-train_index, ]

x_train <- train_data[, -13]
y_train <- train_data$log_area
x_test <- test_data[, -13]
y_test <- test_data$log_area

ridge <- glmnet(x_train, y_train, alpha=0)
cvfit <- cv.glmnet(as.matrix(x_train), y_train, alpha=0)
lambda_0 <- cvfit$lambda.min

print(paste0("Optimal lambda: ", lambda_0))

## [1] "Optimal lambda: 82.3421489652069"

ridge_best <- glmnet(x_train, y_train, alpha=0, lambda=lambda_0)

plot(ridge, xvar = "lambda", label=T)
```



```
# Evaluate the model
ypred <- predict(ridge_best, as.matrix(x_test), s = lambda_0)

ypred <- exp(ypred) - 1
mse_r <- mean(((exp(y_test)-1) - ypred)^2)
RMSE.ridge <- sqrt(mse_r)
MAD.ridge <- mean(abs((exp(y_test)-1) - ypred))

print(paste0("Optimal lambda: ", lambda_0))

## [1] "Optimal lambda: 82.3421489652069"

print(paste0("RMSE: ", RMSE.ridge))

## [1] "RMSE: 12.8652853331844"

print(paste0("RMSE: ", RMSE.ridge))

## [1] "RMSE: 3.56882106231144"

print(paste0("MAD: ", MAD.ridge))

## [1] "MAD: 2.26842480020263"
```

ELASTIC NET

```
set.seed(123)
train_index <- createDataPartition(fire_nrm$log_area, p = 0.7, list = FALSE)
train_elas <- fire_nrm[train_index, ]
test_elas <- fire_nrm[-train_index, ]

x_train_elas <- train_elas[, -13]
y_train_elas <- train_elas$log_area
x_test_elas <- test_elas[, -13]
y_test_elas <- test_elas$log_area

elastic <- glmnet(x_train_elas, y_train_elas, alpha=0.5)
cv_elastic <- cv.glmnet(as.matrix(x_train_elas), y_train_elas, alpha=0.5)
lambda_elastic <- cv_elastic$lambda.min

elas_best <- glmnet(x_train_elas, y_train_elas, alpha=0.5, lambda=lambda_elastic)

plot(elastic, xvar = "lambda", label=T)
```



```
# Evaluate the model
ypred_elas <- predict(elas_best, as.matrix(x_test_elas), s = lambda_elastic)

ypred_elas <- exp(ypred_elas) - 1
elas_mse <- mean(((exp(y_test_elas)-1) - ypred_elas)^2)
RMSE.elas <- sqrt(elas_mse)
MAD.elas <- mean(abs((exp(y_test_elas)-1) - ypred_elas))

print(paste0("Optimal lambda: ", lambda_elastic))

## [1] "Optimal lambda: 0.113516498395109"

print(paste0("RMSE: ", elas_mse))

## [1] "RMSE: 13.5088284834"

print(paste0("RMSE: ", RMSE.elas))

## [1] "RMSE: 3.62658137652714"

print(paste0("MAD: ", MAD.elas))

## [1] "MAD: 2.28542549806184"
```

RANDOM FOREST

```
set.seed(123)
train_rand <- createDataPartition(fire_nrm$log_area, p = 0.7, list = FALSE)
train_rand <- fire_nrm[train_rand, ]
test_rand <- fire_nrm[-train_rand, ]

model_rand <- randomForest(log_area ~ ., data = train_rand)

pred_rand <- predict(model_rand, newdata = test_rand)

pred_rand <- exp(pred_rand) - 1
mse_rand <- mean(((exp(test_rand$log_area)-1) - pred_rand)^2)
RMSE.rand <- sqrt(mse_rand)
MAD.rand <- mean(abs((exp(test_rand$log_area)-1) - pred_rand))

print(paste0("RMSE: ", mse_rand))

## [1] "RMSE: 12.1263775641386"

print(paste0("RMSE: ", RMSE.rand))

## [1] "RMSE: 3.48229486092461"

print(paste0("MAD: ", MAD.rand))

## [1] "MAD: 2.16273568413039"
```

SUPPORT VECTOR REGRESSION

```
set.seed(123)
train_ind_svr <- createDataPartition(fire_nrm$log_area, p = 0.7, list = FALSE)
train_svr <- fire_nrm[train_ind_svr, ]
test_svr <- fire_nrm[-train_ind_svr, ]

tail(train_svr)

##      X      Y      FFCMC      DMC      DC      ISI
## 588 -1.1361354 -0.2442125 0.07372876 0.8822143 0.8293628 -0.4194354
## 589 0.5991227 0.5650757 0.07372876 0.8822143 0.8293628 -0.4194354
## 590 1.4687518 1.3743839 1.55159292 0.8398904 0.4806821 -1.5422418
## 591 -0.2685064 1.0535097 -1.55159292 -0.8398904 0.4806821 -1.5422418
## 592 1.8293713 -0.2442125 1.55159292 -0.8398904 0.4806821 -1.5422418
## 593 -1.5695600 -0.2442125 0.66157852 0.9556083 0.2766938 0.4874468
## 588 1.2326500 -0.2342472 -0.29787172 -0.87812976 0.0000000
## 591 -0.1153148 1.0666083 0.8044892 -0.87812976 0.3576744
## 592 1.5644580 -0.5920088 -0.7148485 -0.87812976 0.0000000
## 593 1.5644580 -0.7703996 -0.7148485 -0.87812976 0.0000000
## 595 0.4896142 1.5581071 1.53562449 -0.87812976 2.4981519
## 596 1.1795101 -0.1619706 0.0179953 -0.87812976 0.0000000

tune_model_svr <- tune(svm$log_area ~ ., data = train_svr, kernel = "radial", ranges = list(cost = c(0.01, 0.1, 1, 10, 100),
  epsilon = c(0.1, 0.2, 0.5, 0.8, 1)))
print(tune_model_svr$bestModel)

##
## Call:
## svm.best.tune$M2MO <- svm, train.x = log_area ~ ., data = train_svr,
## ranges = list(cost = c(0.01, 0.1, 1, 10, 100), epsilon = c(0.1,
## 0.2, 0.5, 0.8, 1)), kernel = "radial")
##
## Parameters:
## SVM-Type: eps-regression
## SVM-kernel: radial
## cost: 0.1
## gamma: 0.1
## epsilon: 0.8
##
## Number of Support Vectors: 184

model_svr <- svm(log_area ~ ., data = train_svr, kernel = "radial", cost = 0.1, epsilon = 0.8)

pred_svr <- predict(model_svr, newdata = test_svr)

pred_svr <- exp(pred_svr) - 1
mse_svr <- mean(((exp(test_svr$log_area)-1) - pred_svr)^2)
RMSE.svr <- sqrt(mse_svr)
MAD.svr <- mean(abs((exp(test_svr$log_area)-1) - pred_svr))

print(paste0("RMSE: ", mse_svr))

## [1] "RMSE: 12.6357238038991"

print(paste0("RMSE: ", RMSE.svr))

## [1] "RMSE: 3.55467627026556"

print(paste0("MAD: ", MAD.svr))

## [1] "MAD: 2.28995872496435"
```

SUMMARIZATION



Note that the echo = FALSE parameter was added to the code chunk to prevent printing of the R code that generated the plot.