

# Training Parent Nodes

Brian K. Masinde

```
# Working Environment: Clear working environment
rm(list = ls())

# Libraries: Load libraries
library(rpart)
library(here)
library(rpart.plot)
library(caret)
```

## Reusable functions

```
rmse <- function(actual, predicted) {
  sqrt(mean((actual - predicted)^2))
}
```

## Inputs

```
# inputs
base_train <- read.csv(here("data", "base_train.csv"))
base_validation <- read.csv(here("data", "base_validation.csv"))
base_test <- read.csv(here("data", "base_validation.csv"))

# Combining train and validation datasets to one
# Because we are going to use CV to train the models later
# naming it df_base_train2 to remain consistent with df naming
df_base_train2 <- rbind(base_train, base_validation)

cat("number of rows in combined train data:", nrow(df_base_train2), sep = " ")

## number of rows in combined train data: 7184
```

## Wind Model Training & Testing

### Decision trees

```
# Basic training of wind model
# We do not account for regions like the adjusted model.
```

```
base_wind_model <- rpart(wind_max ~ track_min_dist,
  data = df_base_train2,
  method = "anova")
```

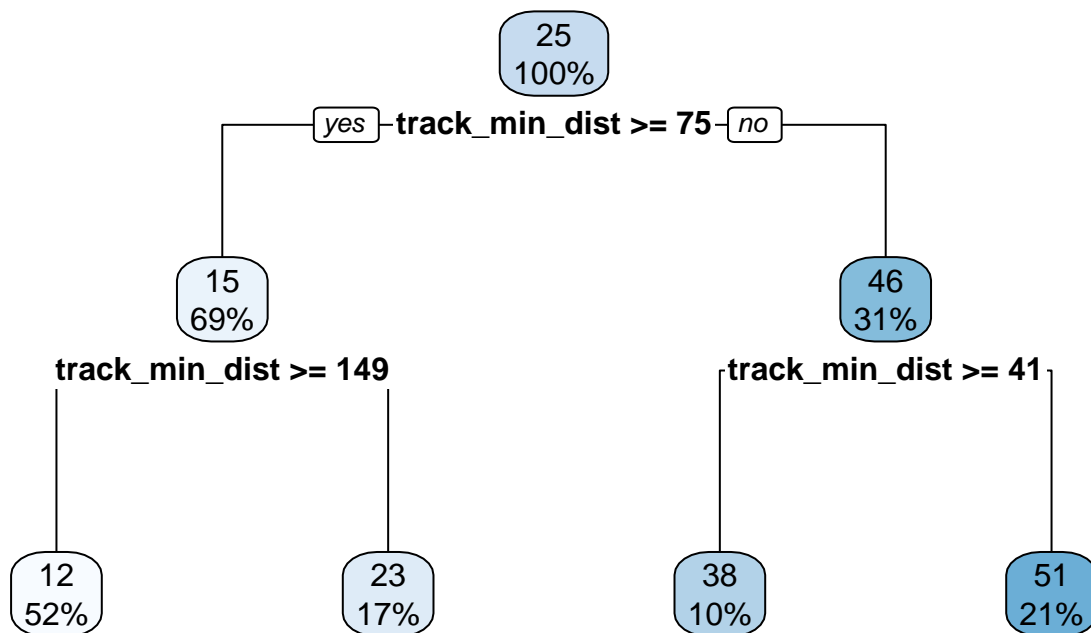
```
dec_wind_pred <- predict(base_wind_model, base_test)
```

```
rmse_dec_wind_pred <- rmse(actual = df_base_train2$wind_max,
  predicted = dec_wind_pred)
```

```
cat("rmse of decision tree of wind model", rmse_dec_wind_pred, sep = " ")
```

```
## rmse of decision tree of wind model 20.9515
```

```
rpart.plot(base_wind_model)
```



Optimizing For Best parameters

```

set.seed(123)
# Define training control
control <- trainControl(method = "cv", number = 5)

# Set tuning grid
grid <- expand.grid(
  cp = seq(0.001, 0.05, by = 0.005) # try several cp values
)

# Train model
dec_base_wind_model_tuned <- train(
  wind_max ~ track_min_dist, data = df_base_train2,
  method = "rpart",
  trControl = control,
  tuneGrid = grid
)

print(dec_base_wind_model_tuned)

```

```

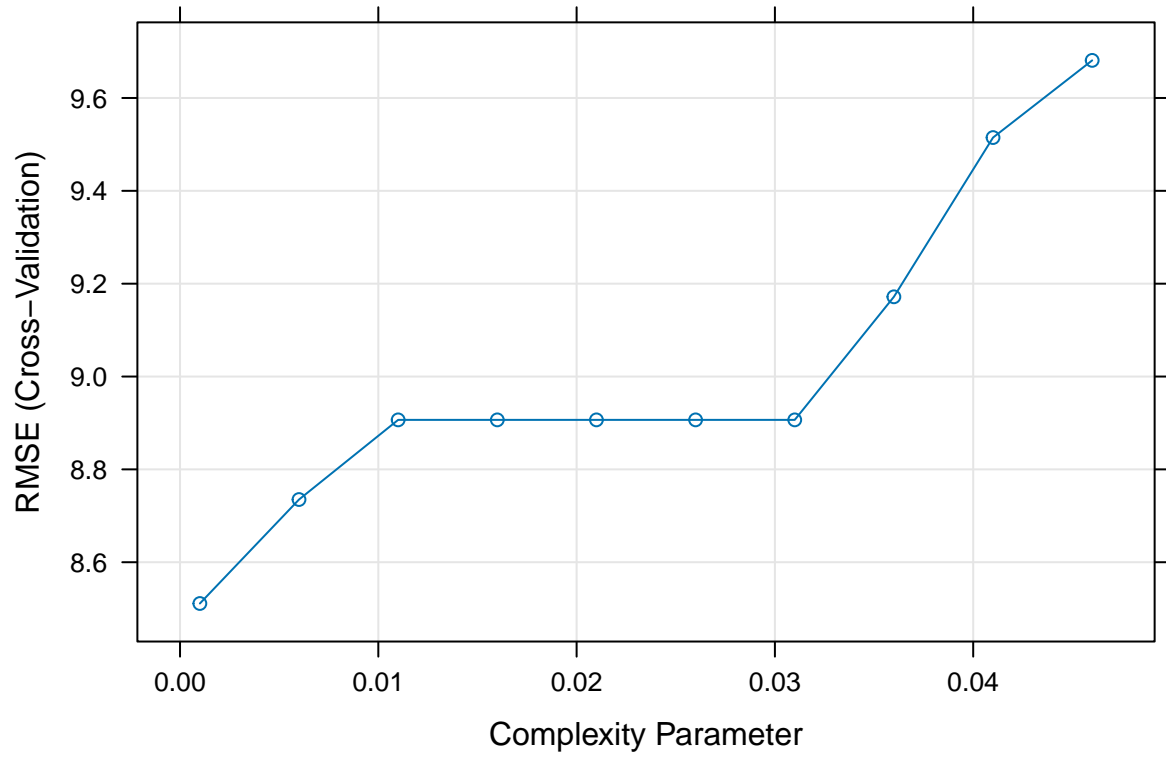
## CART
##
## 7184 samples
##    1 predictor
##
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 5748, 5746, 5747, 5747, 5748
## Resampling results across tuning parameters:
##
##   cp      RMSE      Rsquared    MAE
## 0.001  8.511103  0.7757006  5.823712
## 0.006  8.735074  0.7637364  6.015736
## 0.011  8.906656  0.7543552  6.312006
## 0.016  8.906656  0.7543552  6.312006
## 0.021  8.906656  0.7543552  6.312006
## 0.026  8.906656  0.7543552  6.312006
## 0.031  8.906656  0.7543552  6.312006
## 0.036  9.171699  0.7393462  6.460167
## 0.041  9.514772  0.7200191  6.761375
## 0.046  9.681026  0.7097058  6.872794
##
## RMSE was used to select the optimal model using the smallest value.
## The final value used for the model was cp = 0.001.

```

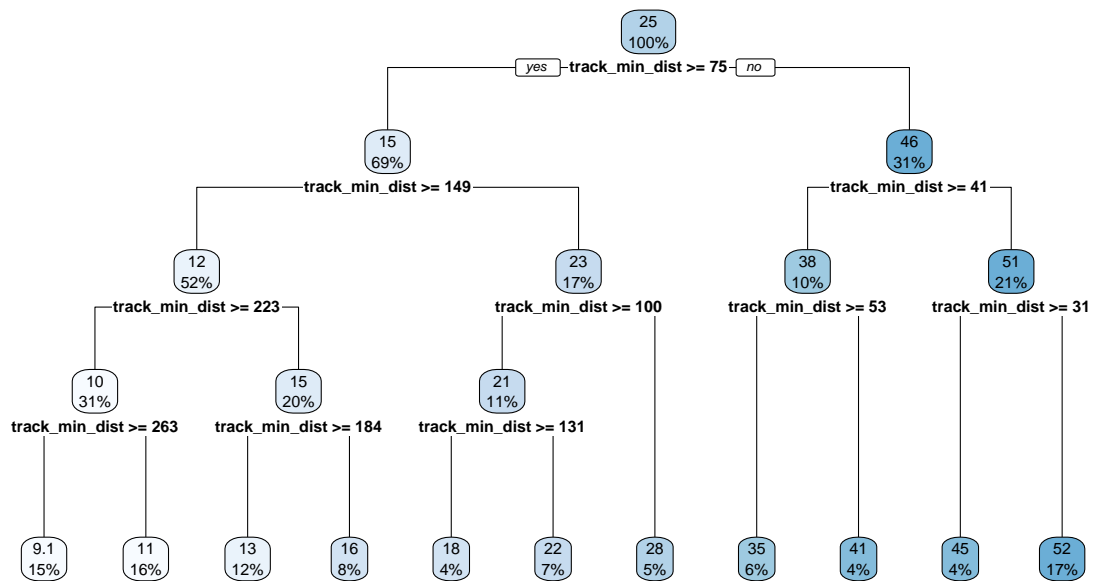
```

plot(dec_base_wind_model_tuned)

```



```
rpart.plot(dec_base_wind_model_tuned$finalModel)
```



```
dec_wind_pred_tuned <- predict(dec_base_wind_model_tuned, base_test)

rmse_dec_wind_pred_tuned <- rmse(actual = df_base_train2$wind_max,
                                  predicted = dec_wind_pred_tuned)

cat("rmse of tuned decision tree of wind model", rmse_dec_wind_pred_tuned, sep = " ")
```

```
## rmse of tuned decision tree of wind model 21.06189
```

## Output

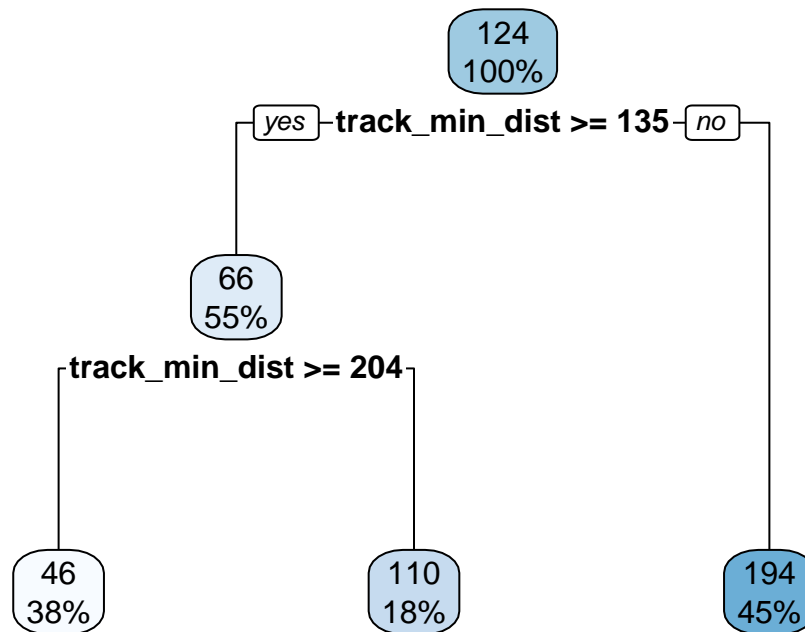
We output the caret tuned decision tree model

```
full_path <- here("unadjusted SCM/new base models")
saveRDS(dec_base_wind_model_tuned,
        file = file.path(full_path, paste0("dec_base_wind_model_tuned", ".rds")))
```

## Rain model Training and Testing

```
base_rain_model <- rpart(rain_total ~ track_min_dist,
                        data = df_base_train2,
                        method = "anova")
```

```
rpart.plot(base_rain_model)
```



## Optimizing For Best parameters

```
set.seed(123)
# Define training control
control <- trainControl(method = "cv", number = 5)

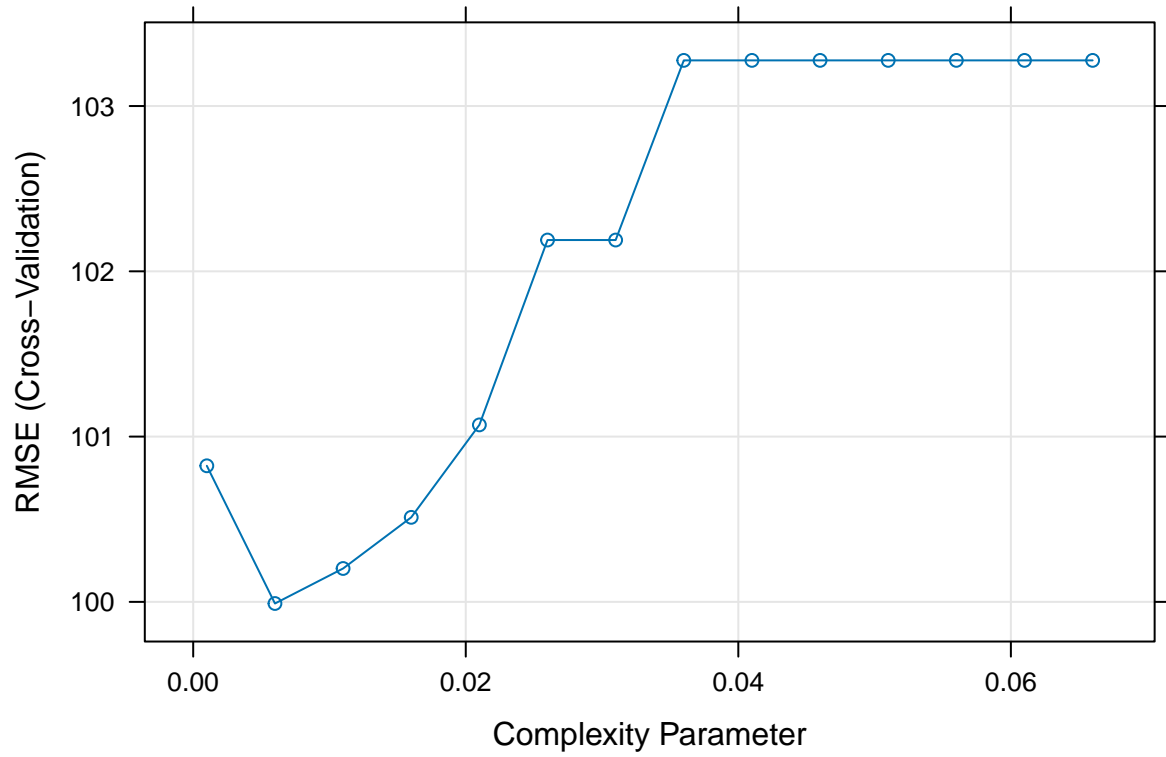
# Set tuning grid
grid <- expand.grid(
  cp = seq(0.001, 0.07, by = 0.005) # try several cp values
)

# Train model
dec_base_rain_model_tuned <- train(
  rain_total ~ track_min_dist, data = df_base_train2,
  method = "rpart",
  trControl = control,
  tuneGrid = grid
)

print(dec_base_rain_model_tuned)
```

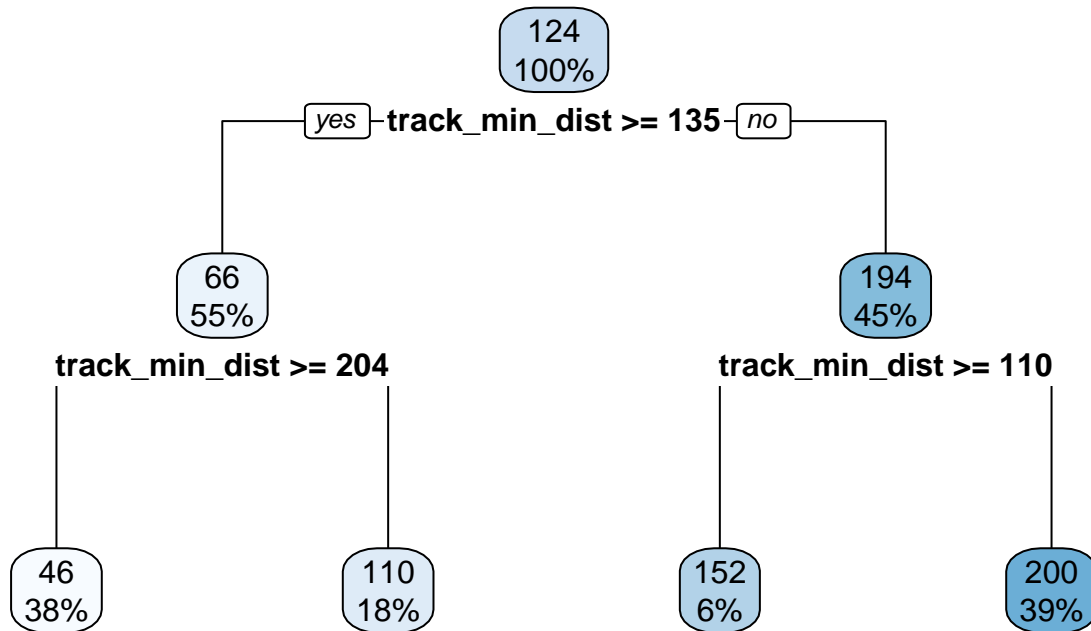
```
## CART
##
## 7184 samples
##    1 predictor
##
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 5748, 5747, 5746, 5747, 5748
## Resampling results across tuning parameters:
##
##    cp      RMSE      Rsquared    MAE
##    0.001  100.82332  0.3041851  68.26986
##    0.006   99.99011  0.3142410  67.81088
##    0.011  100.20194  0.3113975  67.87942
##    0.016  100.51143  0.3070774  68.11902
##    0.021  101.07087  0.2992208  68.40240
##    0.026  102.18927  0.2832561  69.84635
##    0.031  102.18927  0.2832561  69.84635
##    0.036  103.27638  0.2684042  71.32285
##    0.041  103.27638  0.2684042  71.32285
##    0.046  103.27638  0.2684042  71.32285
##    0.051  103.27638  0.2684042  71.32285
##    0.056  103.27638  0.2684042  71.32285
##    0.061  103.27638  0.2684042  71.32285
##    0.066  103.27638  0.2684042  71.32285
##
## RMSE was used to select the optimal model using the smallest value.
## The final value used for the model was cp = 0.006.
```

```
plot(dec_base_rain_model_tuned)
```



```
rpart.plot(dec_base_rain_model_tuned$finalModel)
```





```
dec_rain_pred_tuned <- predict(dec_base_rain_model_tuned, base_test)

rmse_dec_rain_pred_tuned <- rmse(actual = df_base_train2$rain_total,
                                predicted = dec_rain_pred_tuned)

cat("rmse of tuned decision tree of rain model", rmse_dec_rain_pred_tuned, sep = " ")
```

```
## rmse of tuned decision tree of rain model 129.5931
```

## Saving Tuned Rain\_\_Total Decision Tree Model

```
full_path <- here("unadjusted SCM/new base models")
saveRDS(dec_base_rain_model_tuned,
        file = file.path(full_path, paste0("dec_base_rain_model_tuned", ".rds")))
```