

Tree Model Performance

Gwen Rino

4/20/2018

Model Description

This decision tree model uses all the variables except `week_start_date`, which is the unique identifier. It includes the engineered feature `year.season`. I ran the model with missing values imputed two ways, first by median values and second by k nearest neighbor values. I used the `caret` package to cross validate these models using the bootstrapping method.

```
# With knn imputation
tree_2 <- train(total_cases ~ . -week_start_date,
               data = dengue.knn,
               method = "rpart")

# With median value imputation
tree_3 <- train(total_cases ~ . -week_start_date,
               data = dengue.med,
               method = "rpart")
```

Model Evaluation

Both the cross validations returned best models with the tuning parameter $cp = 0.016$. However, in both these cases the MAE was a bit greater than 20 – better than the Naive and Linear models, but not by much. There was little difference between the MAEs of the `rpart` models with different imputation methods, which makes sense because a tree model shouldn't be very sensitive to this kind of preprocessing.

Most of the nodes in the visualized decision trees concern elements of time rather than weather. Again, it is clear that time is critical to understanding this data set. I need to learn about time series!

```
print(tree_2)

## CART
##
## 936 samples
## 26 predictor
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
## Summary of sample sizes: 936, 936, 936, 936, 936, 936, ...
## Resampling results across tuning parameters:
##
##   cp          RMSE      Rsquared    MAE
##   0.01609336  40.49876  0.4159504  21.52882
##   0.02550193  41.06601  0.3947639  22.24855
##   0.19421475  45.17067  0.4673377  24.88030
##
## RMSE was used to select the optimal model using the smallest value.
```

```

## The final value used for the model was cp = 0.01609336.
print(tree_3)

## CART
##
## 936 samples
## 26 predictor
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
## Summary of sample sizes: 936, 936, 936, 936, 936, 936, ...
## Resampling results across tuning parameters:
##
##   cp          RMSE      Rsquared   MAE
## 0.01609336  38.81746  0.4687525  21.42182
## 0.02550193  39.76874  0.4432102  21.91664
## 0.19421475  47.04866  0.4848420  25.77910
##
## RMSE was used to select the optimal model using the smallest value.
## The final value used for the model was cp = 0.01609336.

```