# Respiratory Illness Prediction

Zhirou Zhou

Dec 11, 2019

# Overview

# Variable Transformation

- Numeric: normalization
- Sequential questions

| Any furry animals | How many pets | Pets |
|---|---|---|
| No | NA | 0 |
| Yes | 3 | 3 |

- Virus and bacteria testing: 'tlda'

| Alias | visit_id | Adeno | B.pert |
|---|---|---|---|
| C5E05 | 2 | FALSE | FALSE |
| C5E05 | 4 | TRUE | FALSE |

- Vaccination record: 'vaccines'

| Alias | pCGA | DTaP (2 months) | DTaP (4 months) |
|---|---|---|---|
| C01D8 | 44 | 0 | 0 |
| C01D8 | 52 | 1 | 0 |
| C01D8 | 61 | 1 | 1 |

# Missing Values

- One row per subject

  - Mean

  - Most frequent category

- Multiple records per subject

  - Last observation carry forward/backward

  - Vaccines: most frequent category grouping by pCGA (whether most of other babies with same pCGA received that dose of vaccine)

# Feature Selection

- Correlation

|        | Weight | HC    | Length | Age   | Temp  | BMI    |
|--------|--------|-------|--------|-------|-------|--------|
| Weight | 1      | 0.938 | 0.947  | 0.910 | 0.040 | -0.025 |
| HC     |        | 1     | 0.944  | 0.920 | 0.057 | -0.036 |
| Length |        |       | 1      | 0.927 | 0.048 | -0.068 |
| Age    |        |       |        | 1     | 0.041 | -0.035 |
| Temp   |        |       |        |       | 1     | -0.013 |
| BMI    |        |       |        |       |       | 1      |

# Feature Selection

- Collapse infrequent classes

| #Pets | #Obs | #Pets | #Obs |
|------:|-----:|------:|-----:|
| 0 | 2020 | 0 | 2020 |
| 1 | 155 | 1 | 155 |
| 2 | 97 | 2 | 97 |
| 3 | 43 | 3 | 43 |
| 4 | 11 | 4 | 11 |
| 5 | 2 | 5 | 13 |
| 6 | 2 | | |
| 8 | 2 | | |
| 10 | 6 | | |
| 28 | 1 | | |

# Validation

- 10-fold cross validation

- Blocked by Alias

- Over-sampling to balance the target variable

# Tuning Parameter

- Grow a large tree $T_0$ by stopping the splitting process when some minimum node size is reached

- Prune $T_0$ using *cost-complexity pruning*

    - Define a subtree $T \subset T_0$
    - $t$ is the number of terminal nodes in $T$

    $$N_m = \#\{x_i \in R_m\} \ , \ \hat{c}_m \ , \ Q_m(T) \text{ (Node Impurity)},$$

    - Define the cost complexity criterion

    $$C_\alpha(T) = \sum_{m=1}^{t} N_m Q_m(T) + \alpha \, t$$

    - For each $\alpha$, there is a unique smallest subtree $T_\alpha$ that minimizes $C_\alpha(T)$

# Results: is_illness



minimum node size = 50, $\alpha = 0.03$

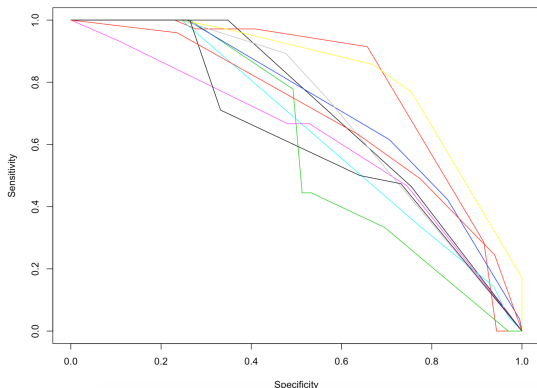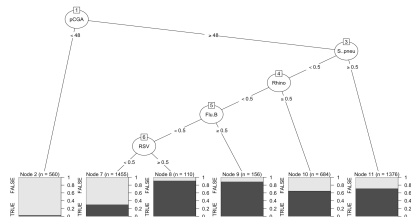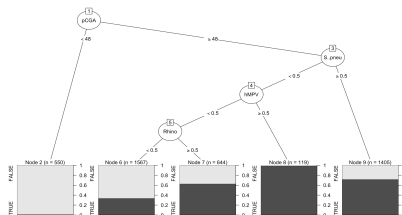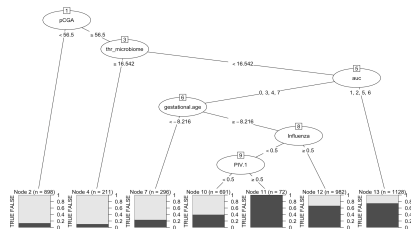# Results: threshold



threshold = 0.5

# Results

|            | is_illness |       |
|------------|------------|-------|
| Prediction | FALSE      | TRUE  |
| FALSE      | 1458       | 782   |
| TRUE       | 771        | 1309  |

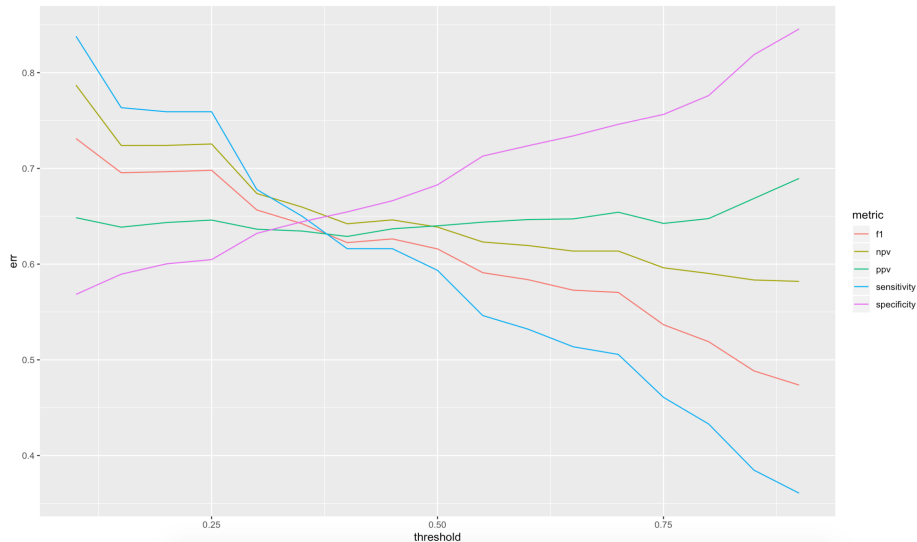| Accuracy | 0.6405 |
|----------|--------|
| PPV      | 0.6293 |
| NPV      | 0.6509 |
| F1       | 0.6277 |
| AUC      | 0.7034 |

# Variation

# Random Forest

For $b = 1$ to $B$:

- Draw a bootstrap sample **Z** of size $N$ from the training data
- Grow a random-forest tree $T_b$ to the bootstrapped data, by recursively repeating the following steps for each terminal node of the tree, until the minimum node size $n_{min}$ is reached

    - Select $m$ variables at random from the $p$ variables
    - Pick the best variable/split-point among the $m$
    - Split the node into two child nodes

- Output the ensemble of trees $\{T_b\}_1^B$

- Classification prediction: $\hat{C}_b(x) =$ the prediction of the $b$th random-forest tree, then

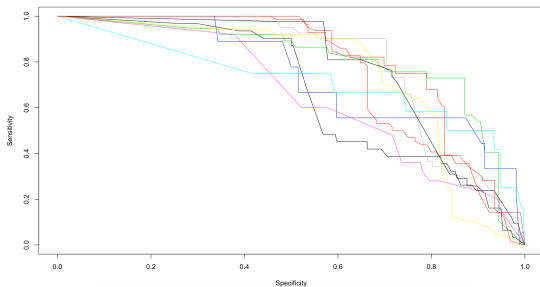$$\hat{C}_{rf}^B(x) = \text{majority vote}\{\hat{C}_b(x)\}_1^B$$

# Results



threshold = 0.375

# Results

| | is_illness | |
|---|---|---|
| Prediction | FALSE | TRUE |
| FALSE | 1609 | 828 |
| TRUE | 620 | 1280 |

| | Tree | Forest |
|---|---|---|
| Accuracy | 0.6405 | 0.6661 |
| PPV | 0.6293 | 0.6737 |
| NPV | 0.6509 | 0.6602 |
| F1 | 0.6277 | 0.6387 |
| AUC | 0.7034 | 0.7352 |

# Q&A