

# Slope, Altitude and planted Area are most useful class separation features. Logistic regression model outperforms the decision tree model.

## Beating The Beetle

Paddington Chiguvare, Masana Khosa

### 1 Intro

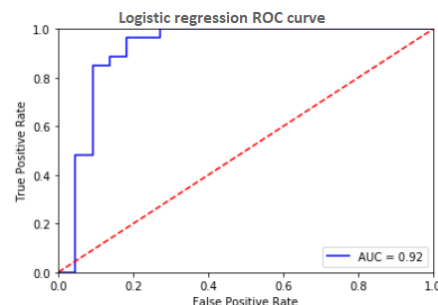
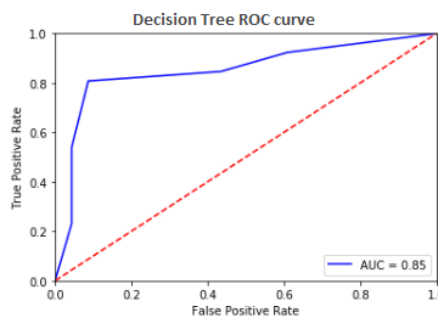
- Eucalyptus plantations in South Africa have been infected by Eucalyptus Beetles.
- A dataset containing comprehensive information about infected and healthy compartments was provided.
- The aim is to develop classification models to classify compartments as healthy or unhealthy and predict sites likely to be threatened.

### 2 Methods

- Logistic regression and decision tree classification models were developed.
- logistic regression model is given by:  $\log\left(\frac{p(X)}{1-p(X)}\right) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$
- To ensure stability, cross-validation was used to train the models. The models were tested using the testing set.

- PCS framework was used to test reliability, reproducibility and transparency.
- Area Under ROC curves were used to measure performance.

### 3 Results



- Logistic regression model has a better performance.

### Extra figures

#### Performance comparison

|                 | Predicted Positive | Predicted Negative |
|-----------------|--------------------|--------------------|
| Actual Positive | 16                 | 6                  |
| Actual Negative | 6                  | 21                 |

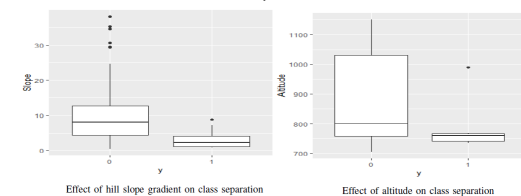
|                 | Predicted Positive | Predicted Negative |
|-----------------|--------------------|--------------------|
| Actual Positive | 16                 | 6                  |
| Actual Negative | 0                  | 27                 |

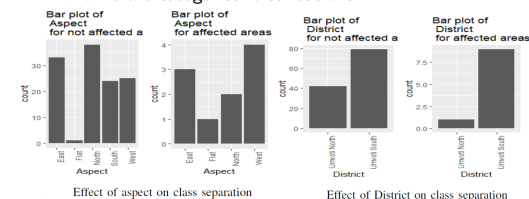
|                     | Precision | Recall | Accuracy |
|---------------------|-----------|--------|----------|
| Decision Tree       | 0.75      | 0.75   | 0.88     |
| Logistic regression | 0.91      | 0.86   | 0.75     |

#### Useful features

- Slope and Altitude are numerical features useful in class separation since their box plot notches do not overlap.



- Aspect and District are categorical features useful in class separation because they have different categorical distributions.



### Deployment of Models

- The two classification models were deployed using shiny and R.

#### Predict Compartment Condition

Planted Area: 7.43

Site index: 17.6

x0 20 Percent (Percentage of compartment in slope class): 100

x21 35 Percent (Percentage of compartment in slope class): 0

x36 50 Percent (Percentage of compartment in slope class): 0

Land type (Number used to group "homogenous" areas of soil, climate etc.): 220

Altitude: 1003

frost Risk rating 1-9 (low - high): 1

Winter Rain (Rainfall A,M,J,J,A,S % of total): 16.818

hillslope gradient: 13.2947

Logistic regression Prediction

Compartment to unhealthy

Decision tree prediction