# Predicting Malware Infection of Windows Machines

Pavel Zimin, PhD
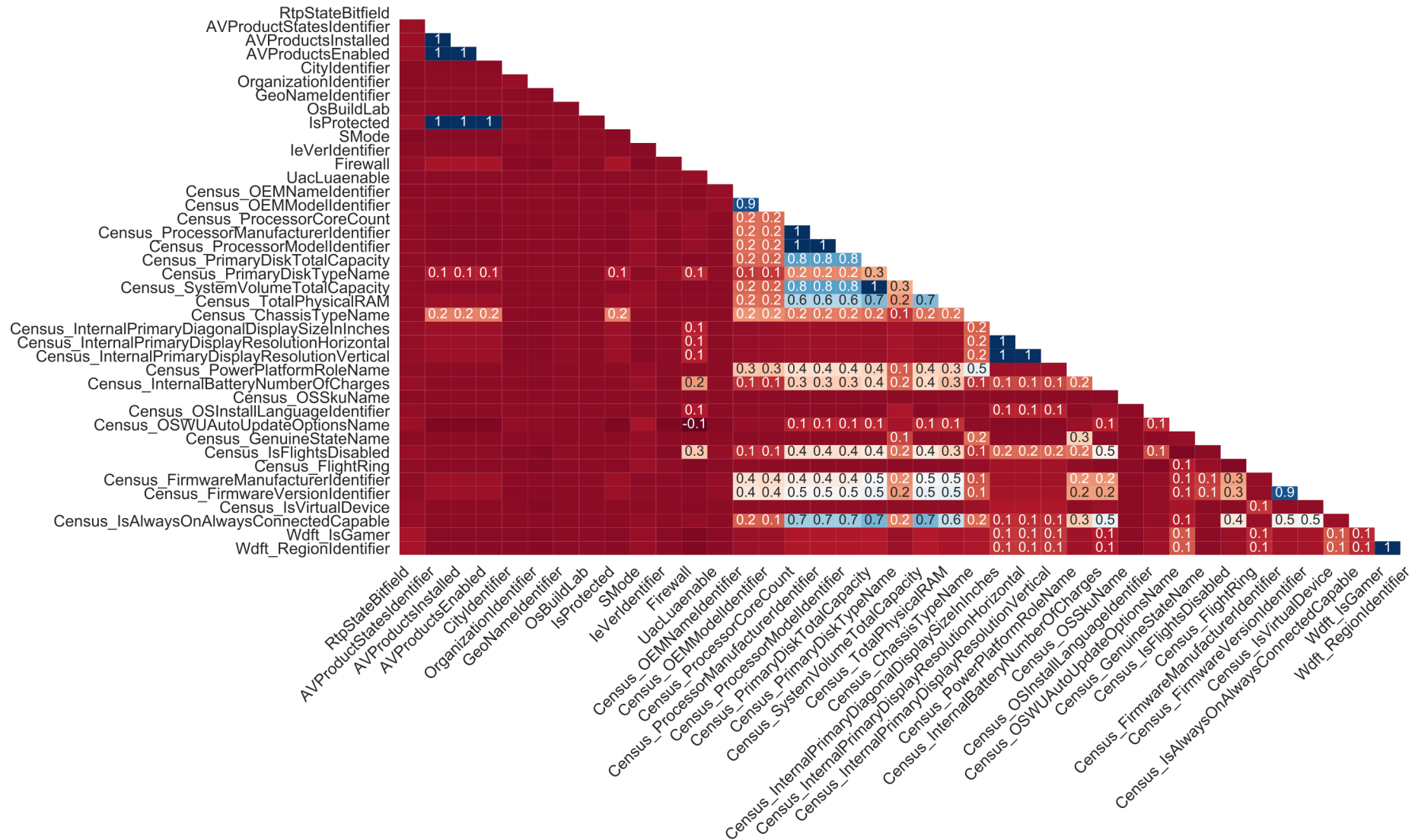
Mentor: Ramkumar Hariharan, PhD

# Problem Statement

- Computer infection by malware constitutes a serious security problem

- The ability to predict the chances of malware infection before they occur would benefit consumers and businesses
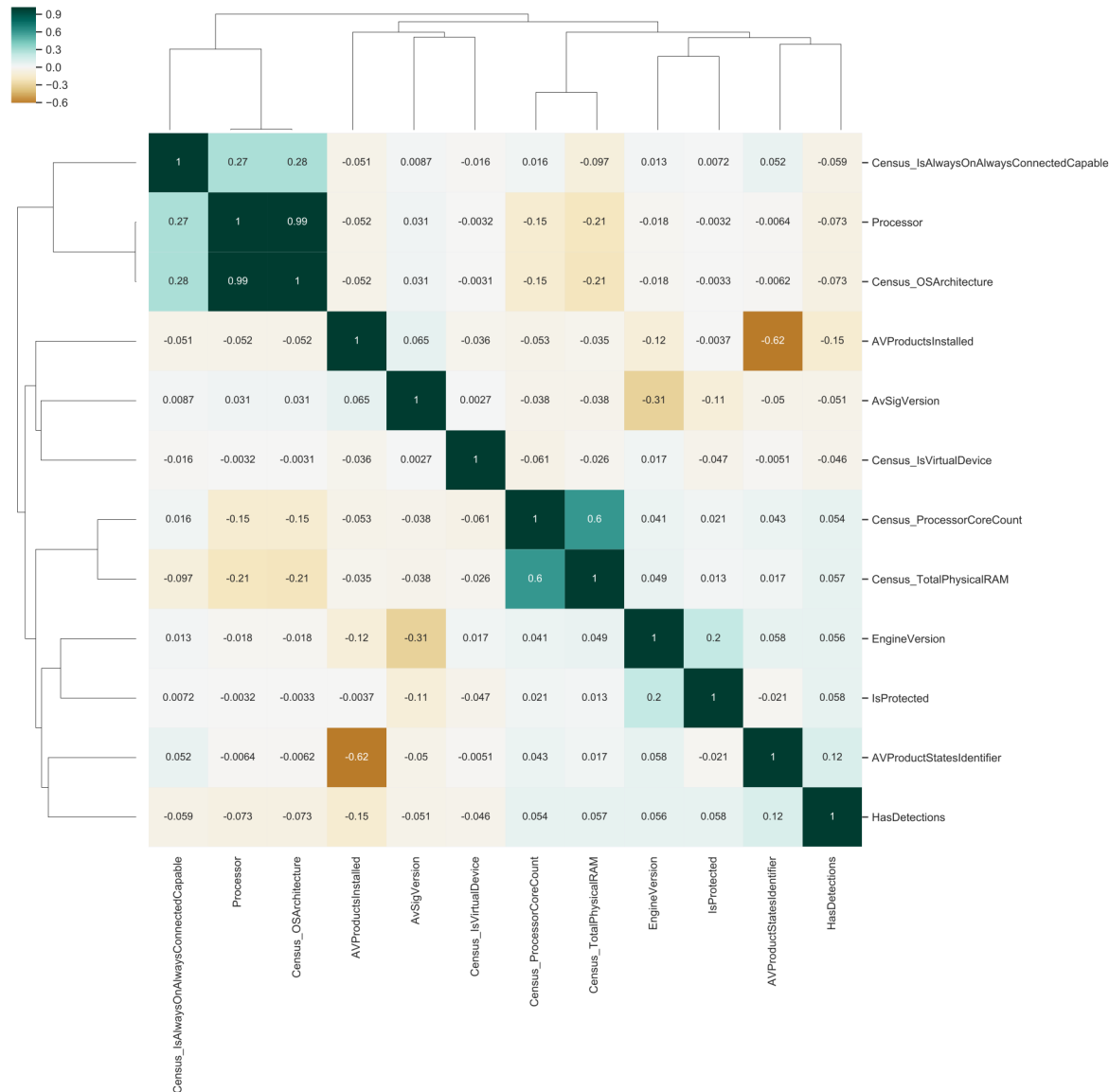
**Business Use Scenario:**

- This project would benefit software manufacturers who would be able to incorporate the model into their software that would allow for additional security measures aimed at preventing the infection by malware
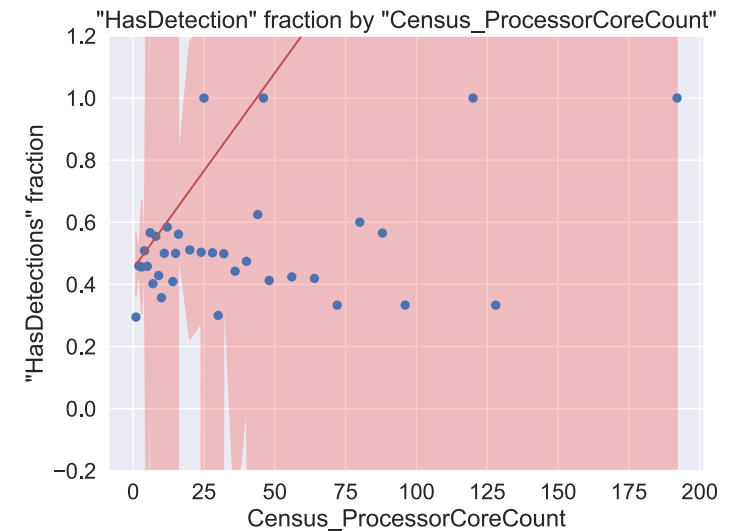
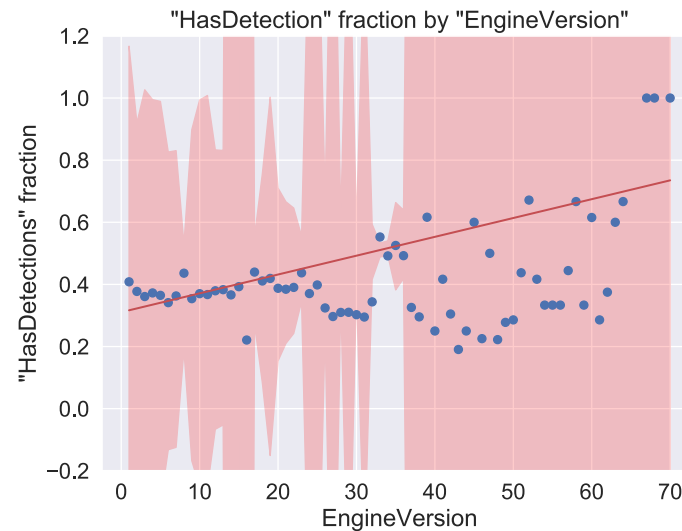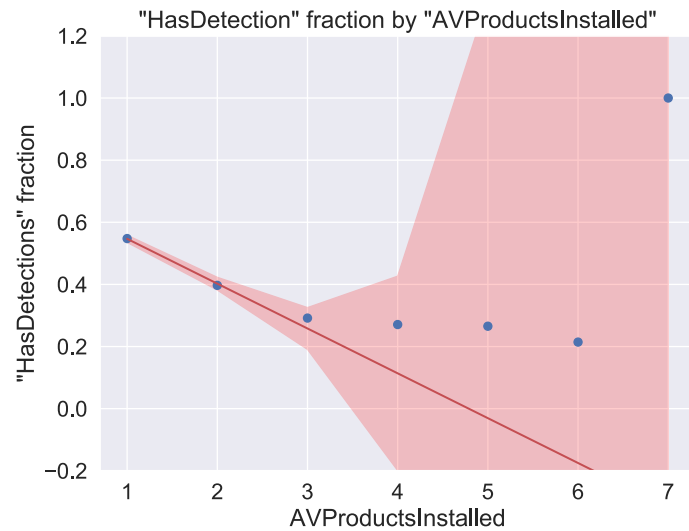# Correlation heatmap for missing values

# Correlation Cluster Heatmap



- The figure shows only the features with the highest positive or negative Pearson correlation coefficients with the target variable
- Overall low correlation of features with the target variable
- The following features show the highest correlation with target variable:
  - 'AVProductsInstalled'
  - 'AVProductStatesIdentifier'
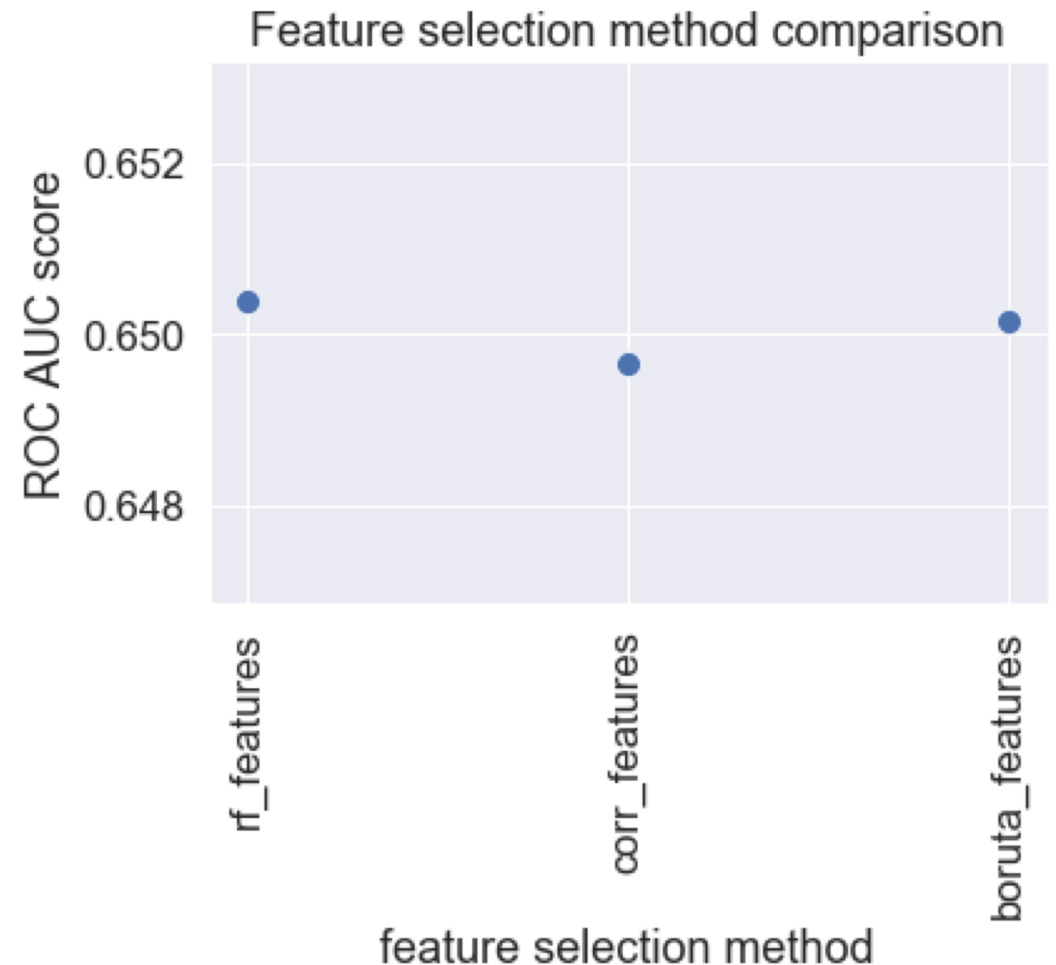  - 'Processor'
  - 'Census_OSArchitecture'

# Frequencies of target variable by selected features



Weighted least squared models were fitted to the data with the weights of the value count for each data point.
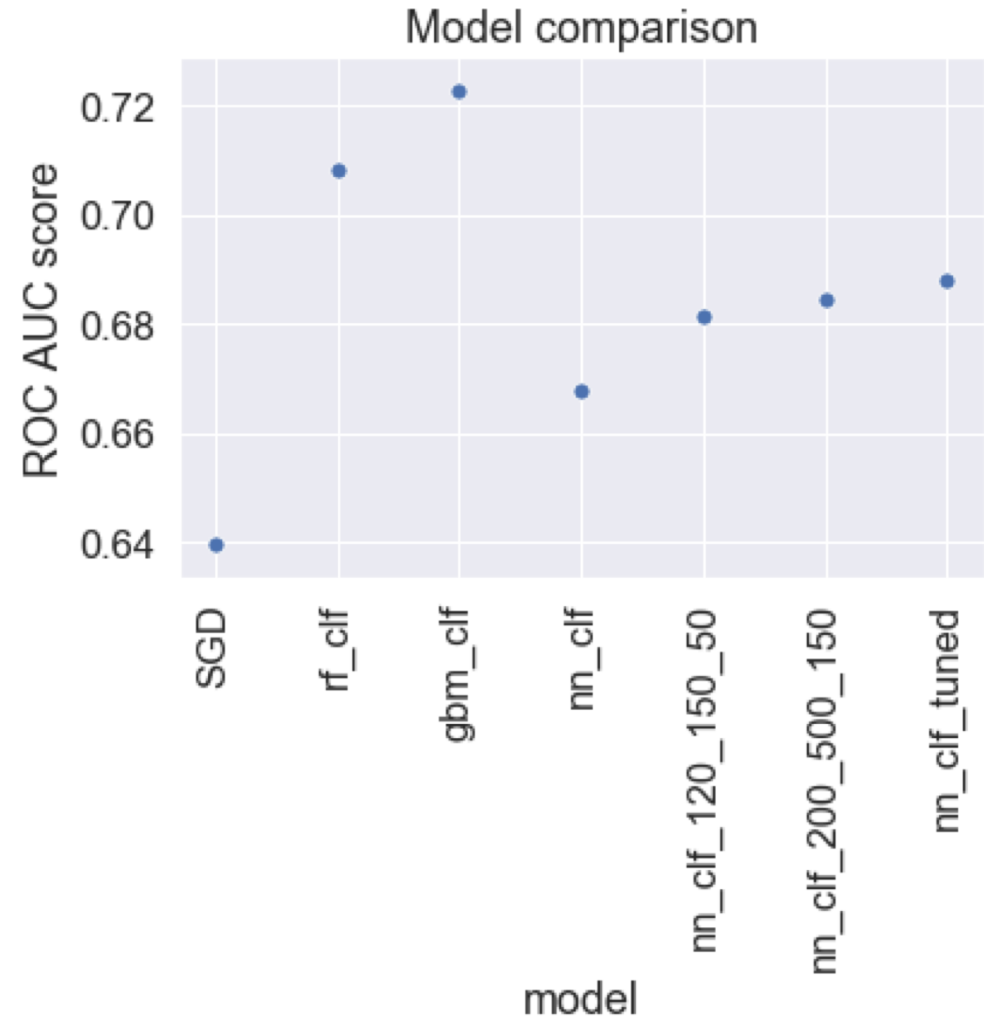
# Feature Selection

- 3 methods of feature selection were performed:
    - Random Forest's feature importances
    - Elimination of highly correlated features
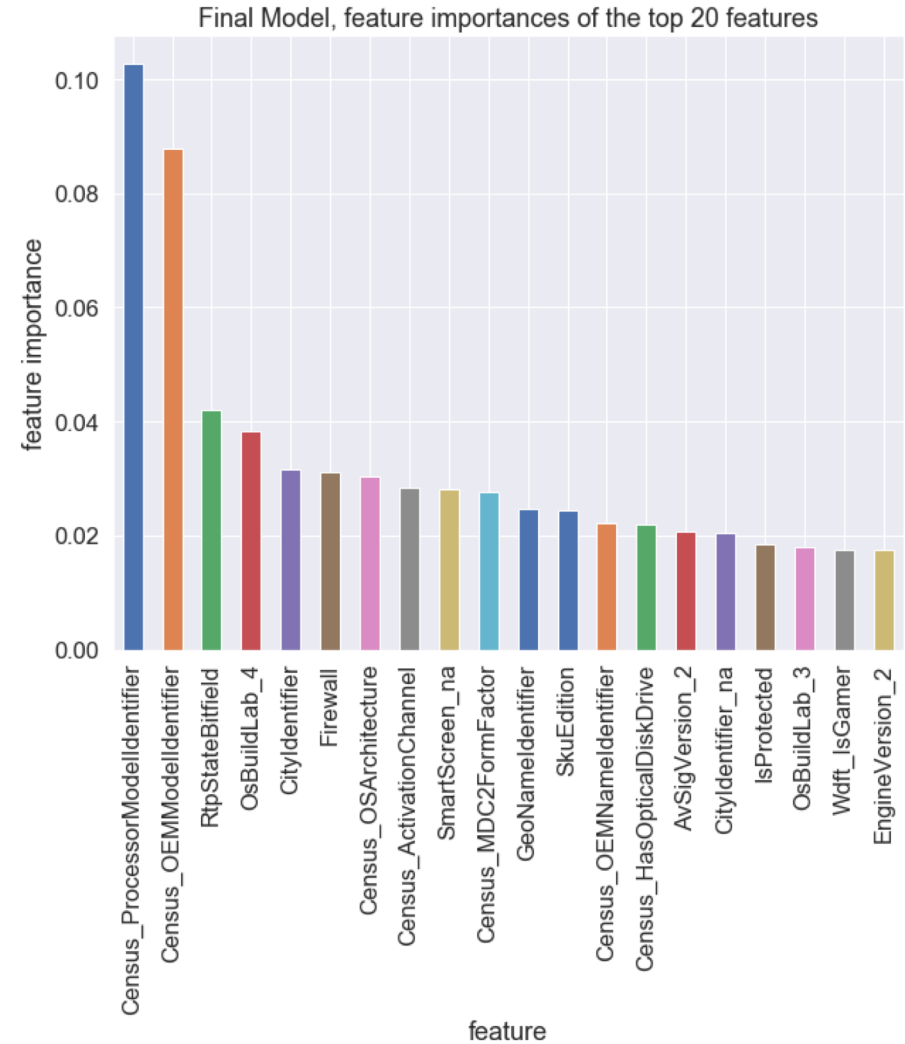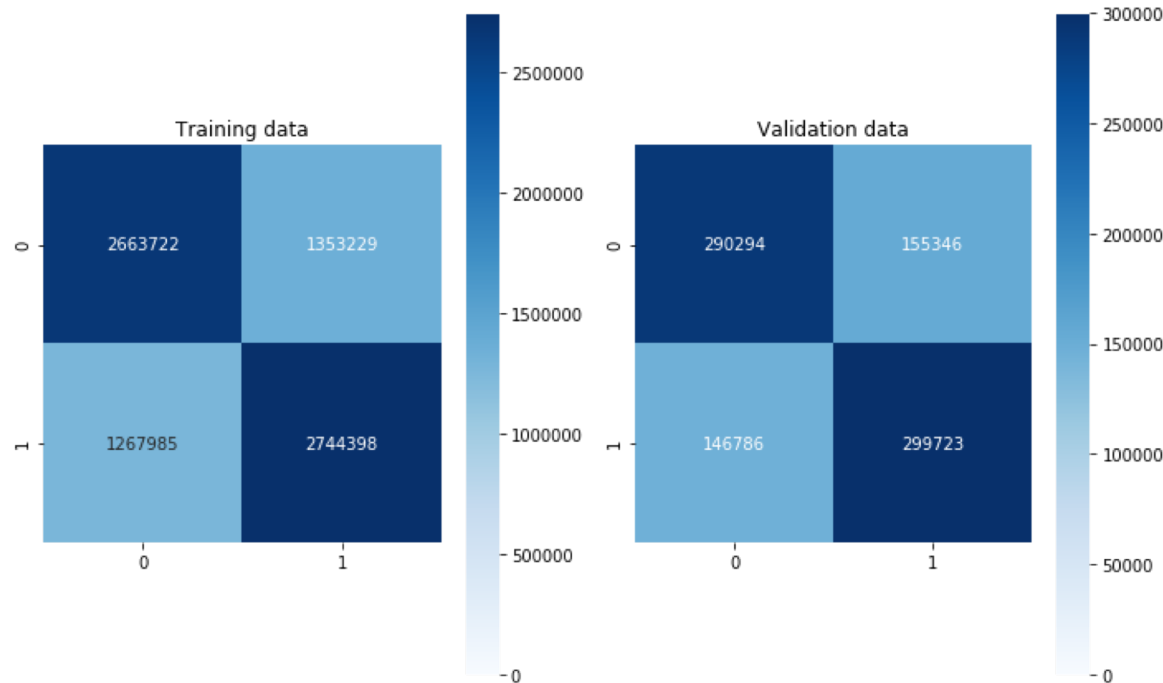    - Boruta algorithm

# Model Selection

- Logistic Regression (SGD)
- Random Forest (rf_clf)
- Gradient Boosting (gbm_clf)
- Neural networks (nn_)



Model comparison

# Final Model Evaluation and Feature Importances

- ROC AUC score 0.74 on test data

# Conclusions

- Related features tend to have the same observation with missing values.

- Features show low overall correlation with the target variable.

- Weighted least squared models show significant relationship between selected features and the target variable.

- Features were selected based on random forest's feature importances.

- Best performing model was Gradient Boosting with ROC AUC score 0.74.

- The most important features are: Census_ProcessorModelIdentifier, Census_OEMModelIdentifier, RtpStateBitfield, OsBuildLab_4, CitiIdentifier, Firewall.