

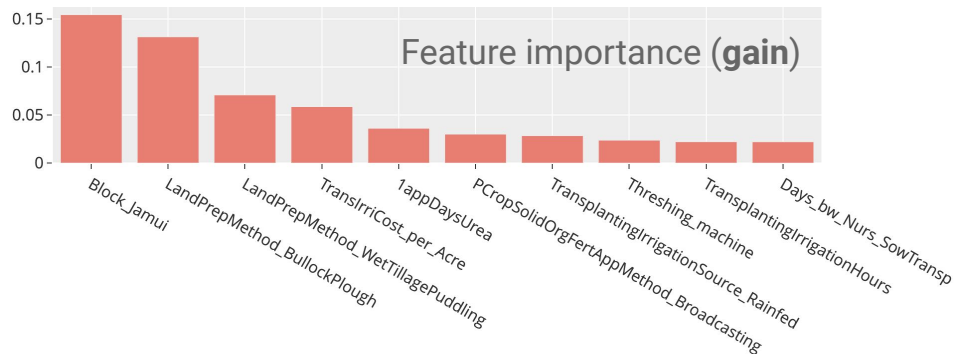
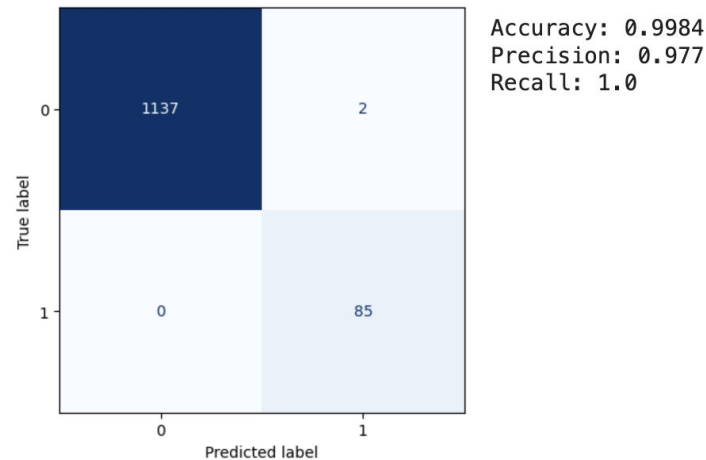
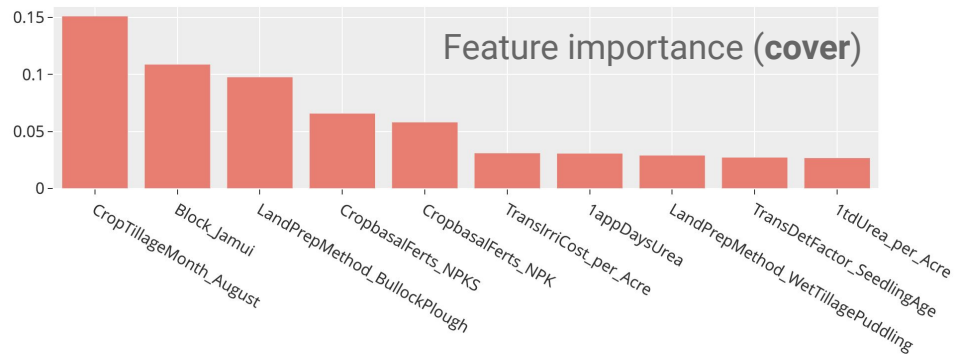
Classification

Outliers vs. non-outliers on yield per acre

A dark blue diagonal gradient bar that starts from the bottom left and extends towards the top right, covering the lower half of the slide.

Goal: predicting which rows are outliers (1) and which rows are not (0) in the test set

- **New binary variable:** “1” for the 21 outliers, “0” for others (determined by $\text{Yield_per_Acre} > 5000$)
- **Class imbalance:** duplicated outlier rows * 10
 - → 231:3849 outlier : non-outlier ratio
- **Train-test split:** only used the train set, and did a 70-30% split
 - number of class 1 in train set: 170
 - number of class 1 in test set: 61
- **Columns dropped:** ID, Set, Yield, Yield_per_Acre, Group_Outlier
 - are there any others I should have removed?
- **Model:** XGBoost Classifier
 - parameters (almost default, did not do much to optimize it): $\text{random_state}=0$, $\text{scale_pos_weight}=1$, $\text{min_child_weight}=1$, $\text{max_depth}=7$, $\text{learning_rate}=0.1$, $\text{gamma}=0.1$, $\text{colsample_bytree}=0.9$, $\text{subsample}=0.7$, $\text{eval_metric}=\text{"error"}$

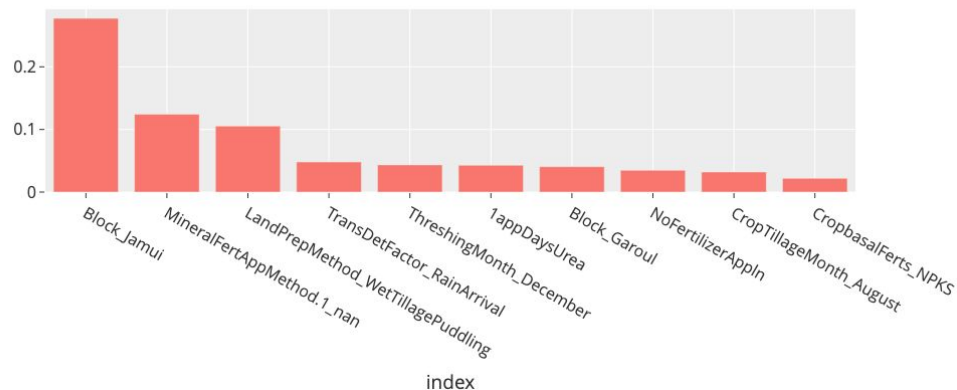


Model 1

deduplicated outlier rows * 10; training set size = 2856 (146 class 1), train set size = 1224 (85 class 1)

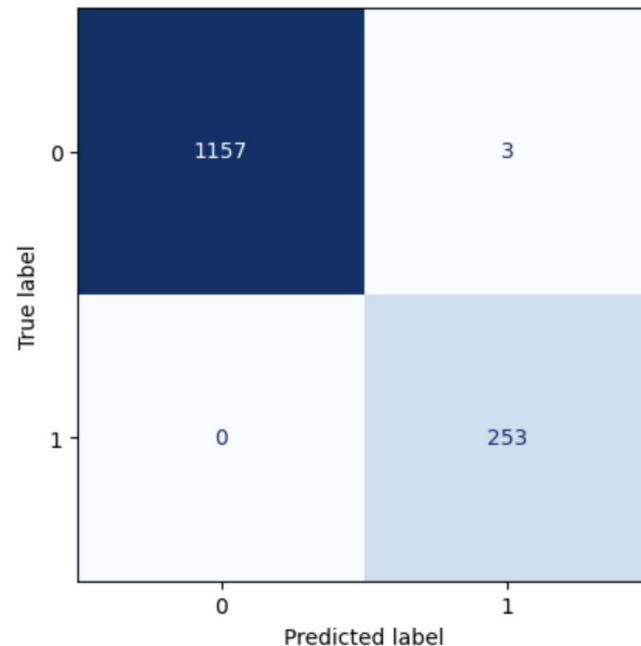
Number of class 1 predictions in the actual test set: 6 (out of 1290) (all Jamui; all LandPrepMethod BullockPlough & WetTillagePuddling = True)

Feature importance (gain)



Number of class 1 predictions in the actual test set: 11
(which include those identified by Models 1 and 2)

Accuracy: 0.9979
Precision: 0.9883
Recall: 1.0



Model 3

Number of class 1 predictions in the actual test set: 11 (out of 1290)

Number of class 1 predictions in the actual test set: 11 (out of 1290)

Ran the model to predict outlier/not-outlier on the test set, exported df as
“preprocessed_with_outlier_classif.csv”