

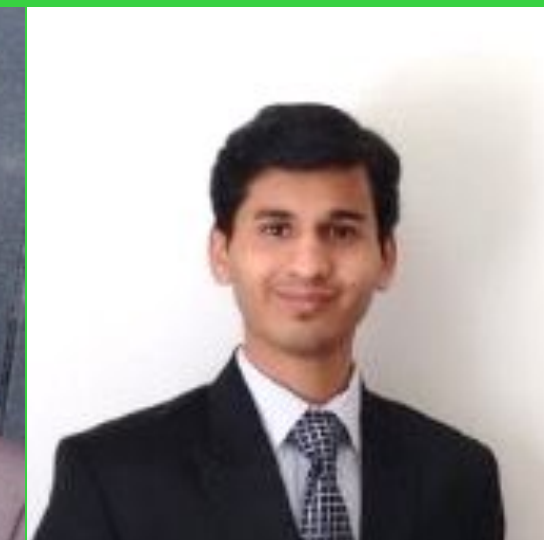
GOOD/BAD CAR PURCHASE FOR AUTO DEALERS



Final Presentation for I526



Nihar Khetan



Ghanshyam Malu



Xiao Liang

Problem Statement

When we go to buy a second hand car we expect to get a good car and dealers want profits as well.

How?

Aim

Trying to predict that whether a
an auction is a good buy or not



bought by a dealer at

Removed Redundant Features

Handled Null/Missing Values

Continuous Data - Took Average

Discrete Data - Created new Category NULL

Removed Features with More Than 95% Missing Values

Normalized All Continuous Values

Generated **Balanced Datasets**



Feature Selection

Expert Knowledge

VehOdo
VehicleAge”
“MMRCurrentAuctionCleanPrice”,
“MMRCurrentRetailAveragePrice”,
“WarrantyCost”.

Chi Square Ranks

Unbalanced Data :
Best score for All features

Balanced data :
17 Features

Recursive Feature Elimination

MMRAcquisitionAuctionAveragePrice,
MMRAcquisitonRetailCleanPrice,
MMRCurrentAuctionAveragePrice,
MMRCurrentAuctionCleanPrice,
"WarrantyCost"

Evaluation Criteria

Balanced Data

Accuracy

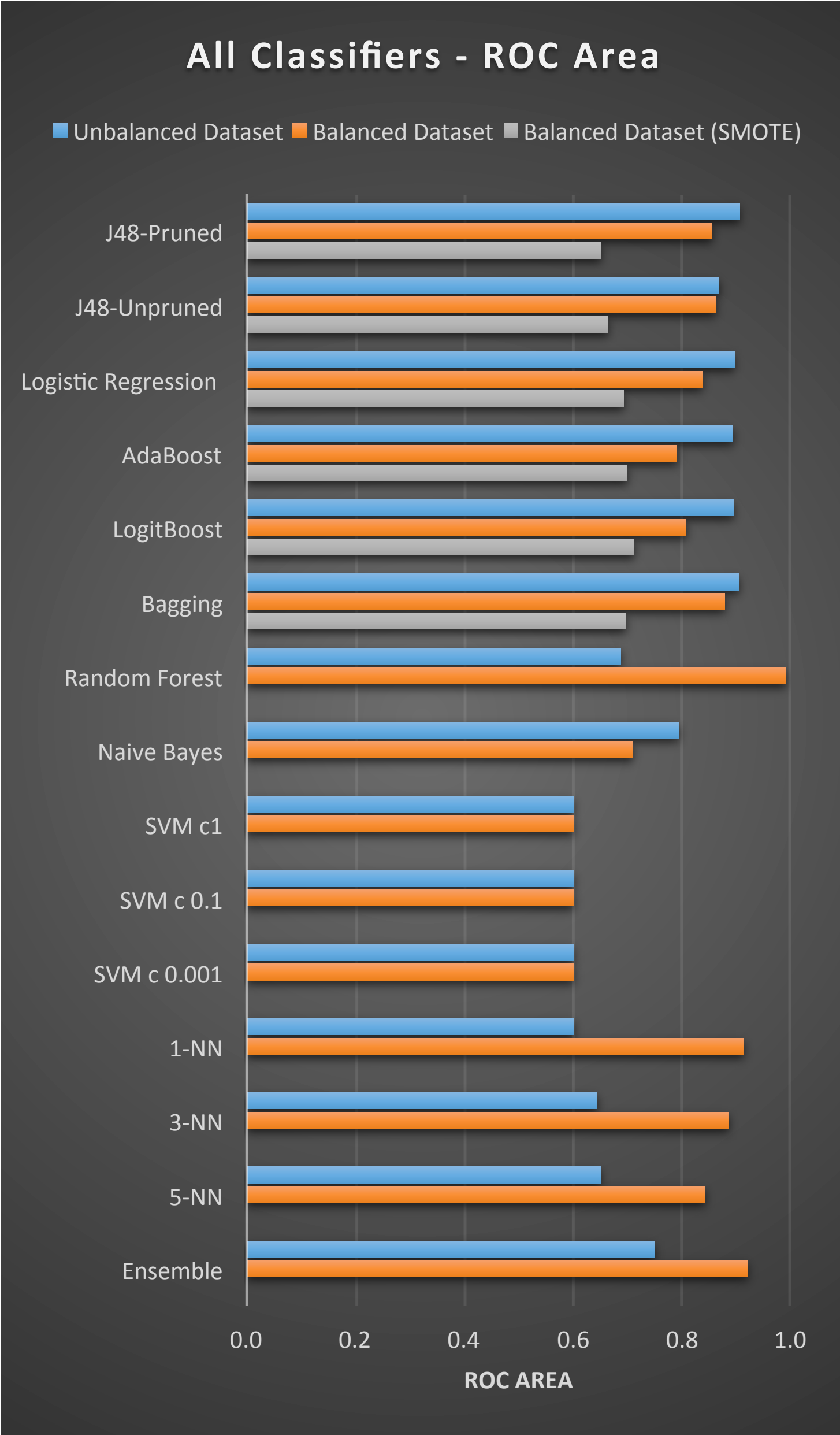
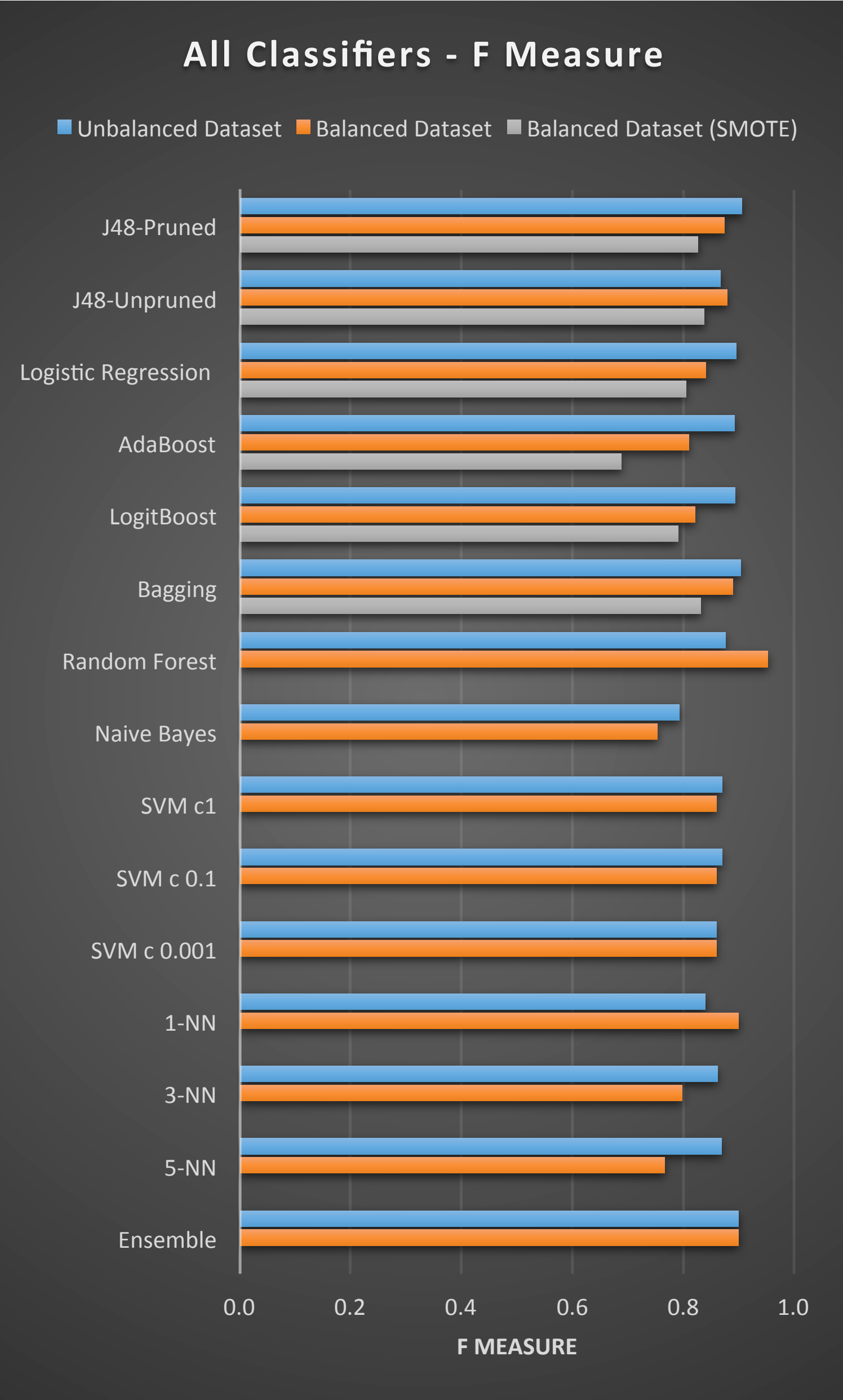
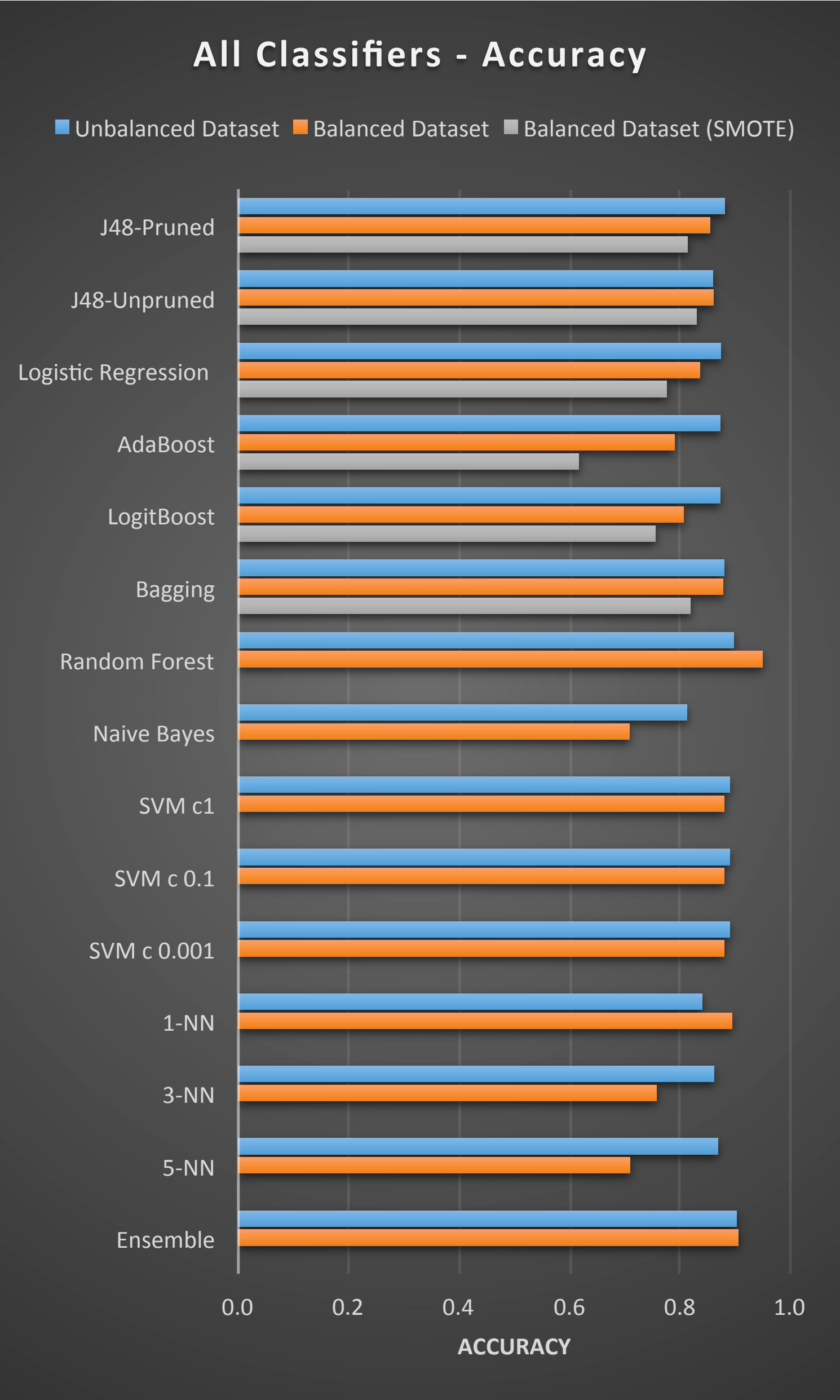
.....

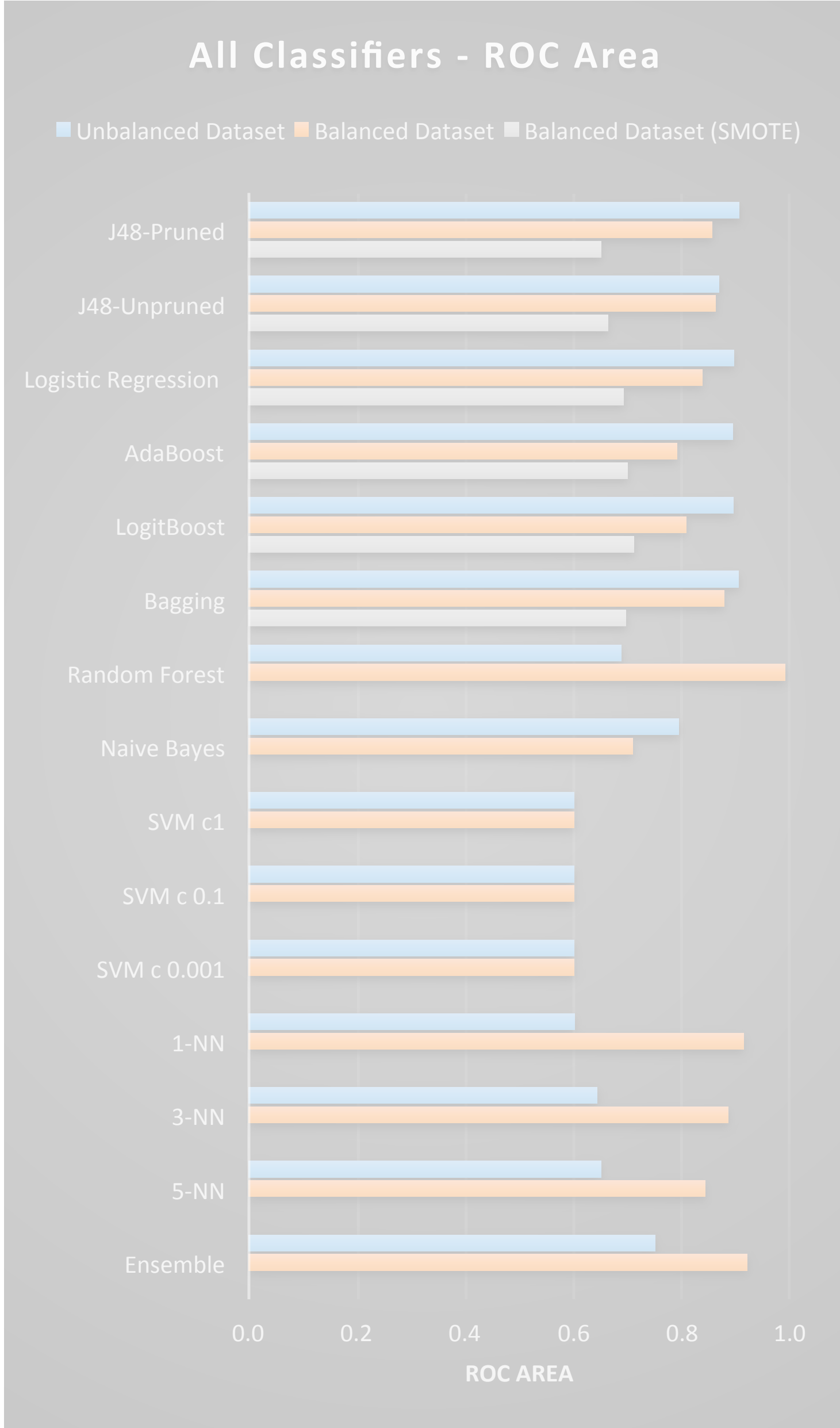
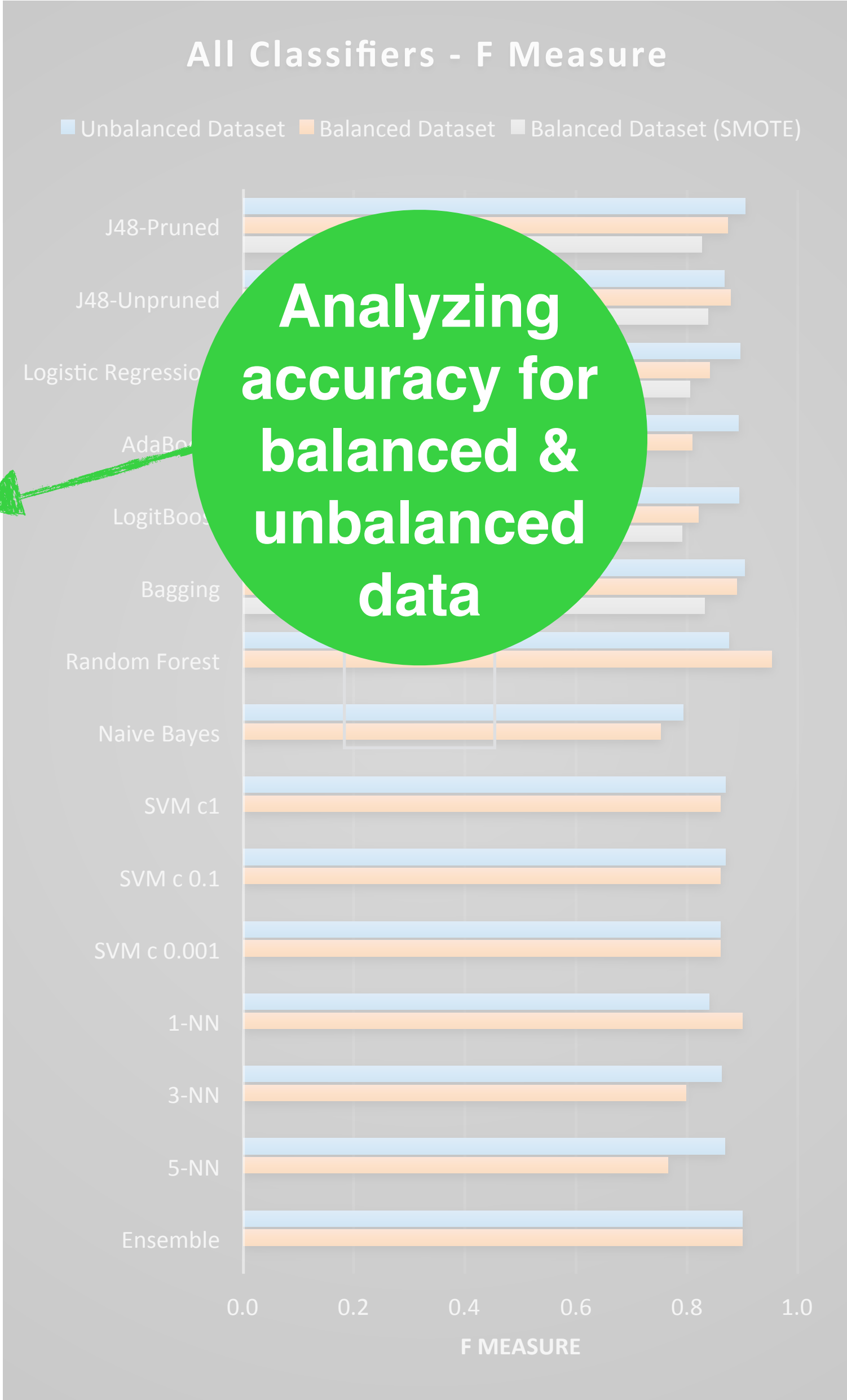
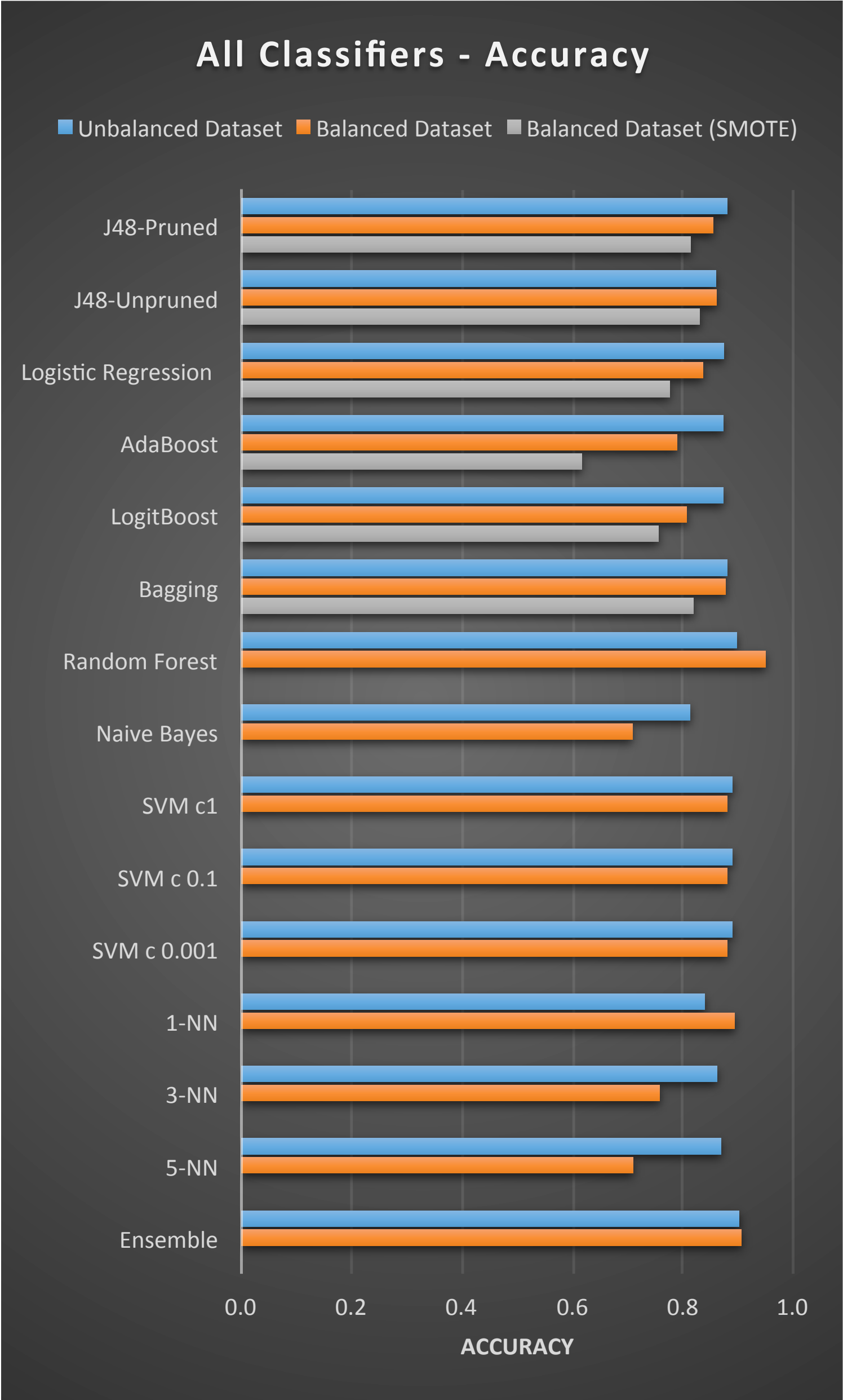
Unbalanced Data

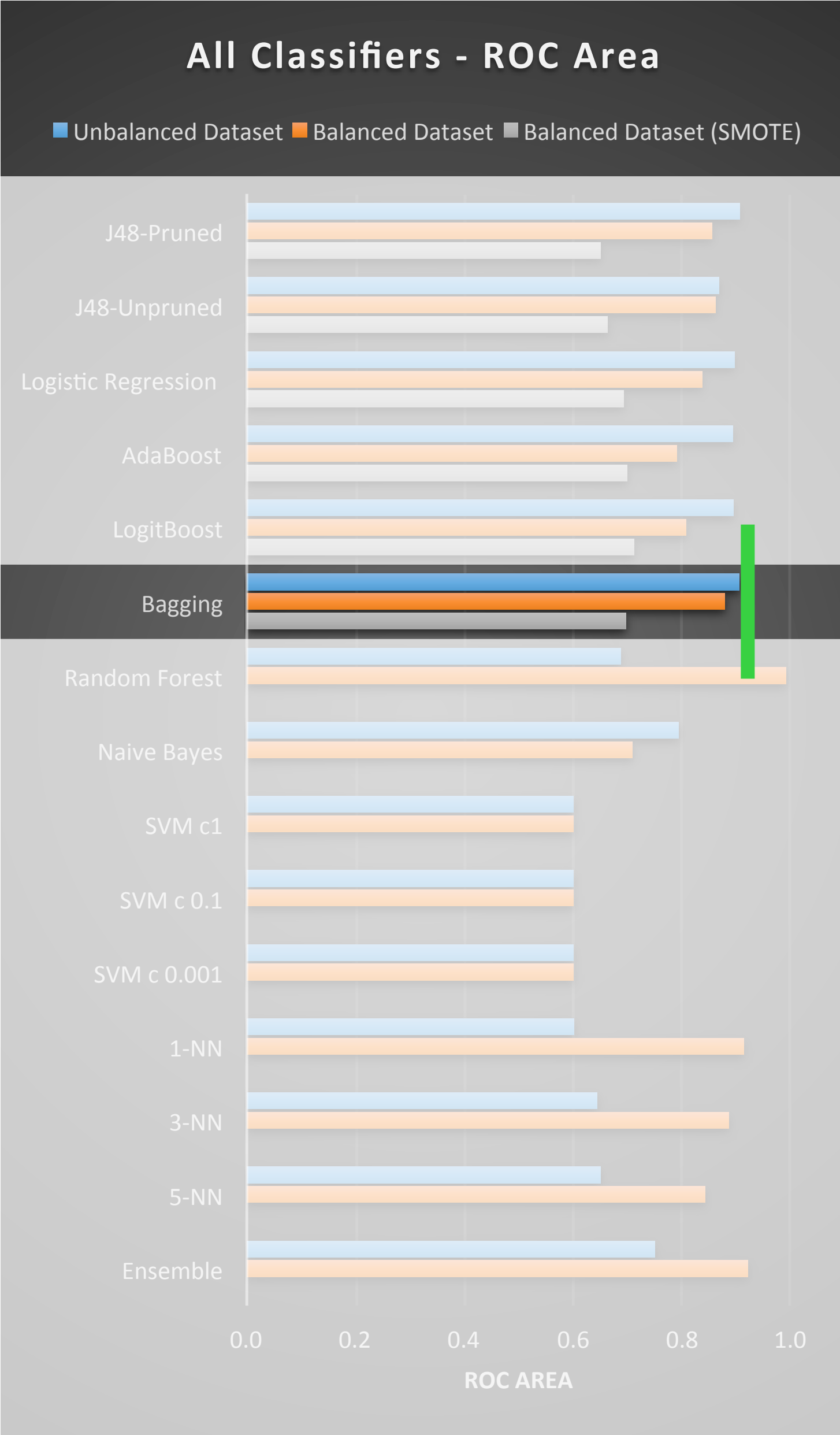
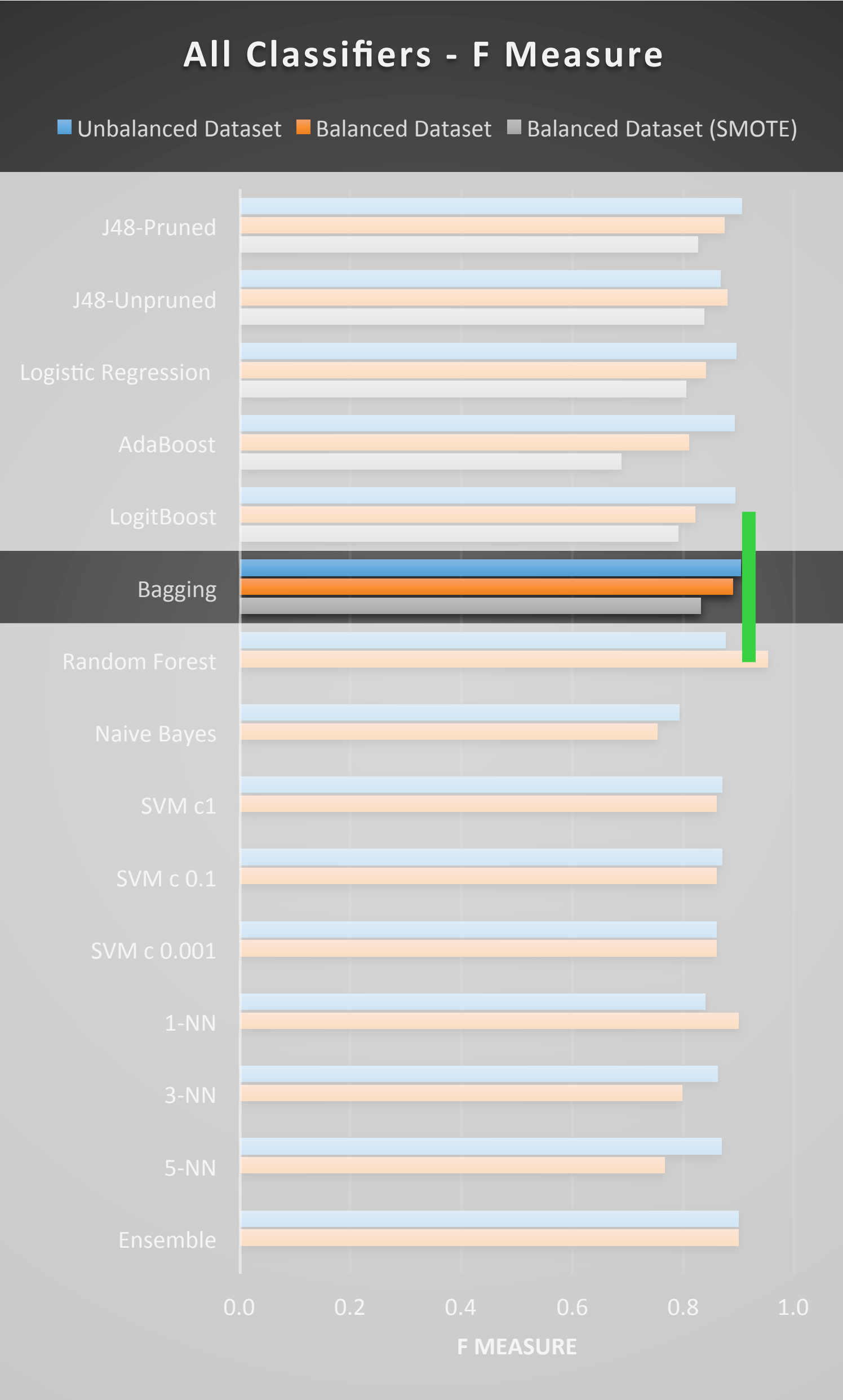
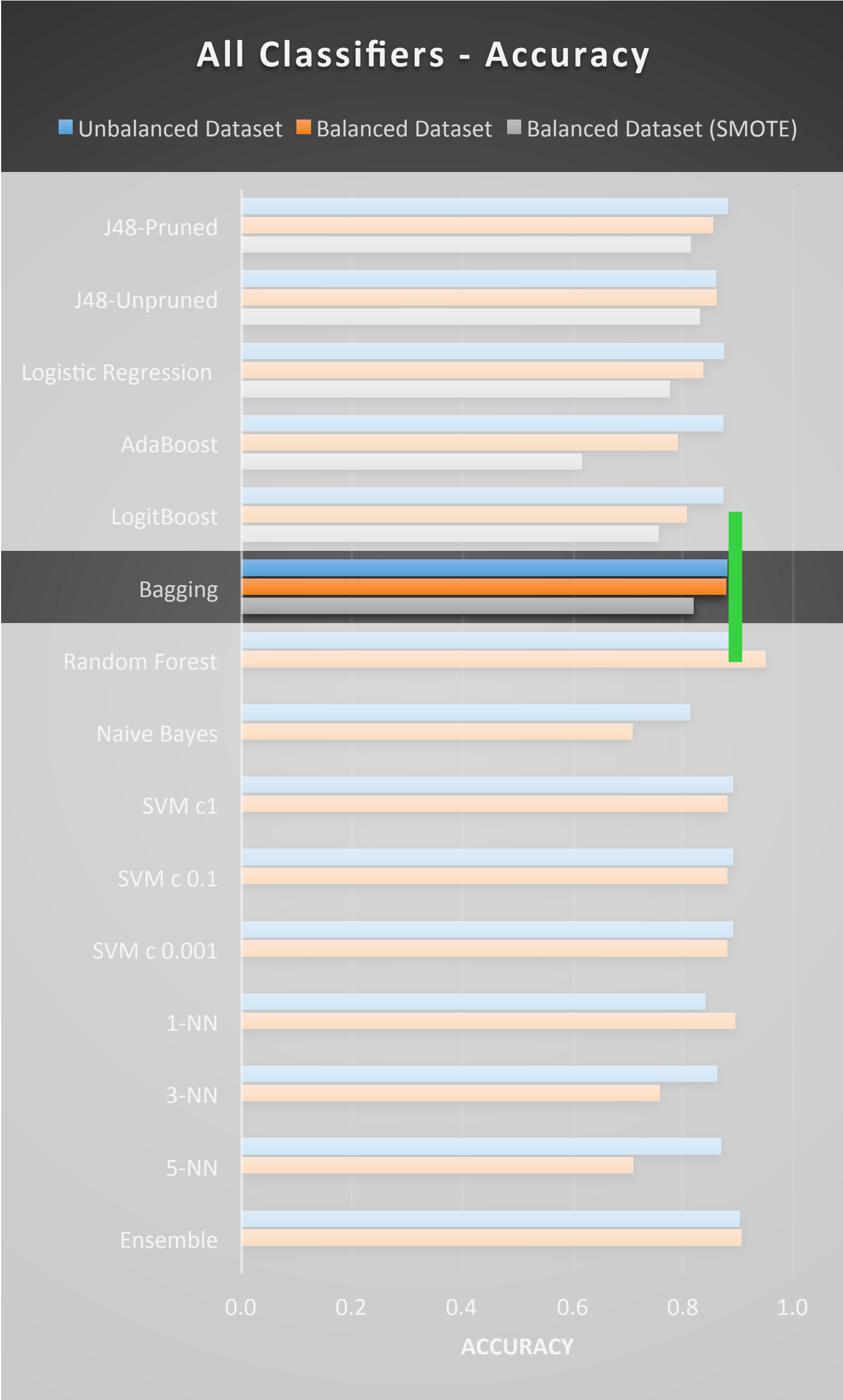
F1 Score

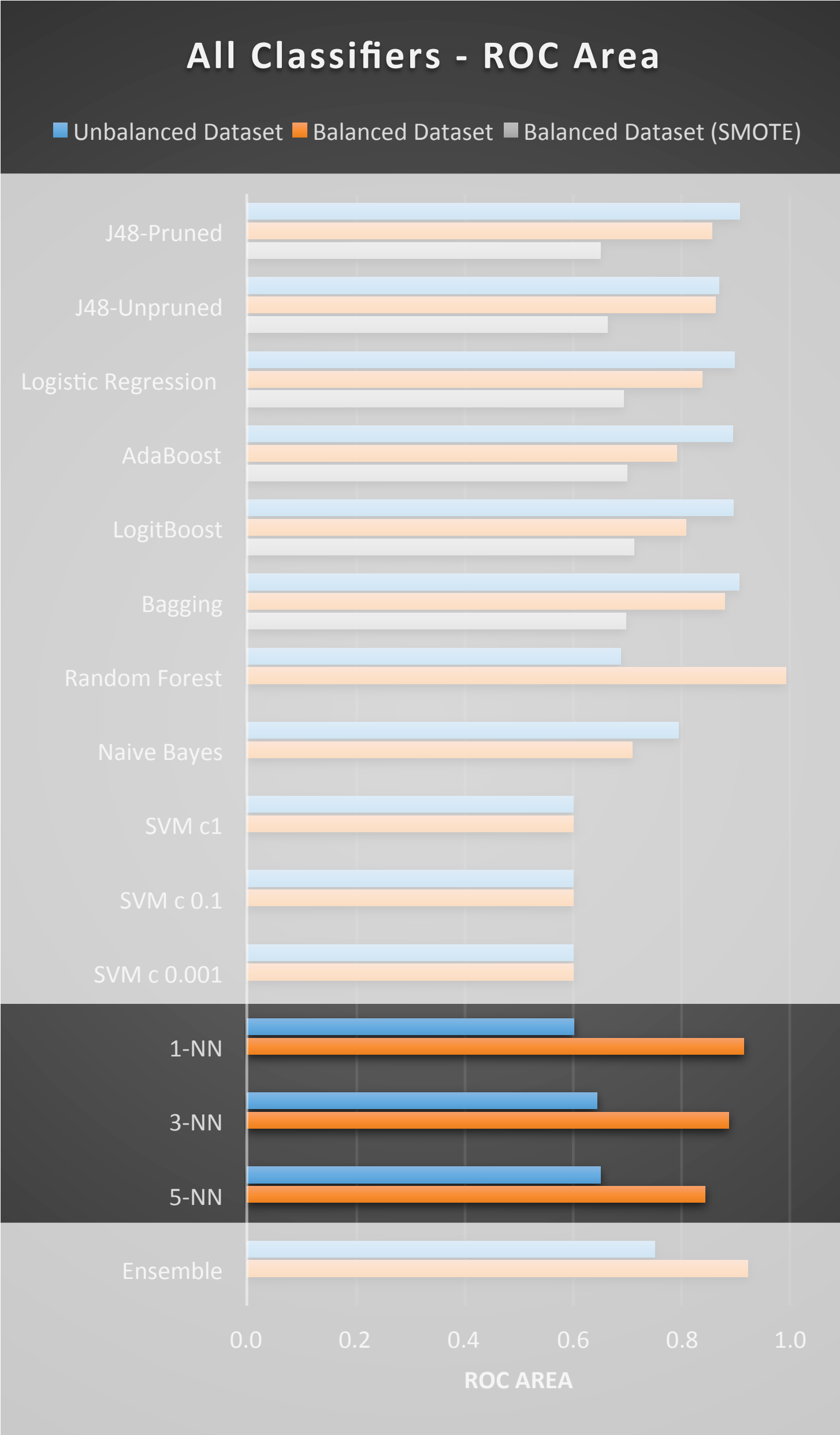
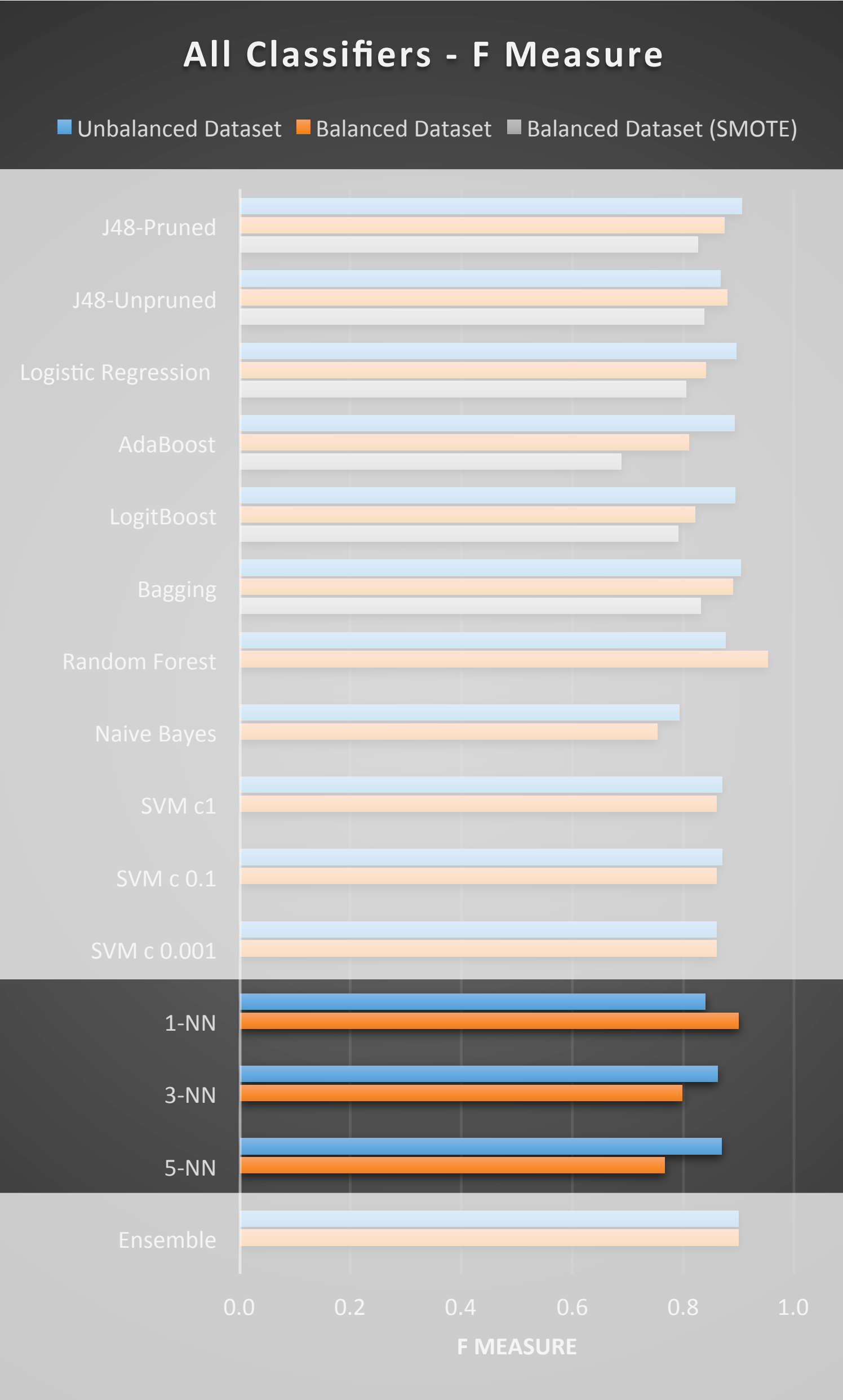
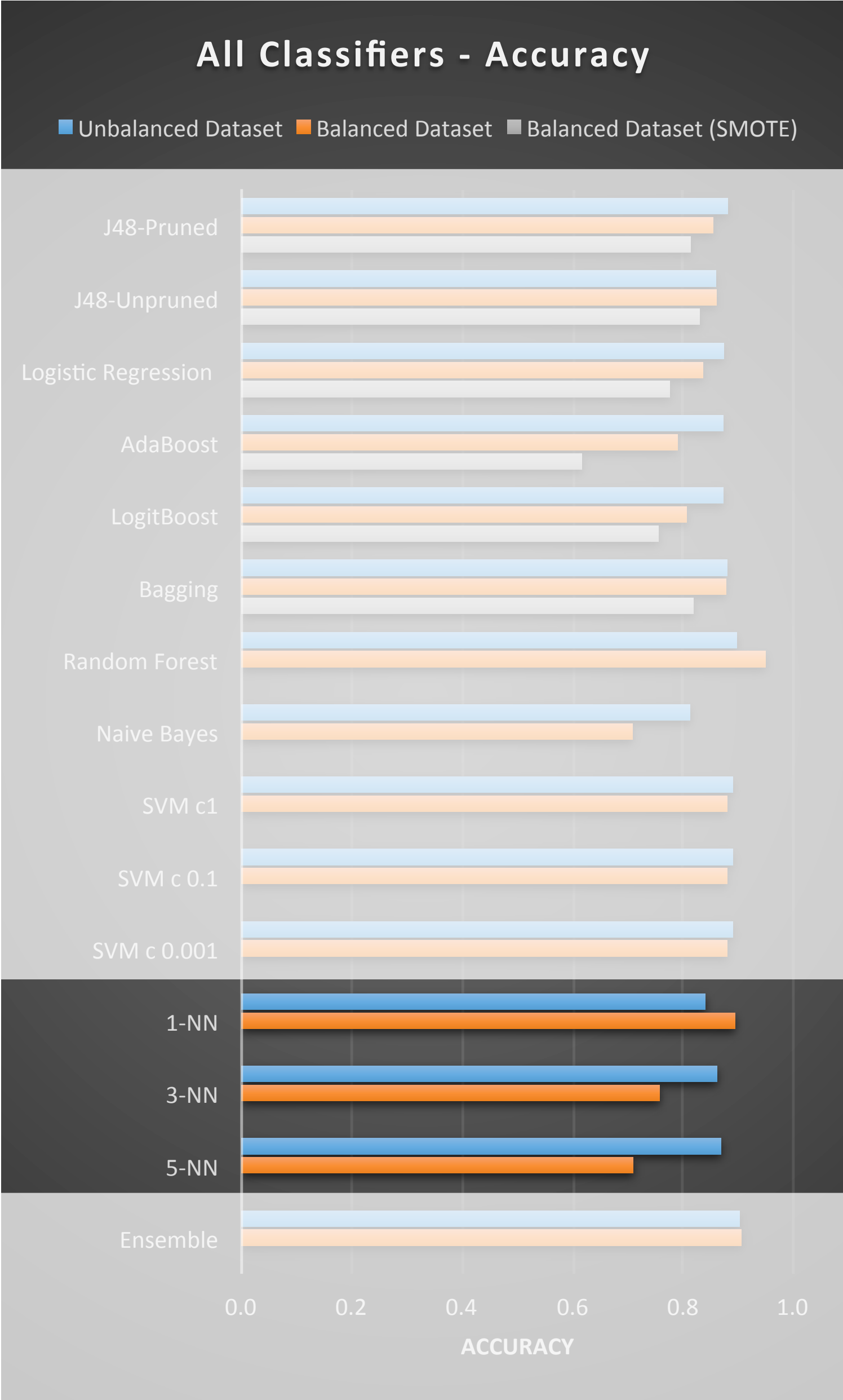
ROC AUC

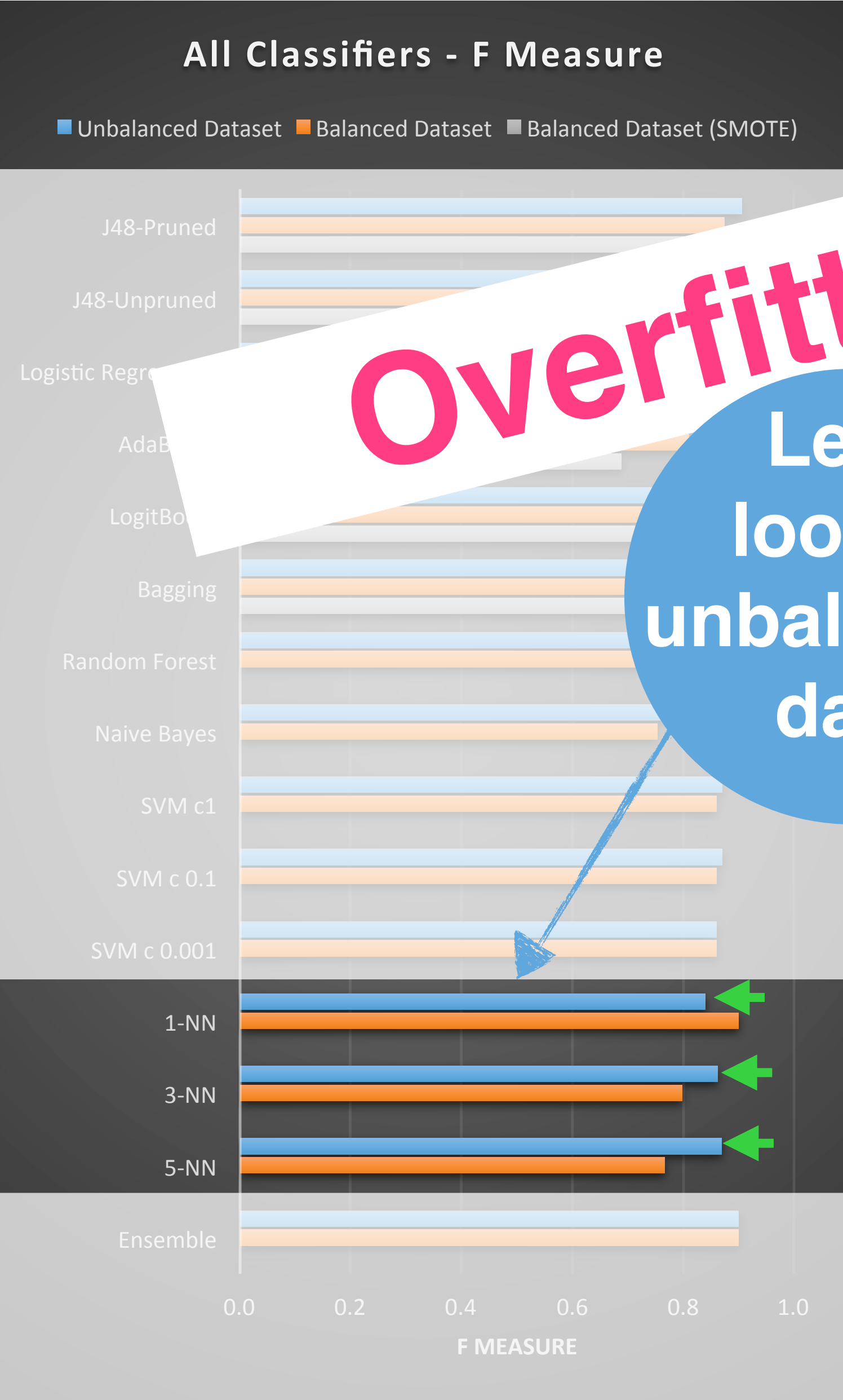
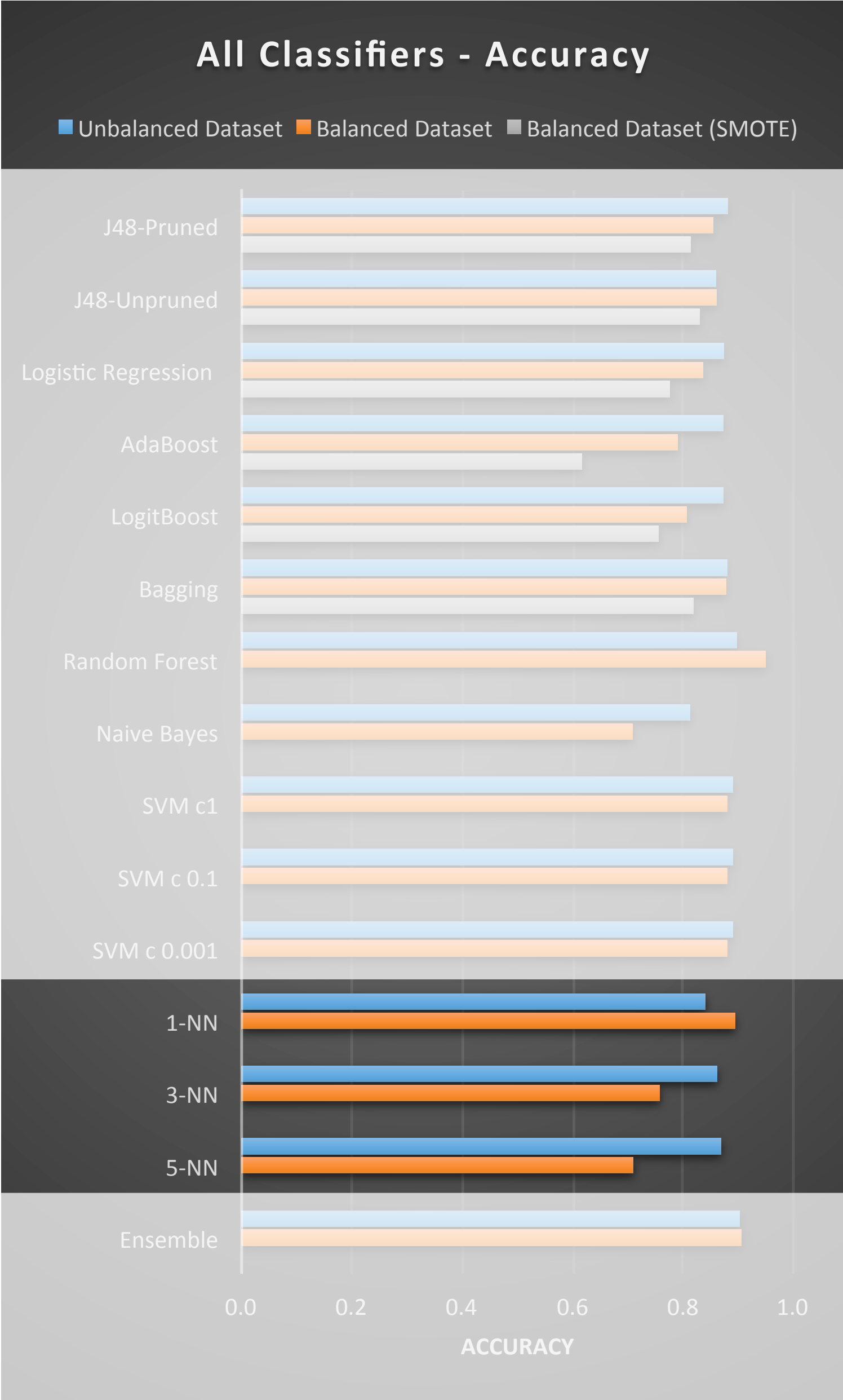
Results & Analysis

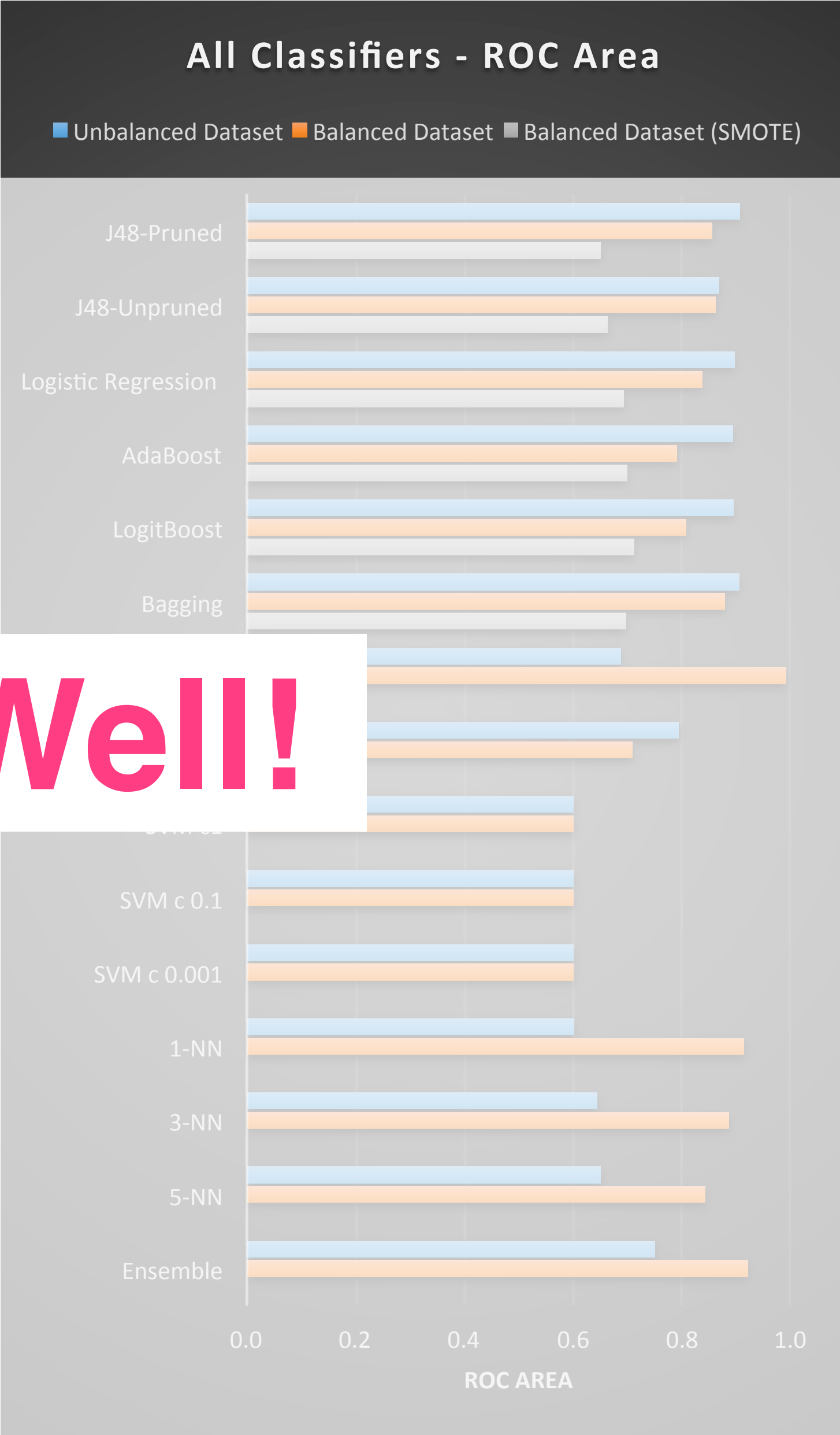
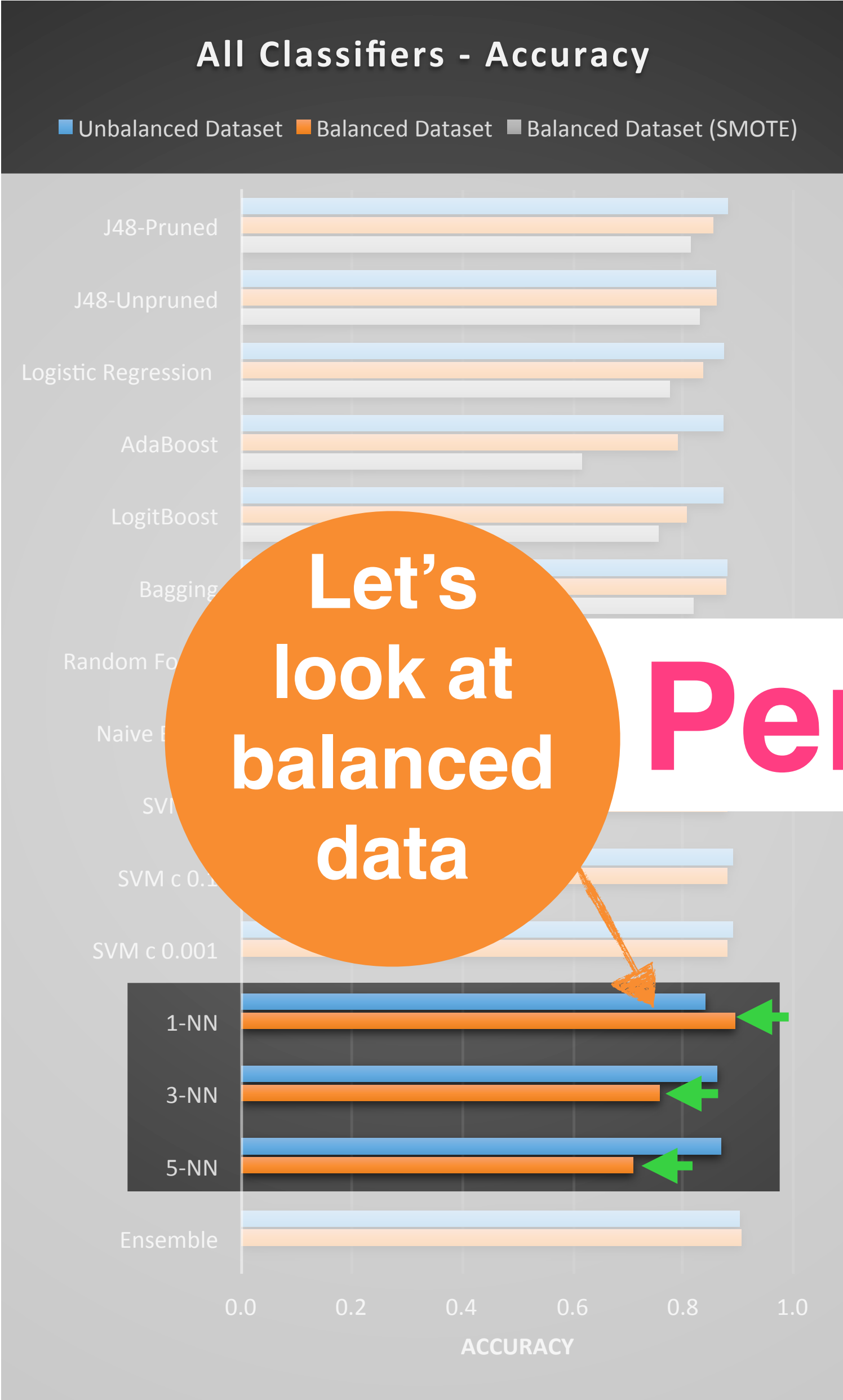






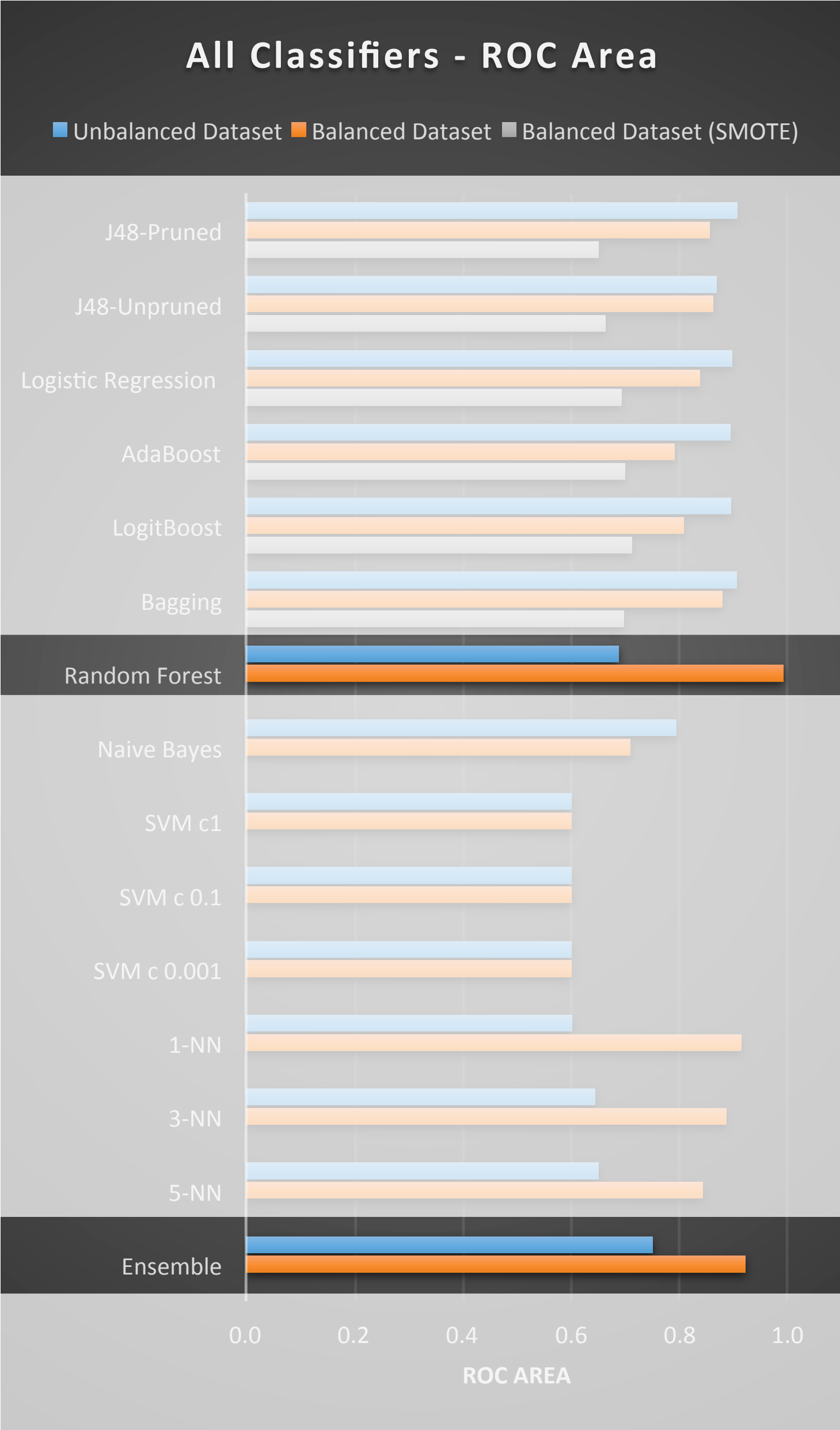
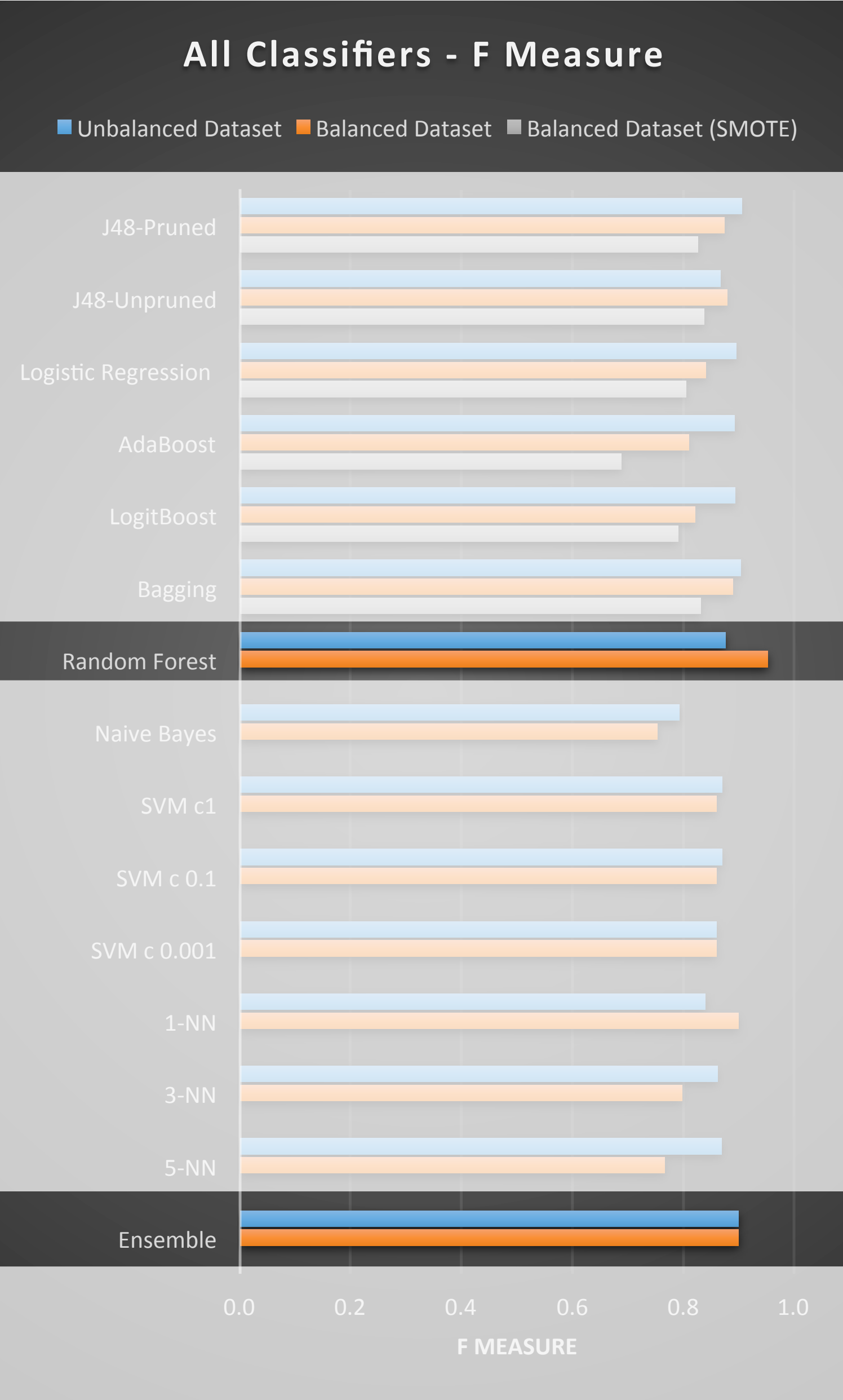
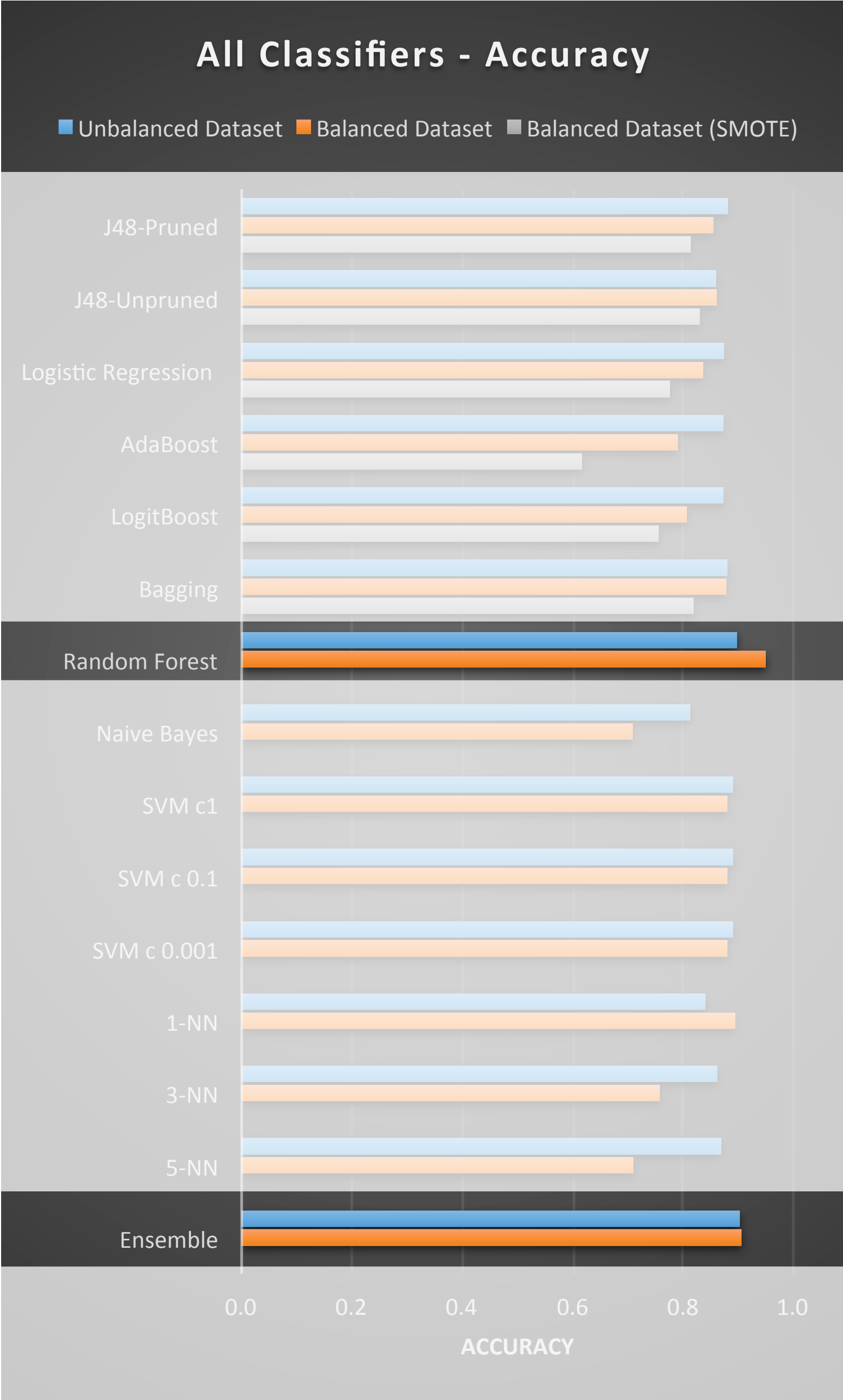


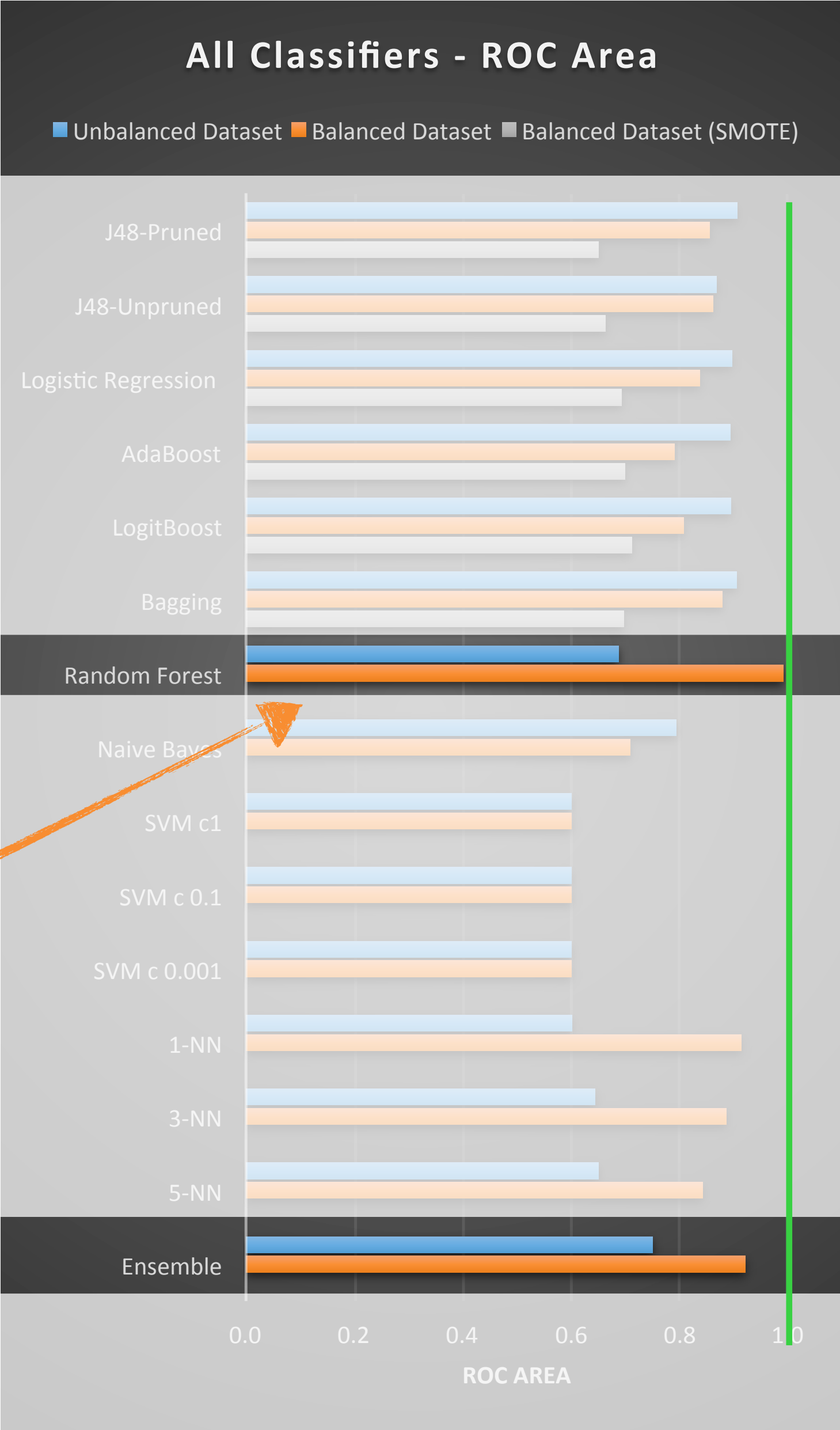
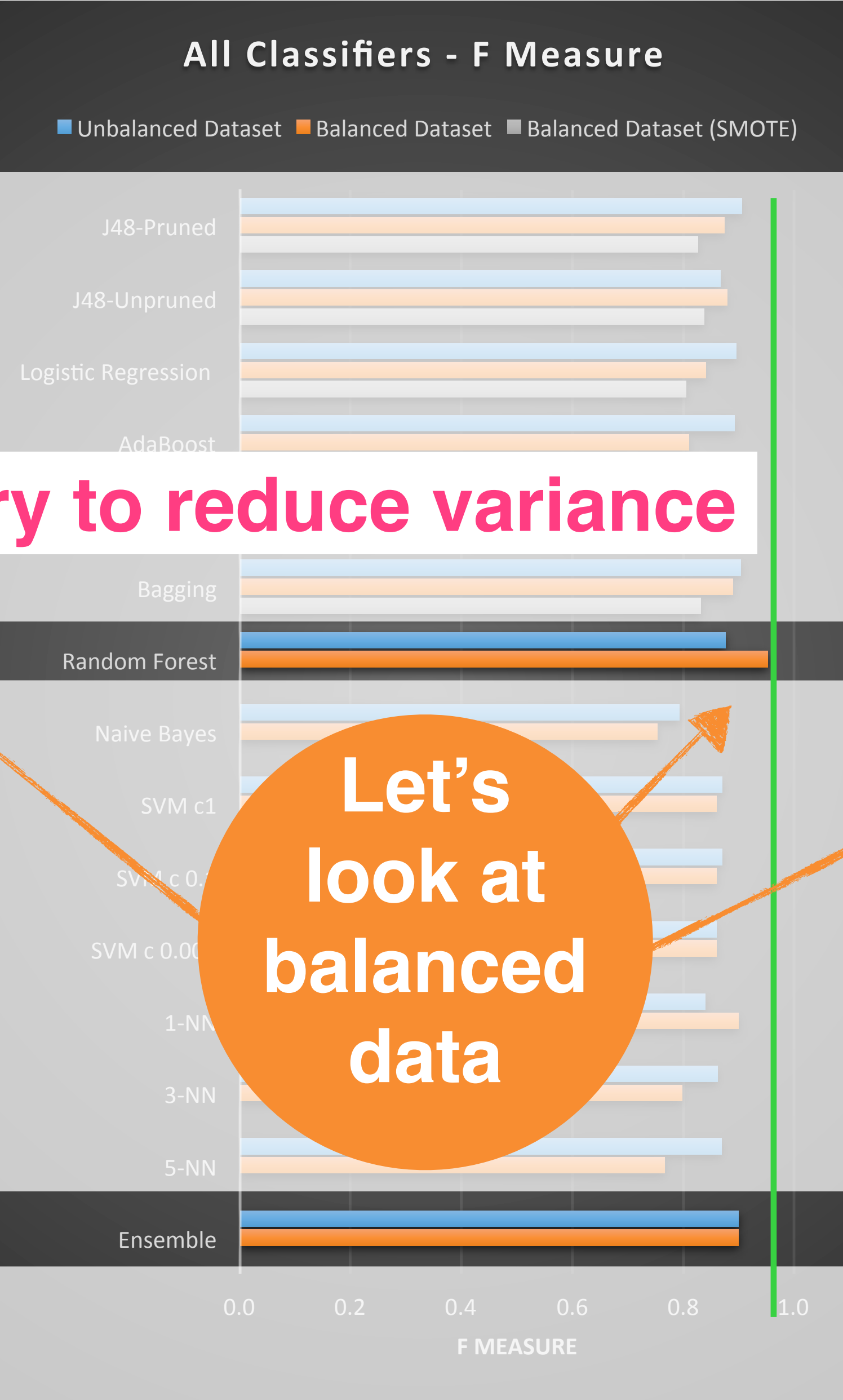
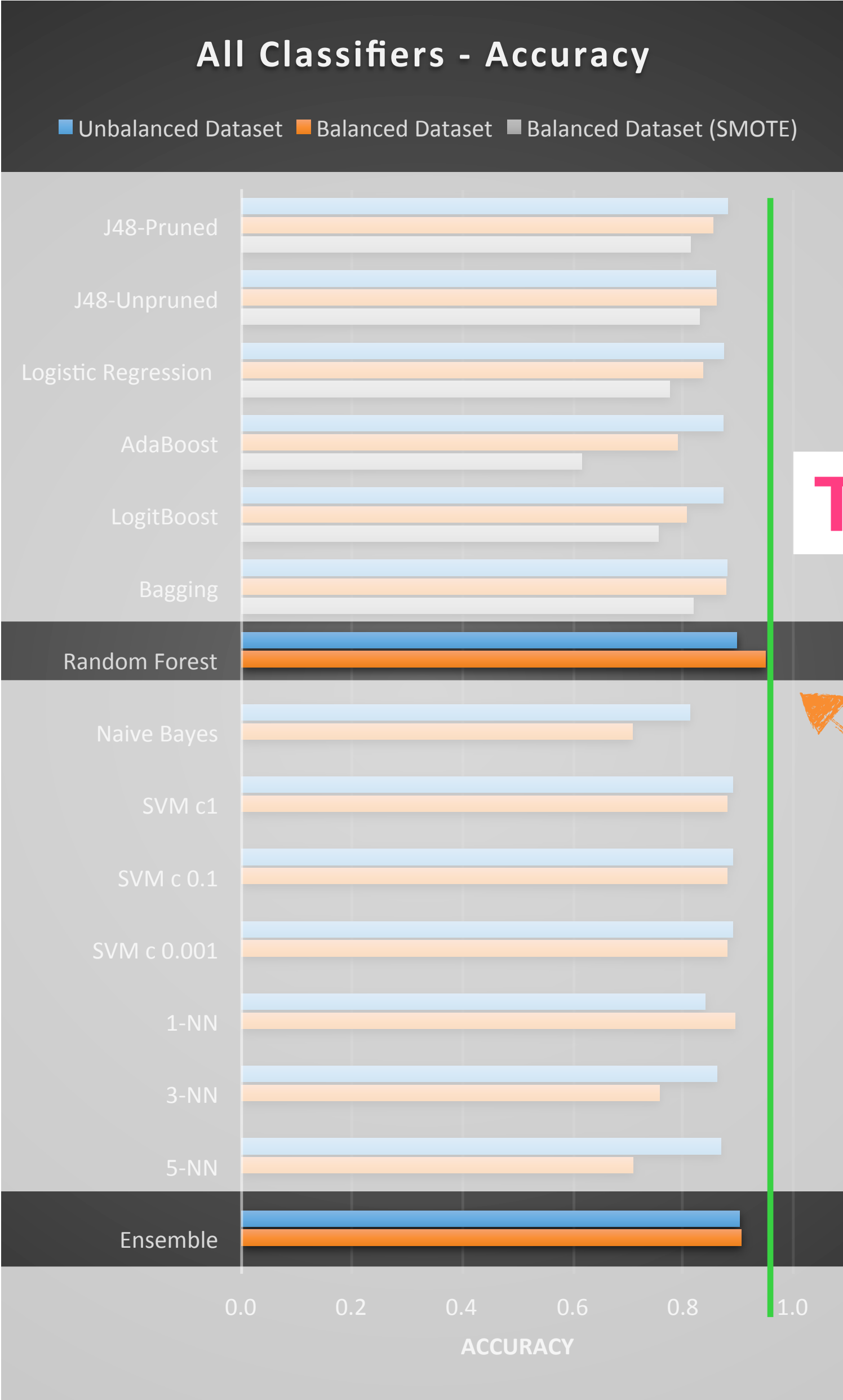




Let's look at balanced data

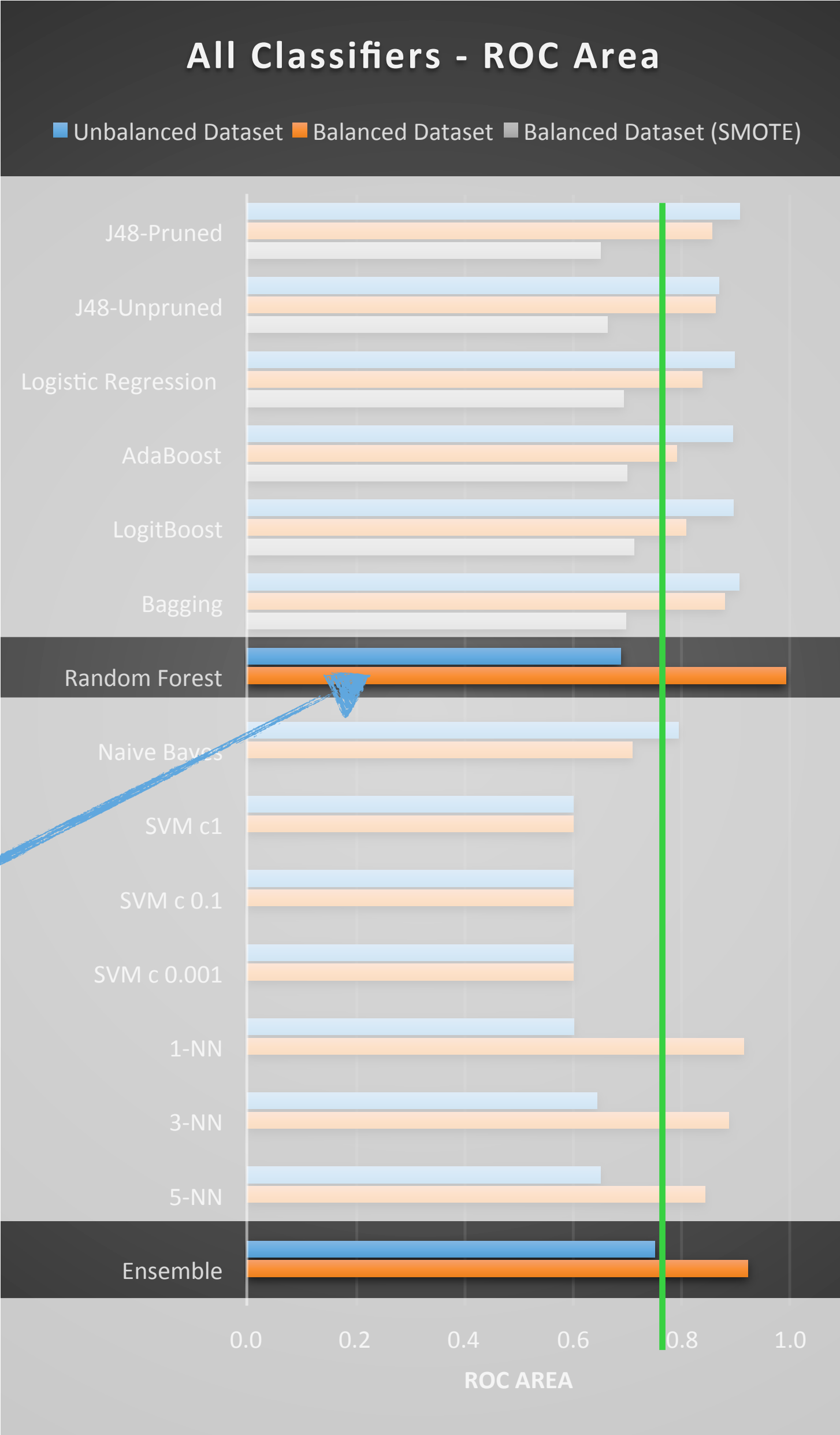
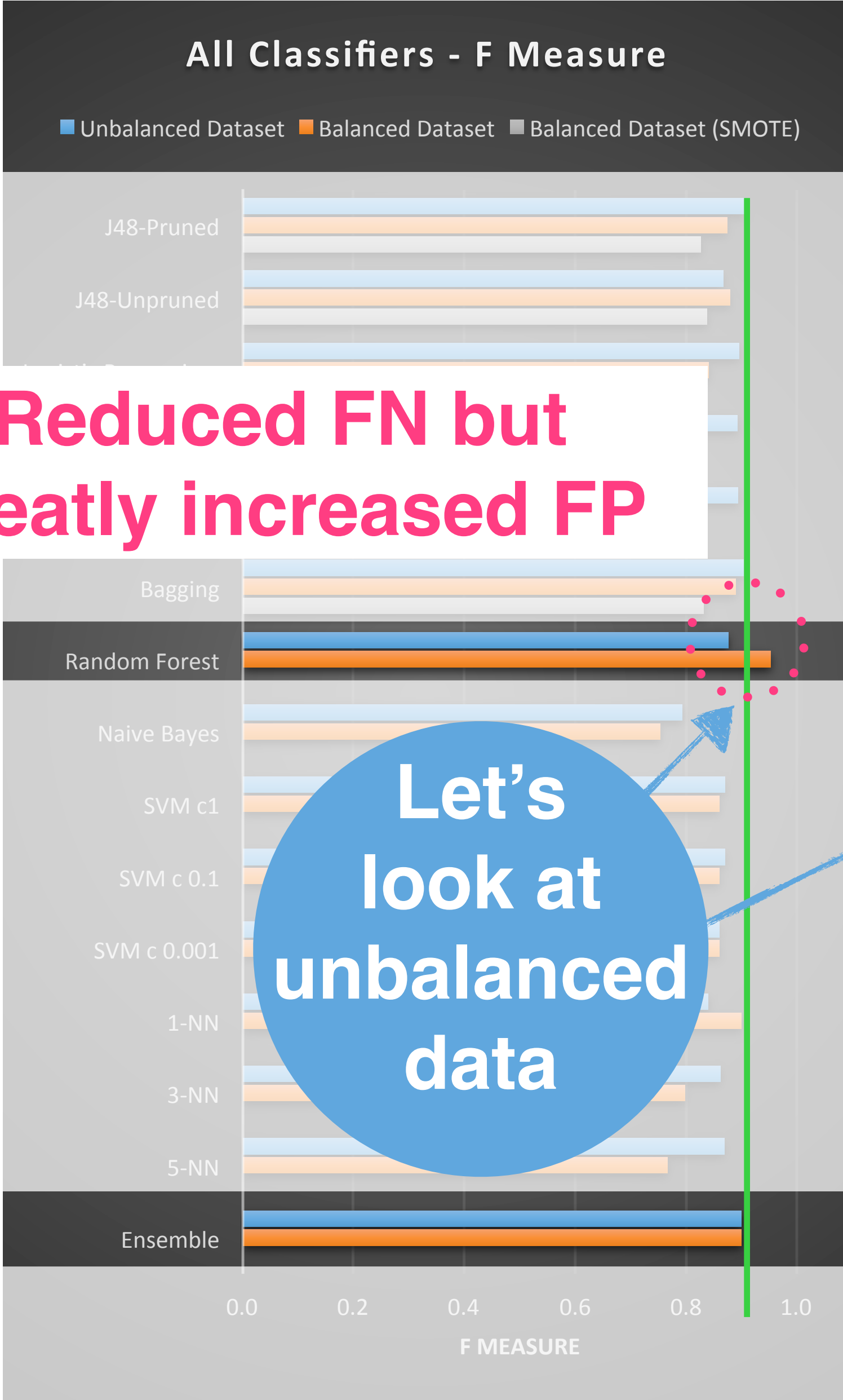
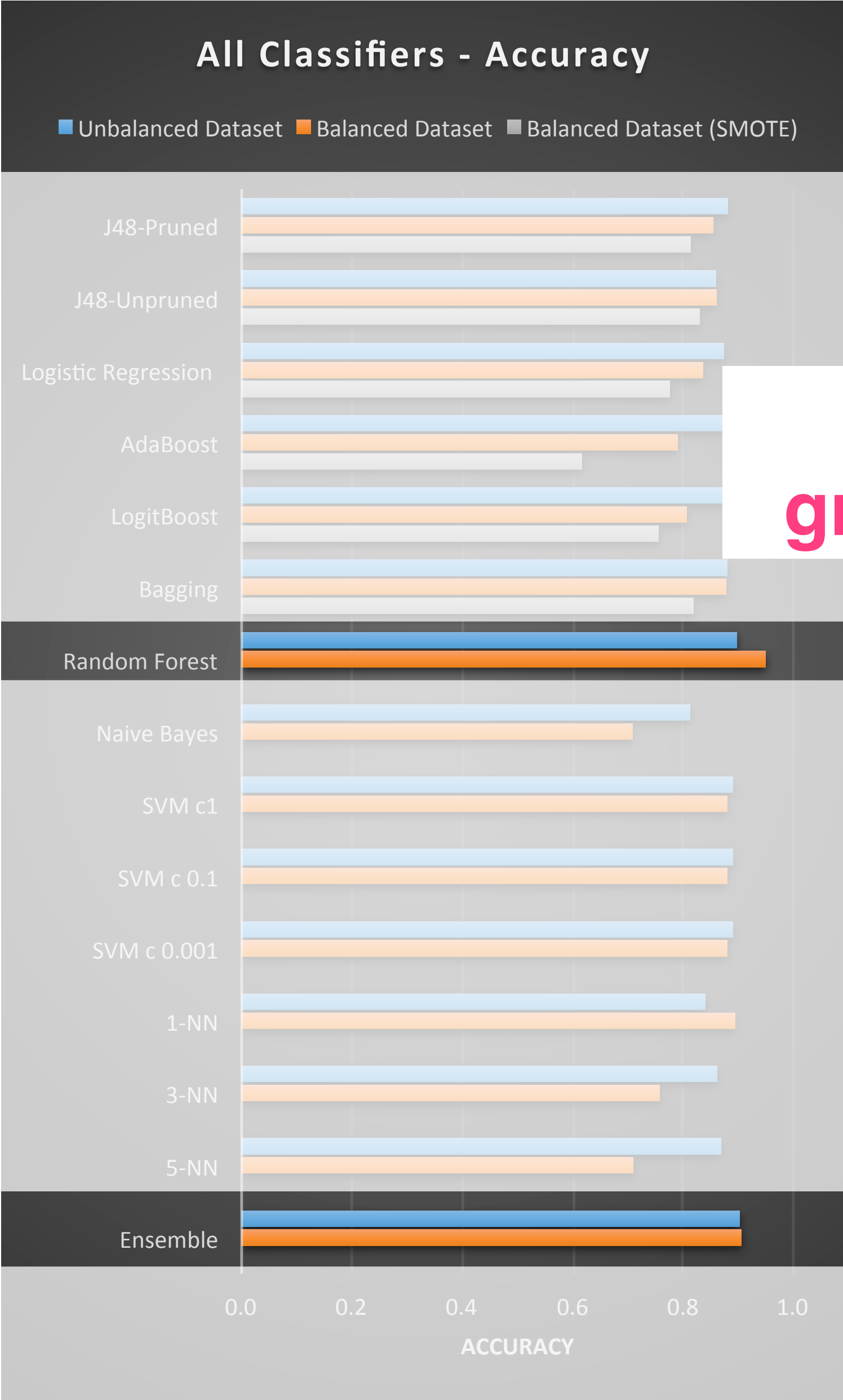
Performs very Well!





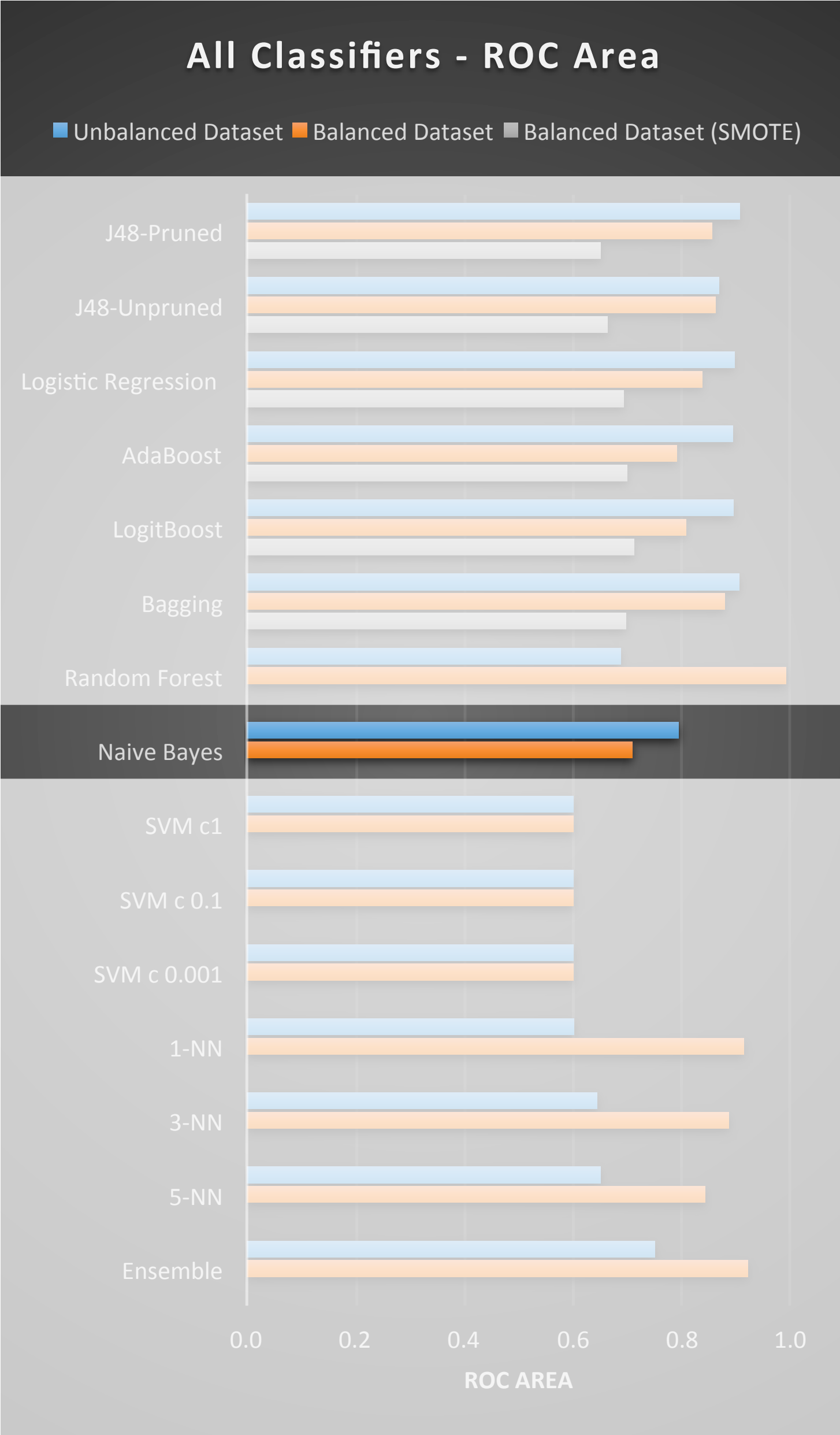
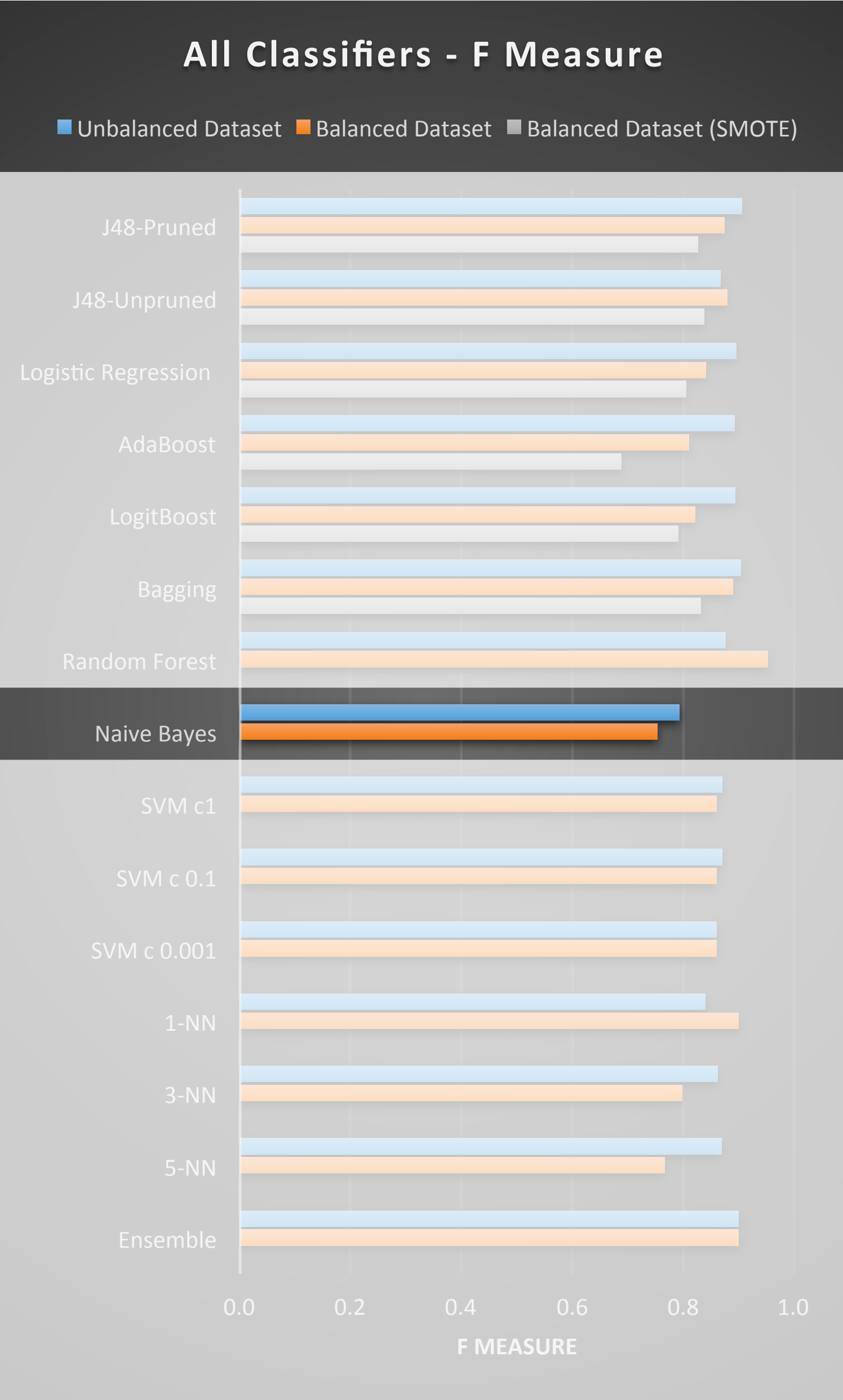
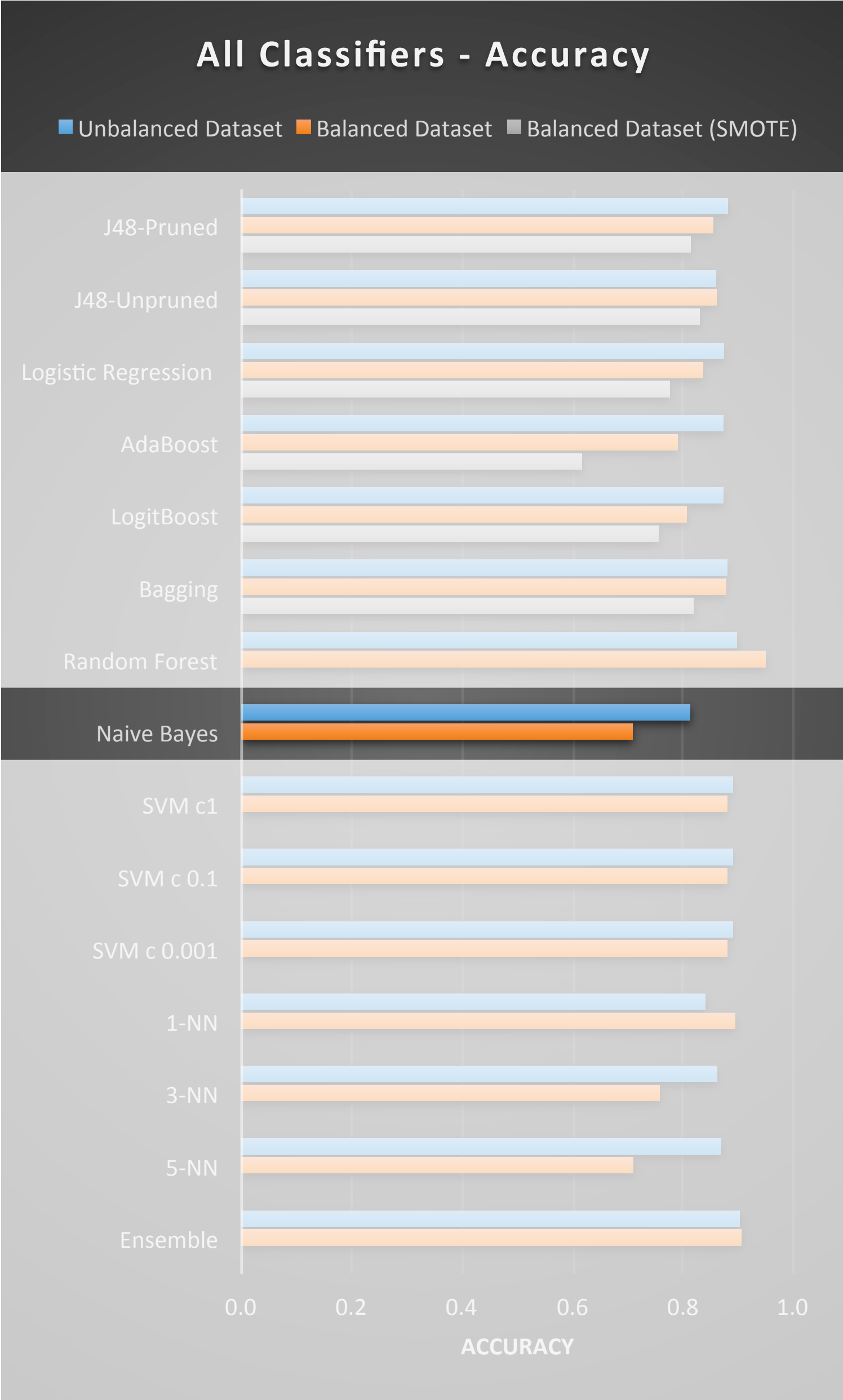
Try to reduce variance

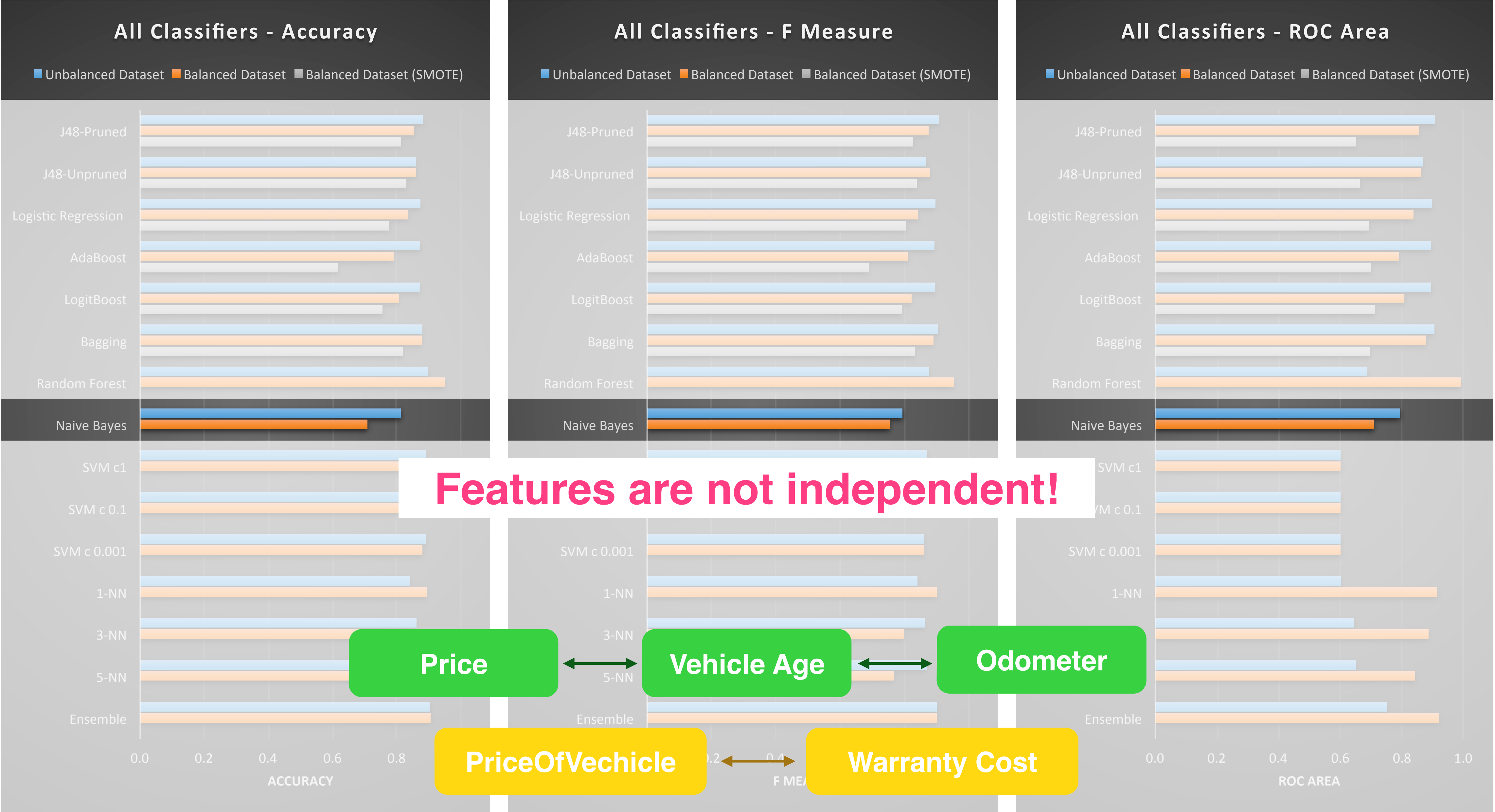
Let's look at balanced data

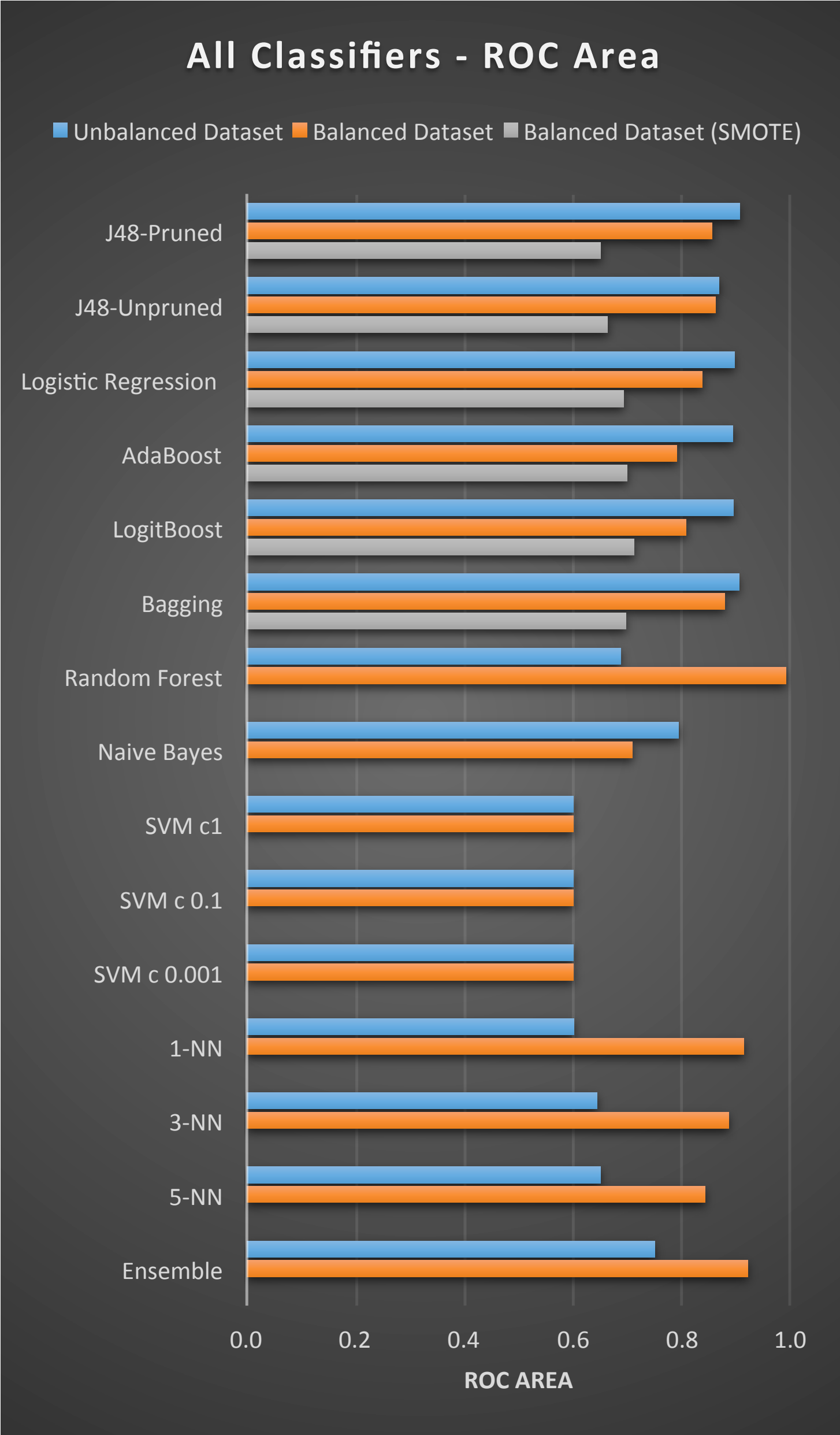
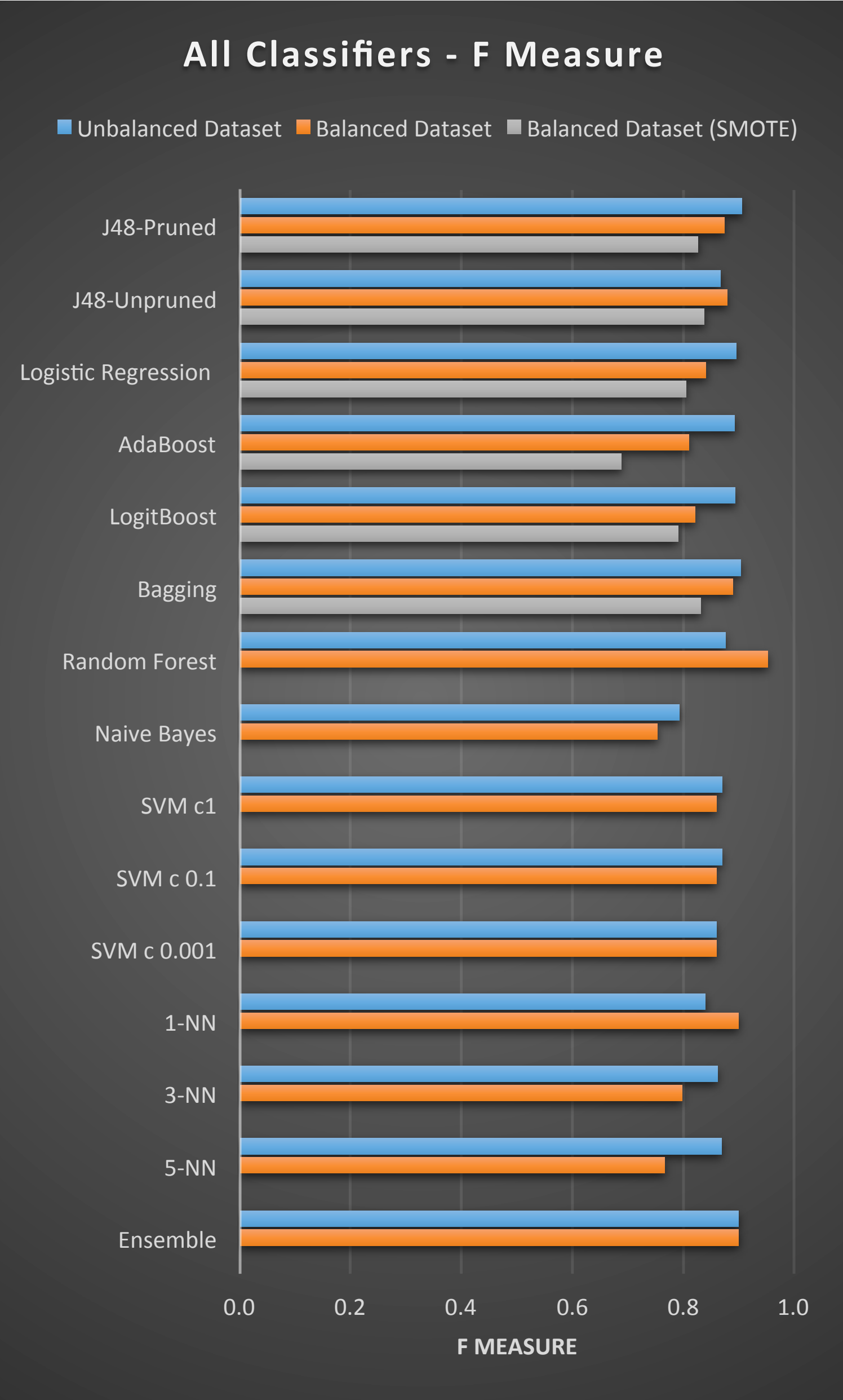
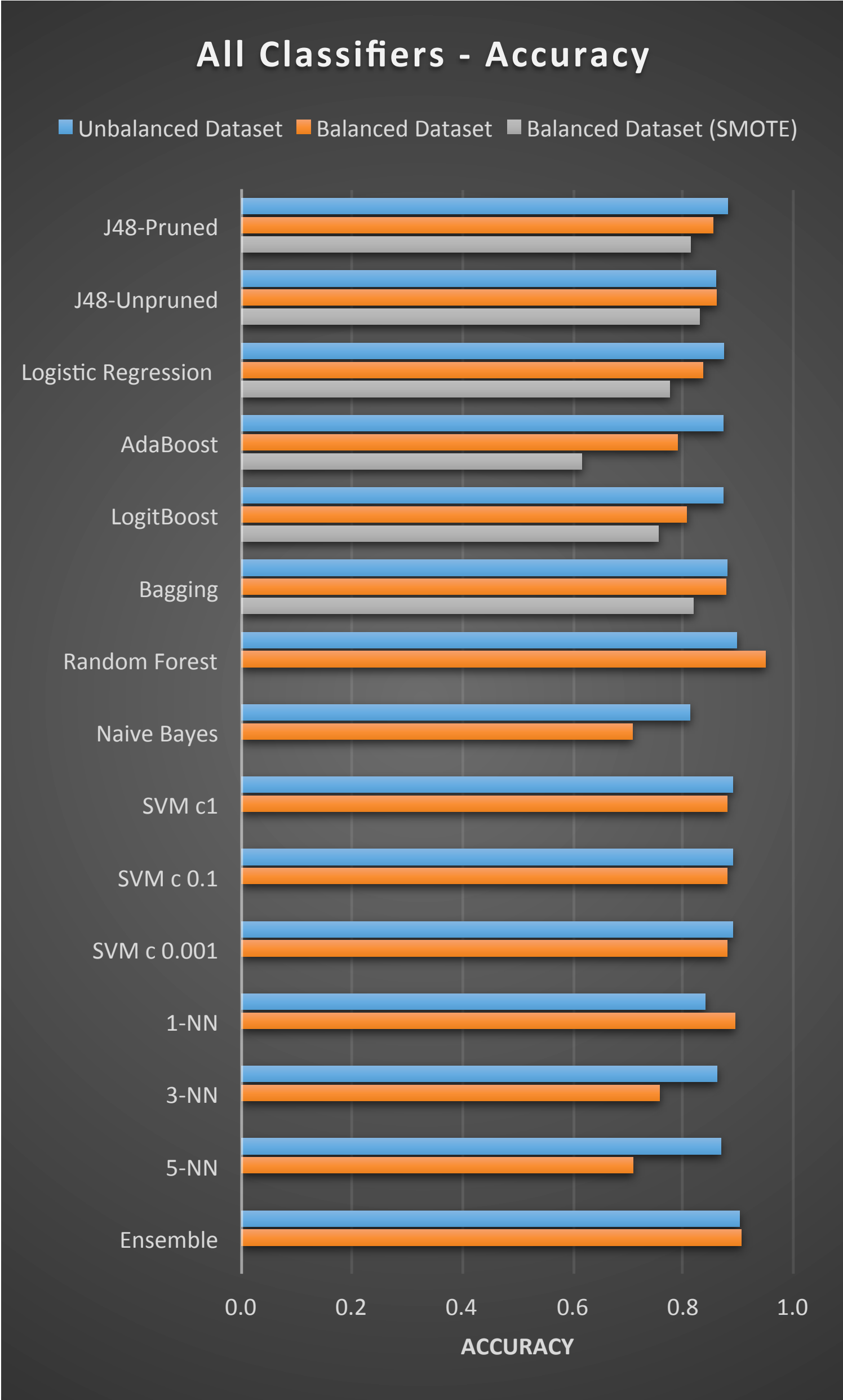


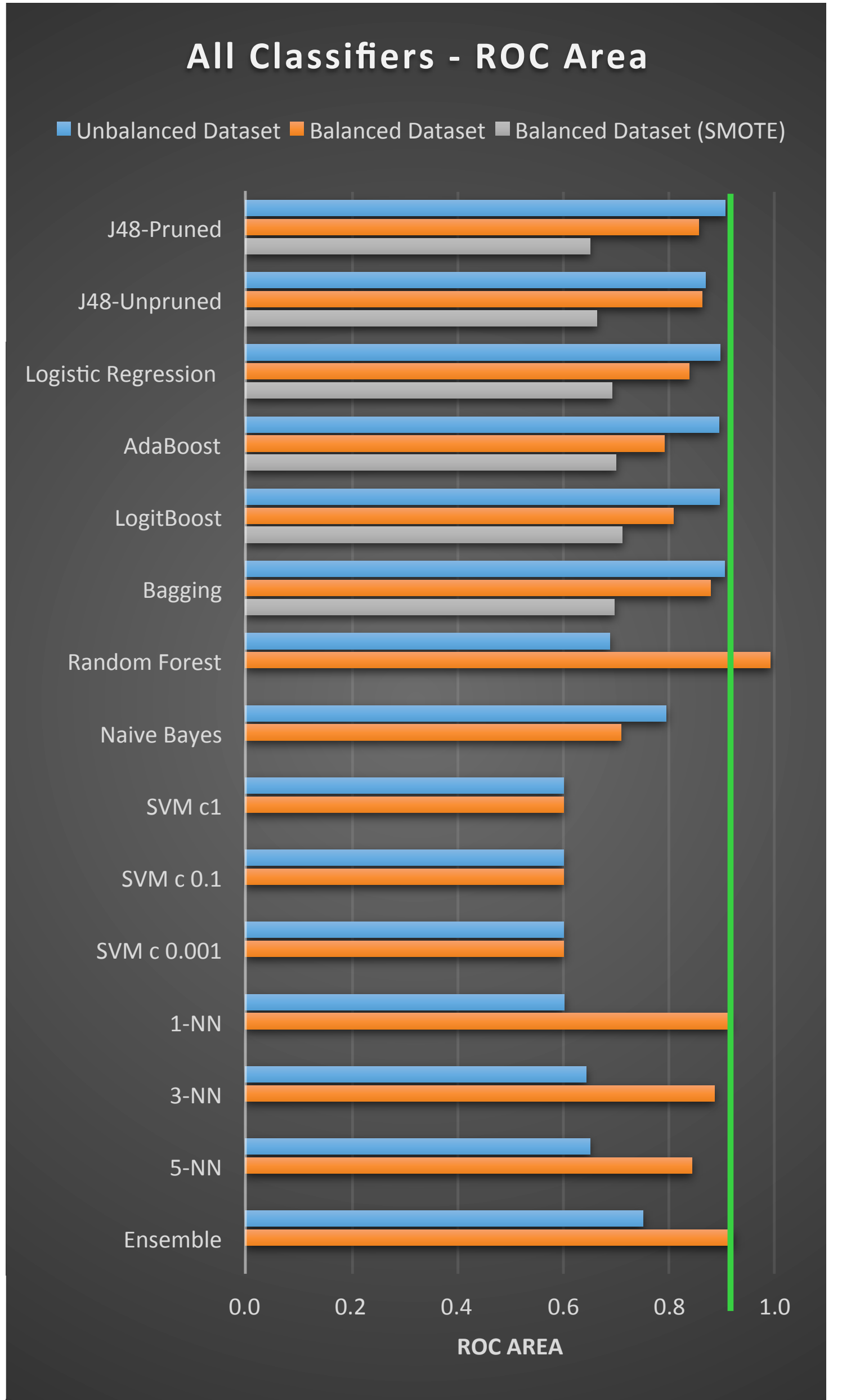
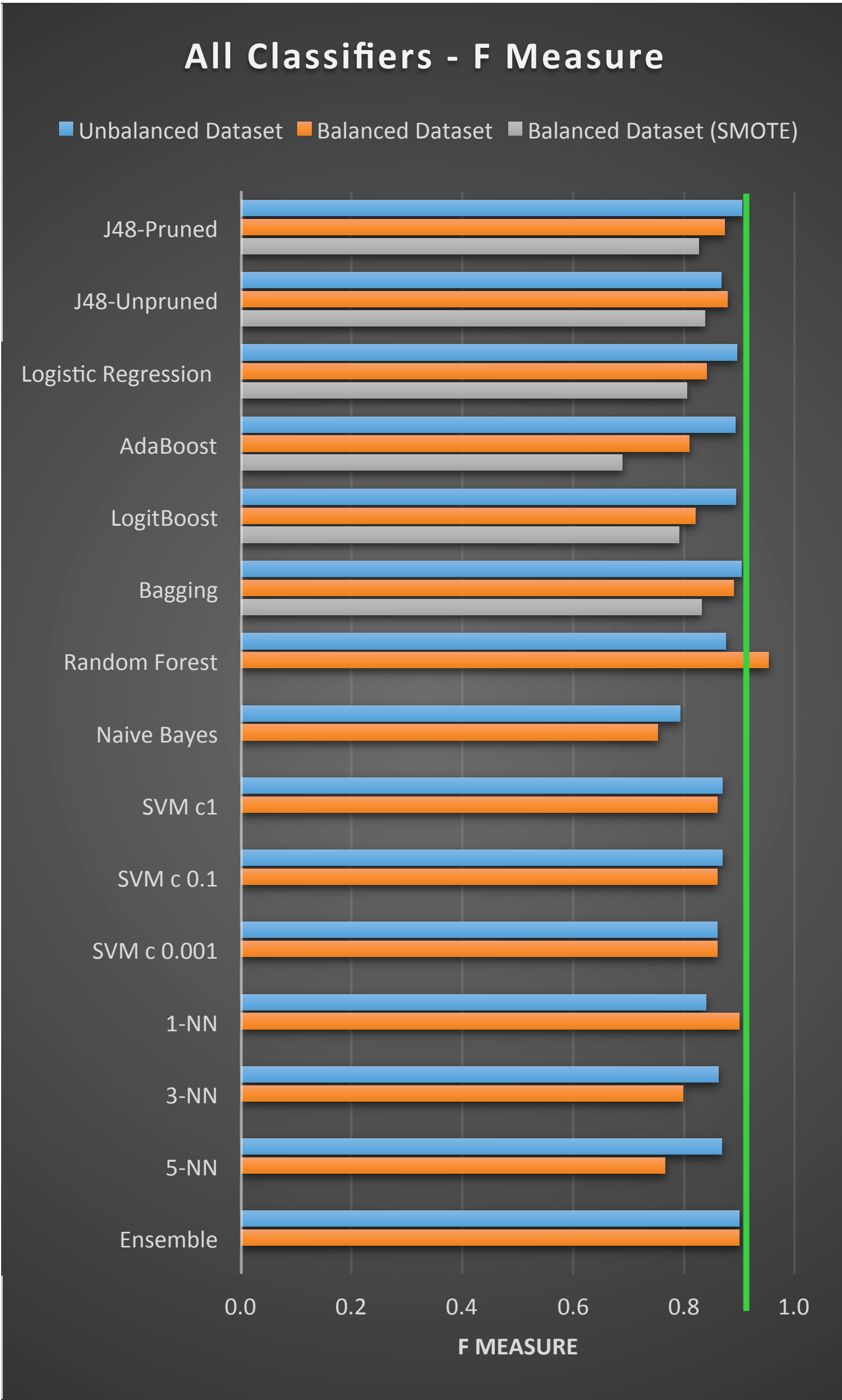
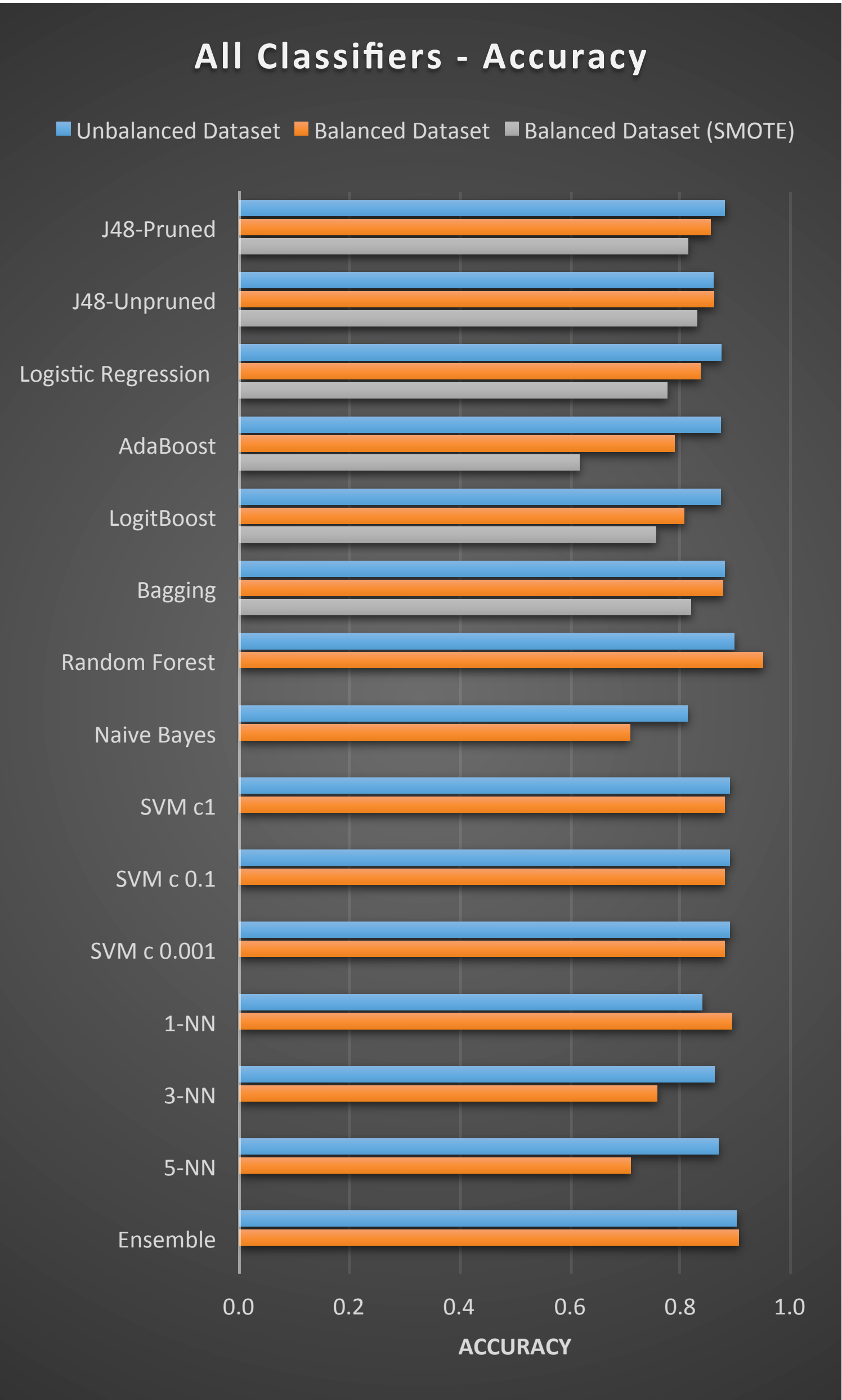
Reduced FN but
greatly increased FP

Let's
look at
unbalanced
data

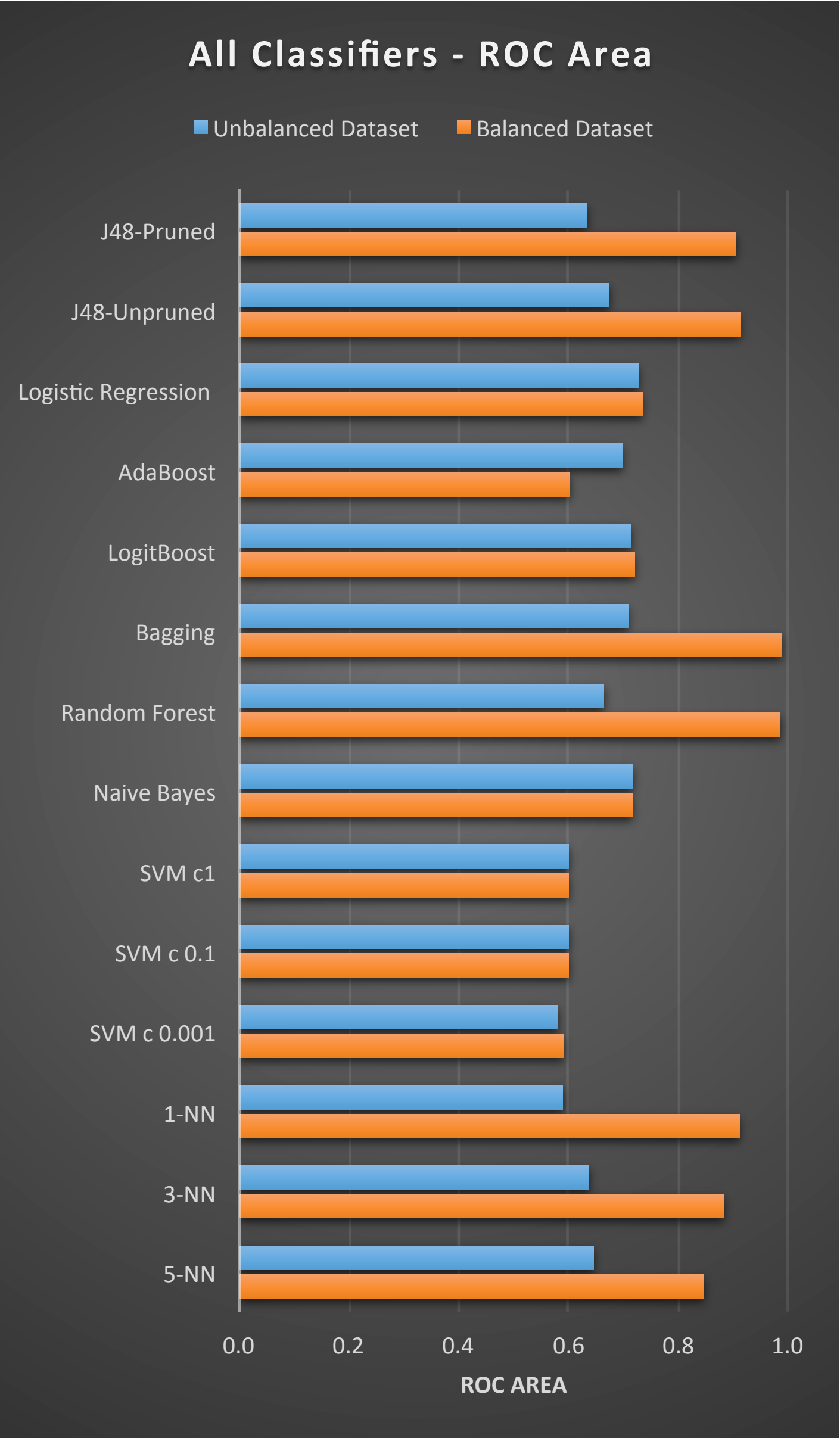
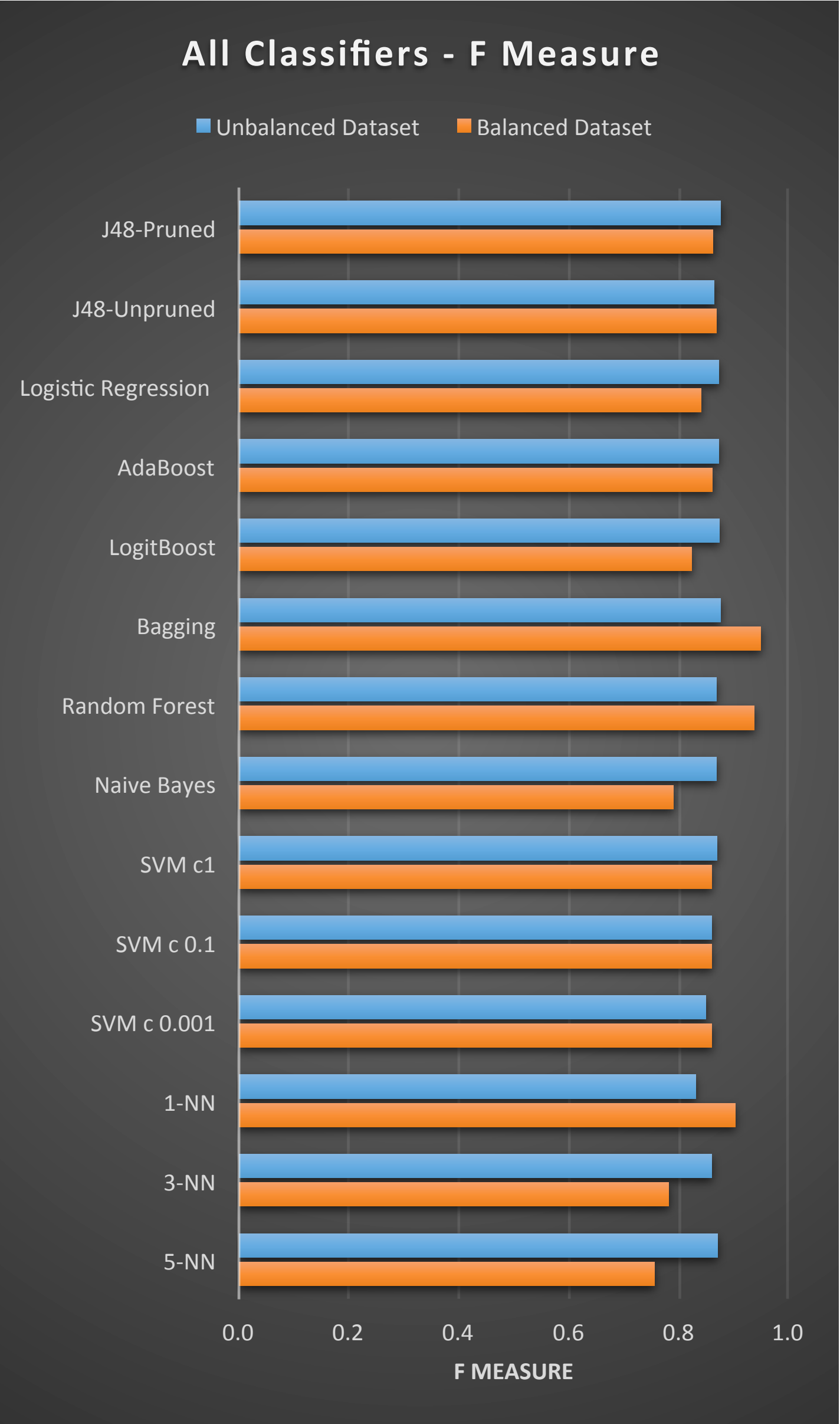
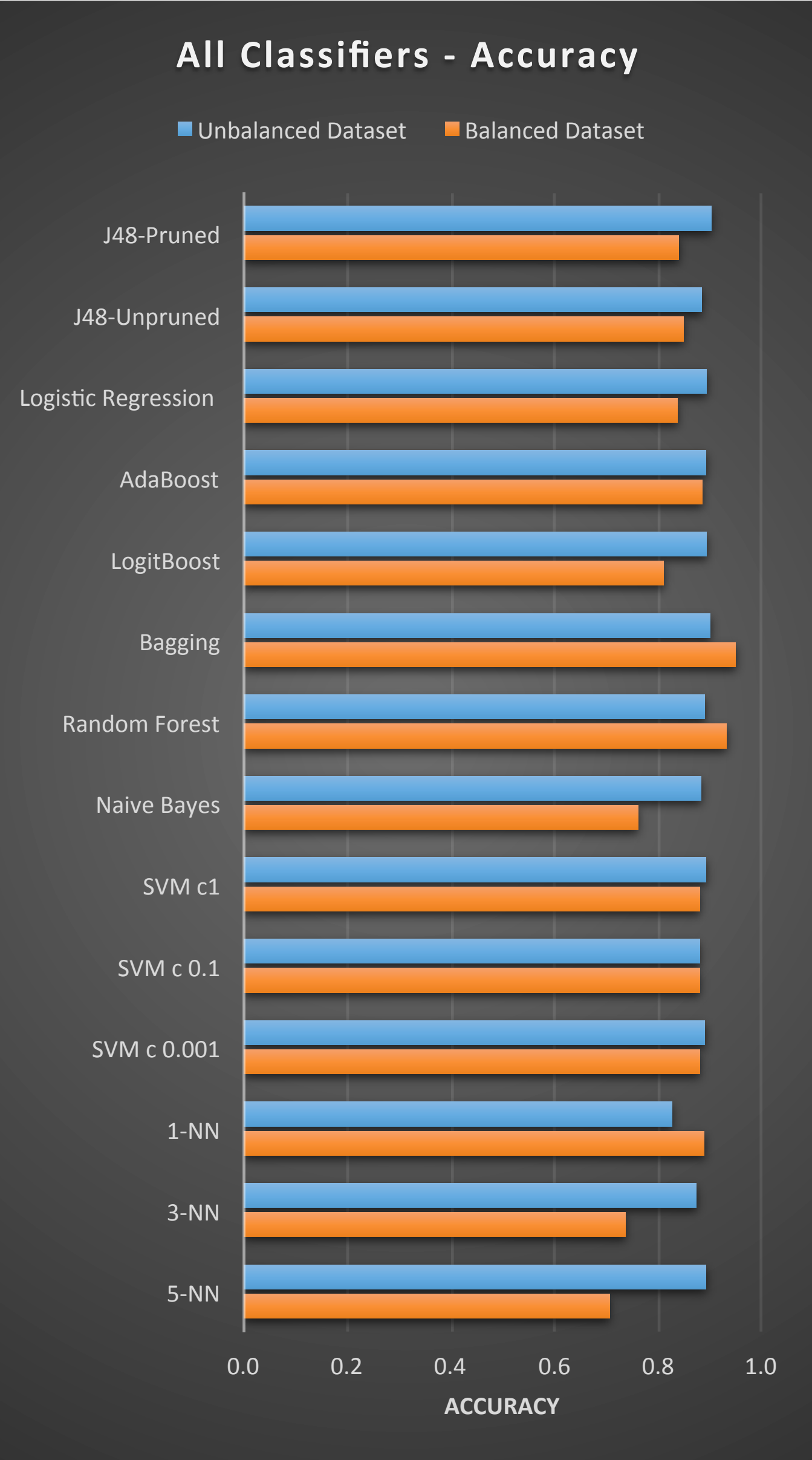








Expert Handpicked Features



Key Learnings

- Important features
 - VehOdo
 - VehAge
 - MMRCurrentAuctionCleanPrice
 - MMRCurrentRetailAveragePrice
 - WarrantyCost
 - Wheel type
- In our dataset classifiers work better on Unbalanced data rather than Balanced data
- Decision tree (J48) and logistic regression perform best
- Oversampling with replacement work better than SMOTE
- Undersampling the majority class did even WORSE