# Class 11 – Tree based methods and Random Forest (Classification)

Pedram Jahangiry

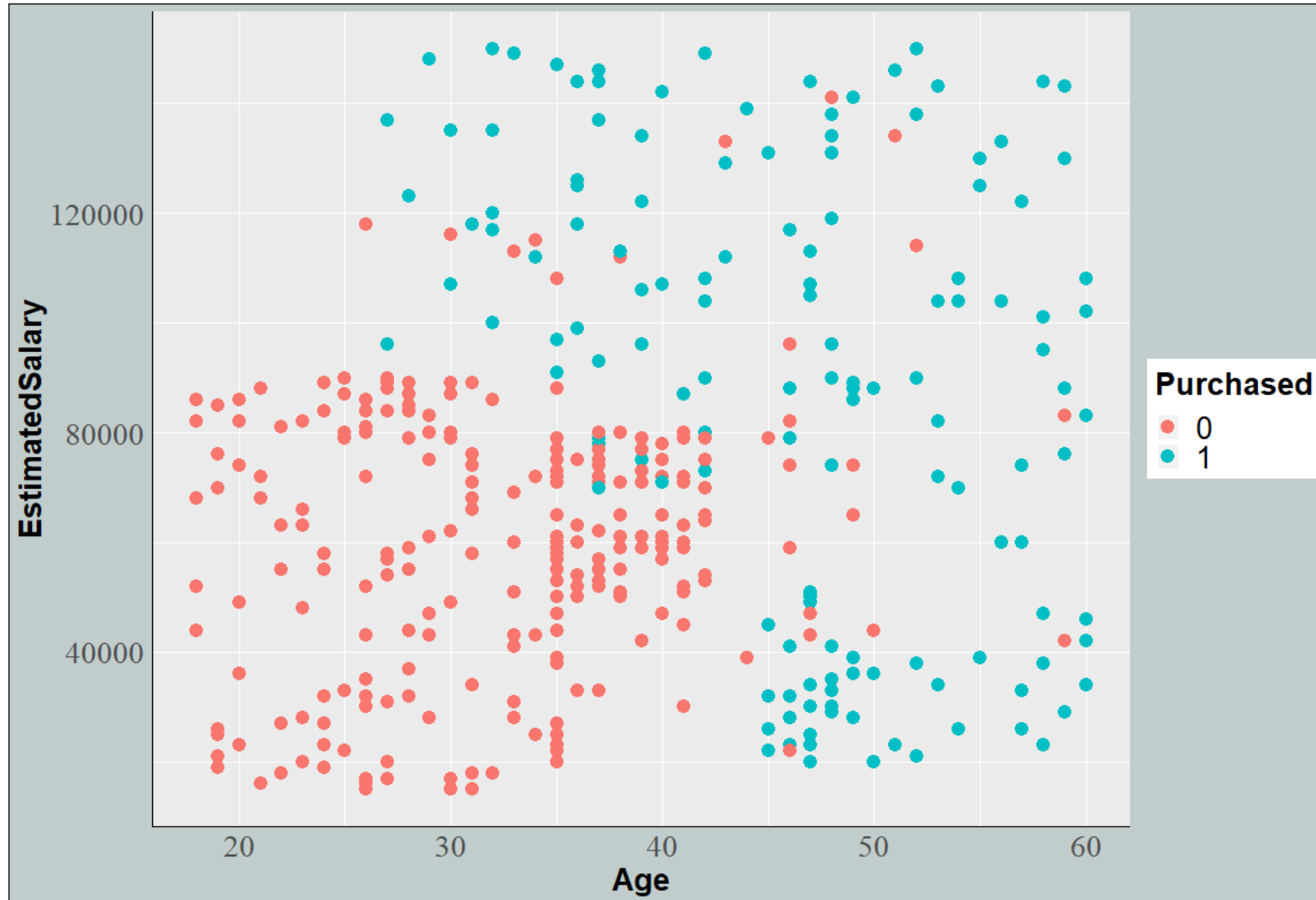Fall 2019

JON M.
HUNTSMAN
SCHOOL OF BUSINESS
**UtahState**University

# Classification Trees

- Very similar to a regression tree, except that it is used to predict a qualitative response rather than a quantitative one.

- For a classification tree, we predict that each observation belongs to the *most commonly occurring class* of training observations in the region to which it belongs.
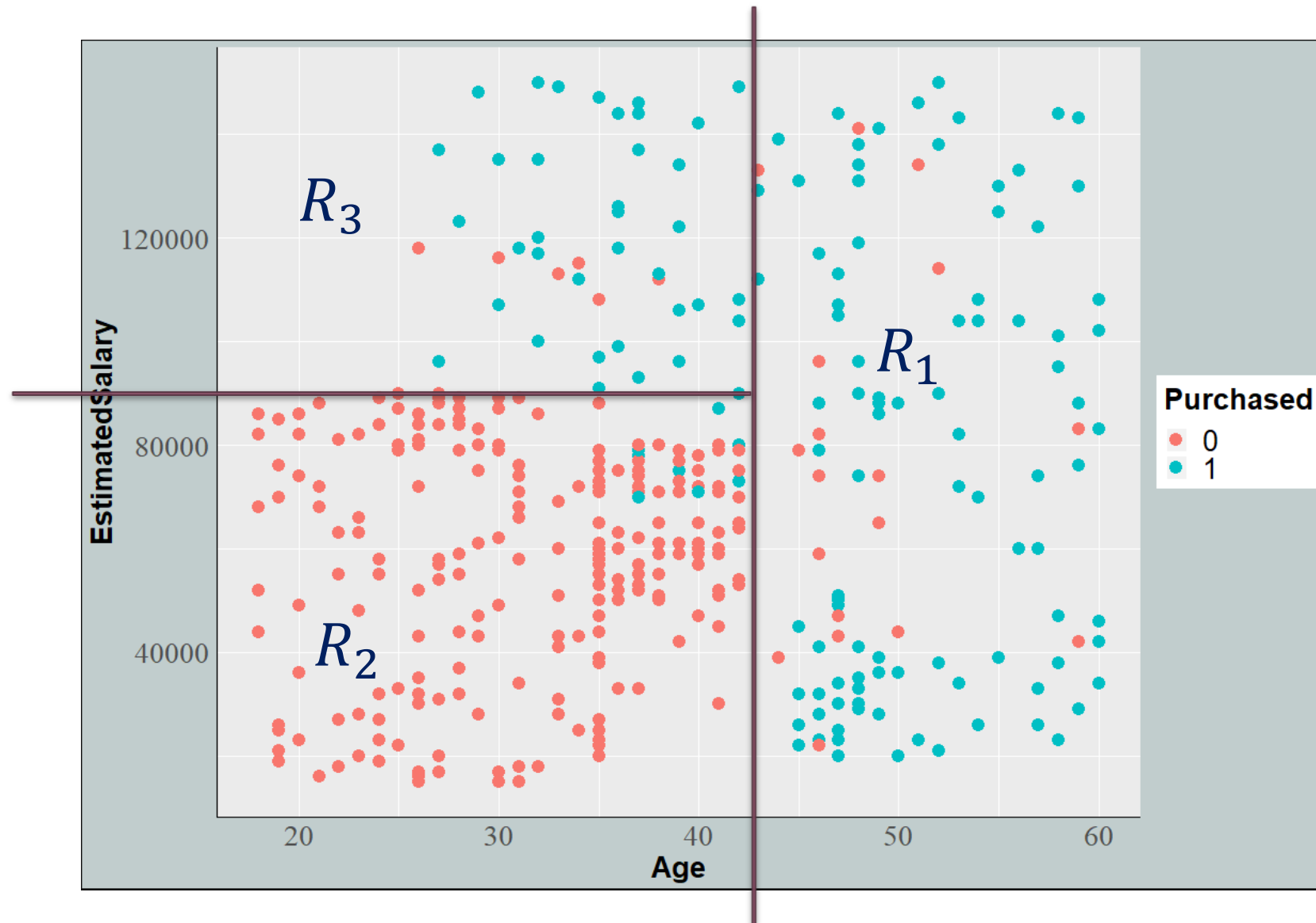
# Tree based classification

# Tree based classification

# Tree based classifiation

# Details of classification Trees

- Just as in the regression setting, we use recursive binary splitting to grow a classification tree.
- In the classification setting, RSS cannot be used as a criterion for making the binary splits
- A natural alternative to RSS is the *classification error rate.* this is simply the fraction of the training observations in that region that do not belong to the most common class:

$$E = 1 - \max_{k}(\hat{p}_{mk}).$$

Here $\hat{p}_{mk}$ represents the proportion of training observations in the $m$th region that are from the $k$th class.

- However classification error is not sufficiently sensitive for tree-growing, and in practice two other measures are preferable.

6

# Gini Index

- The *Gini index* is defined by

$$G = \sum_{k=1}^{K} \hat{p}_{mk}(1 - \hat{p}_{mk}),$$

  a measure of total variance across the $K$ classes. The Gini index takes on a small value if all of the $\hat{p}_{mk}$'s are close to zero or one.

- For this reason the Gini index is referred to as a measure of node *purity* — a small value indicates that a node contains predominantly observations from a single class.

# Cross-Entropy (Deviance)

- An alternative to the Gini index is *cross-entropy*, given by

$$D = -\sum_{k=1}^{K} \hat{p}_{mk} \log \hat{p}_{mk}.$$

- It turns out that the Gini index and the cross-entropy are very similar numerically.

When building a classification tree, either the Gini index or the entropy are typically used to evaluate the quality of a particular split, since these two approaches are more sensitive to node purity than is the classification error rate.

Any of these three approaches might be used when *pruning* the tree, but the classification error rate is preferable if prediction accuracy of the **final pruned tree** is the goal.
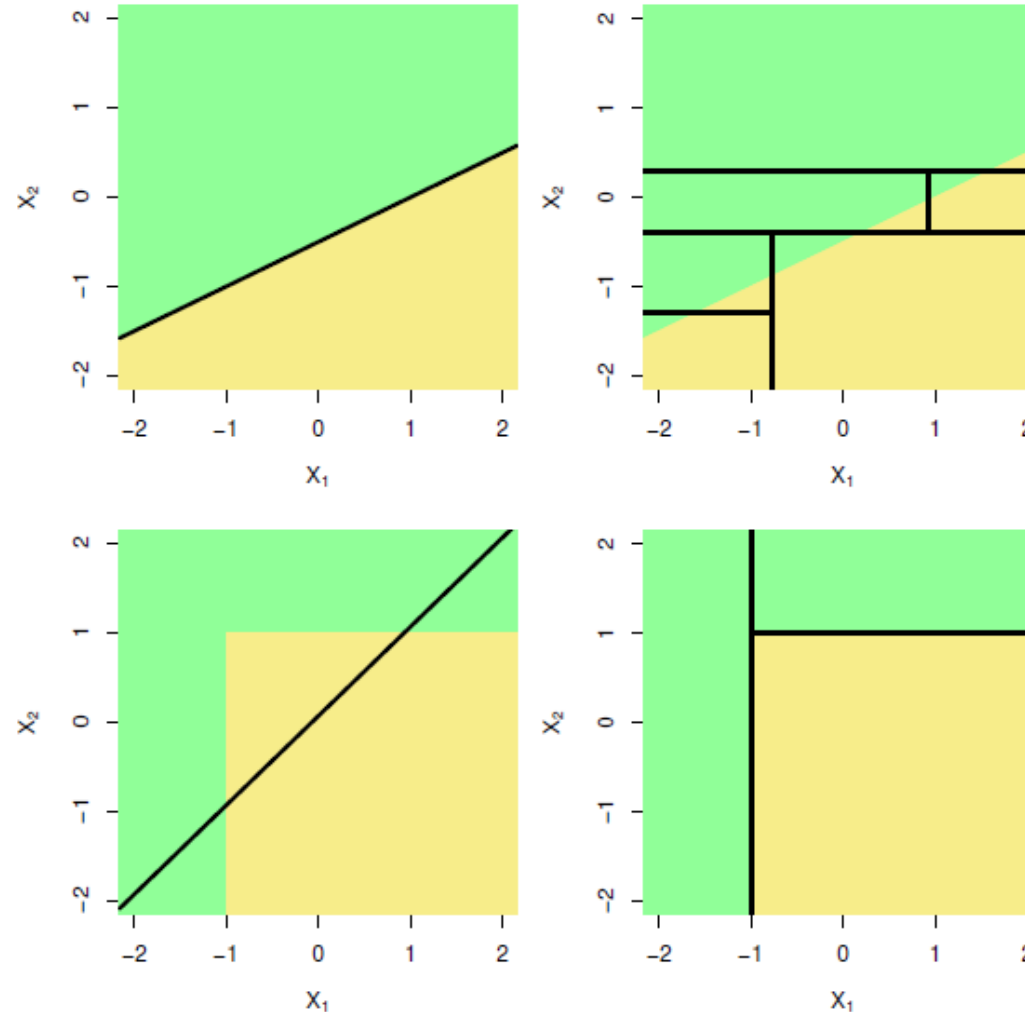
# Bagging Classification Trees

- For classification trees: for each test observation, we record the class predicted by each of the $B$ trees, and take a *majority vote*: the overall prediction is the most commonly occurring class among the $B$ predictions.
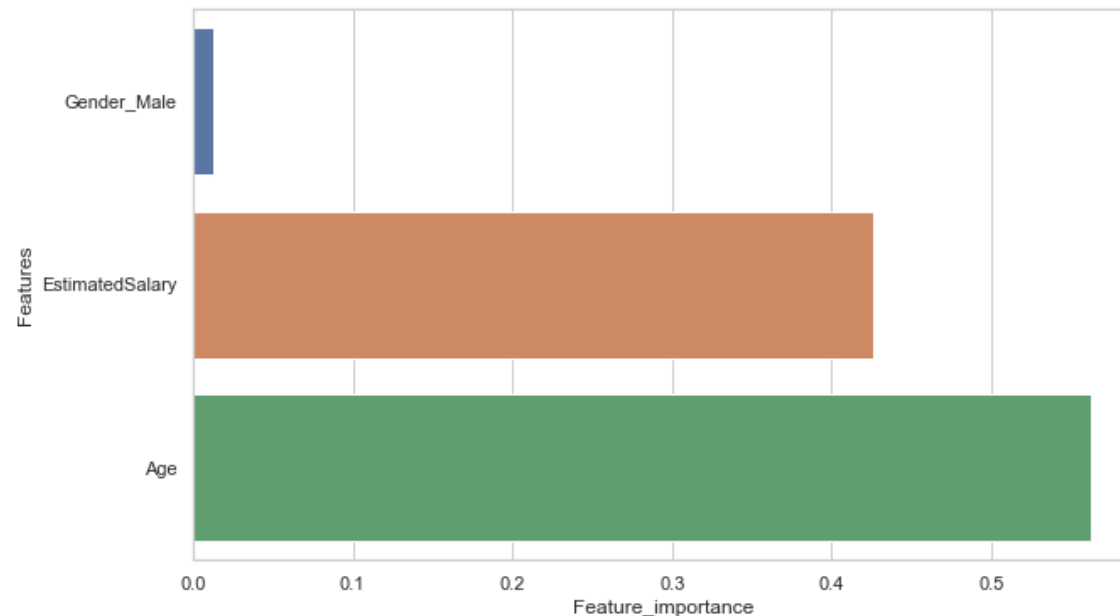
9

# Trees Versus Linear Models

Left column: linear model; Right column: tree-based model

Top Row: True linear boundary
Bottom row: true non-linear boundary.

# Variable importance measure

- For bagged/RF regression trees, we record the total amount that the RSS is decreased due to splits over a given predictor, averaged over all $B$ trees. A large value indicates an important predictor.
- Similarly, for bagged/RF classification trees, we add up the total amount that the Gini index is decreased by splits over a given predictor, averaged over all $B$ trees.

# Random Forest in python

- Find the Random forests Sklearn documentation [here](#)

- Blackbox version of Random Forests (Classification) in python:

```python
# Fitting RF classifier to the Training set

RF_classifier = RandomForestClassifier(n_estimators = 100, criterion='gini')
RF_classifier.fit(X_train, y_train)
```

```
RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                       max_depth=None, max_features='auto', max_leaf_nodes=None,
                       min_impurity_decrease=0.0, min_impurity_split=None,
                       min_samples_leaf=1, min_samples_split=2,
                       min_weight_fraction_leaf=0.0, n_estimators=100,
                       n_jobs=None, oob_score=False, random_state=None,
                       verbose=0, warm_start=False)
```