



CS 4104

APPLIED MACHINE LEARNING

Dr. Hashim Yasin

**National University of Computer
and Emerging Sciences,
Faisalabad, Pakistan.**

BOOSTING & BAGGING



BOOSTING & BAGGING



Boosting & Bagging

4

- Two strategies have been developed for producing optimal trees

Boosting:

- develop *new classification trees* based on the results of previous classification trees

Bagging:

- uses *subsets of the training data* to develop new classification trees

BAGGING

Bagging

6

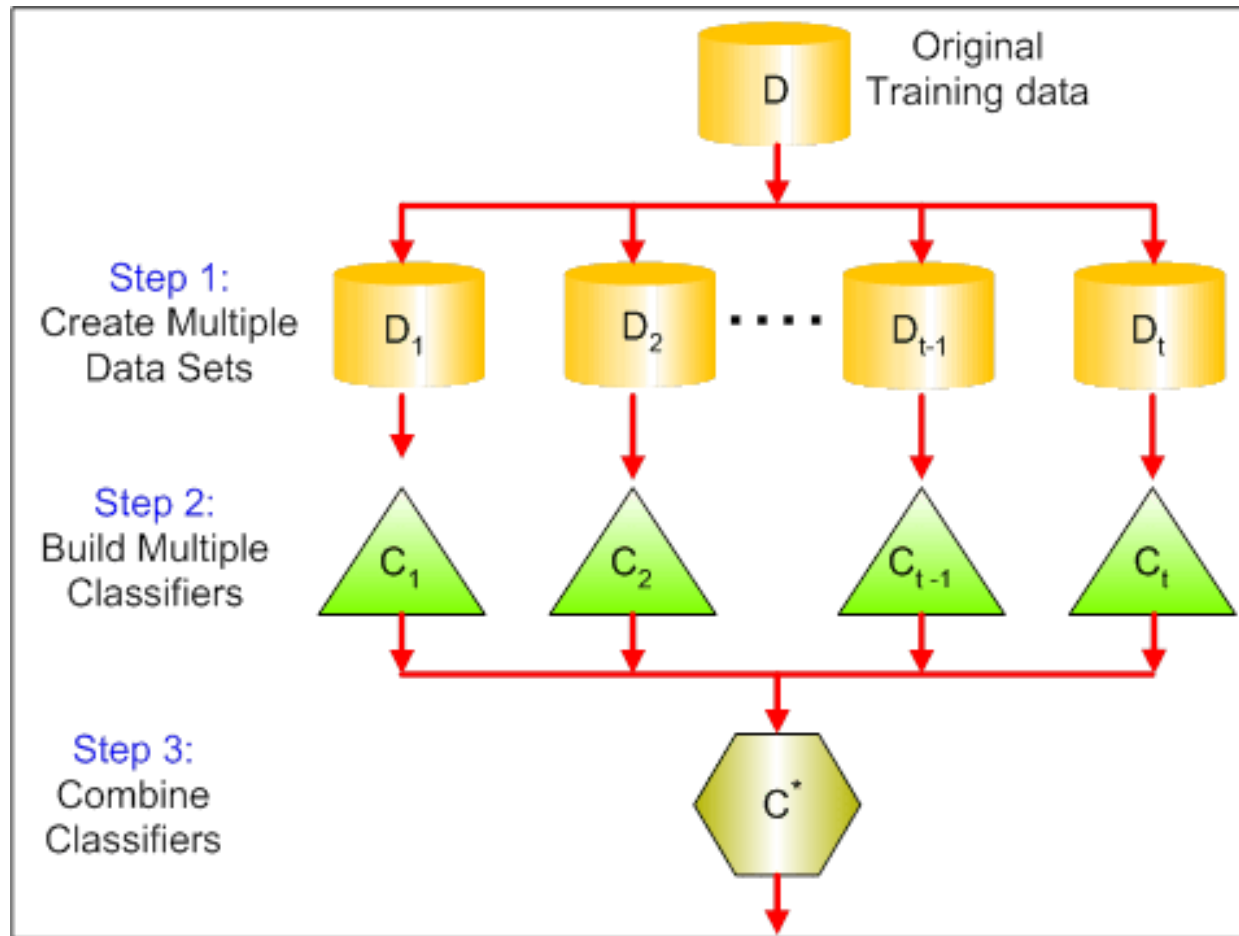
- If *we split the data in random different ways*, decision trees give different results, results into **high variance**.
- We want to **reduce the variance** of a decision tree.

Bagging:

- **Bootstrap aggregating** is a method that result in *low variance*.

Bagging

7



Bagging

8

- For each sample b , we calculate $f^b(x)$, then:

$$\hat{f}_{avg}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^b(x)$$

The prediction at input x when bootstrap sample b is used for training

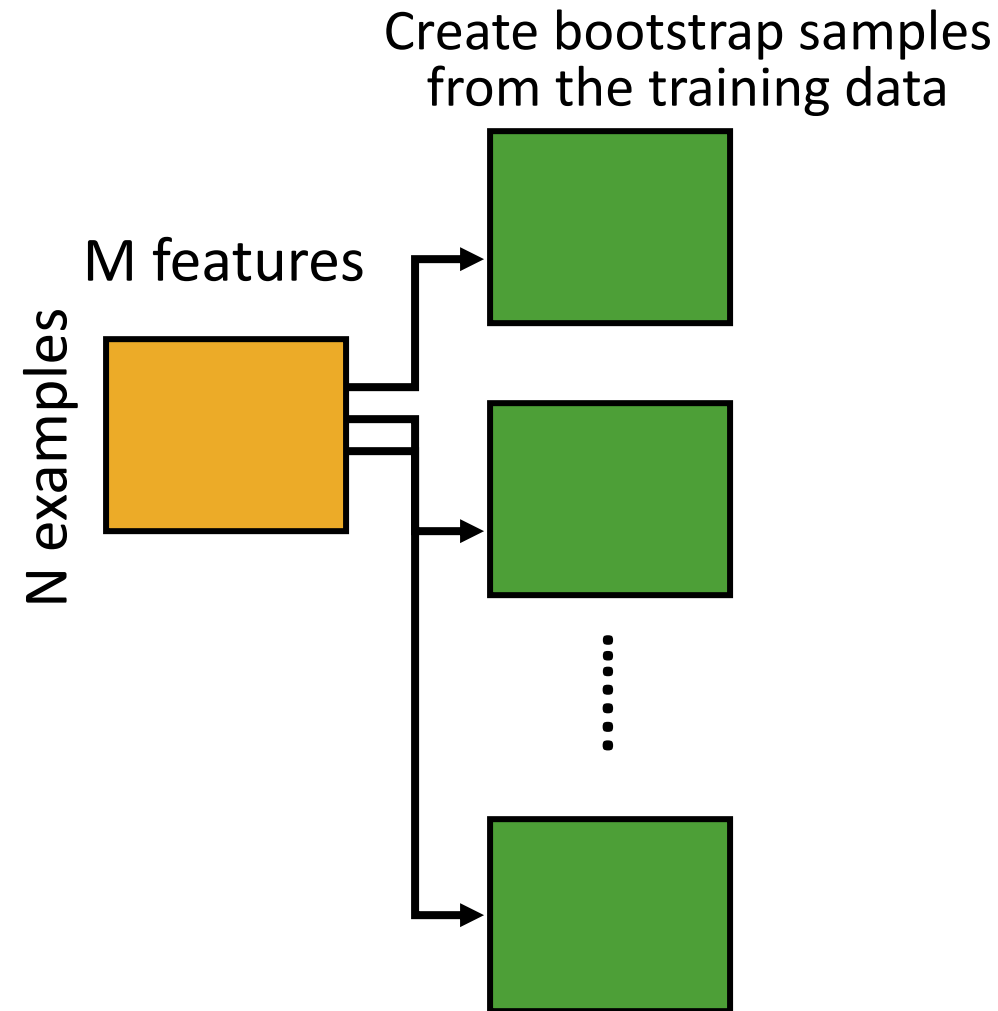
How?

Bootstrap

- Construct B (hundreds) of trees (no pruning)
- Learn a classifier for each bootstrap sample and average them
- Very effective method

Bagging

9



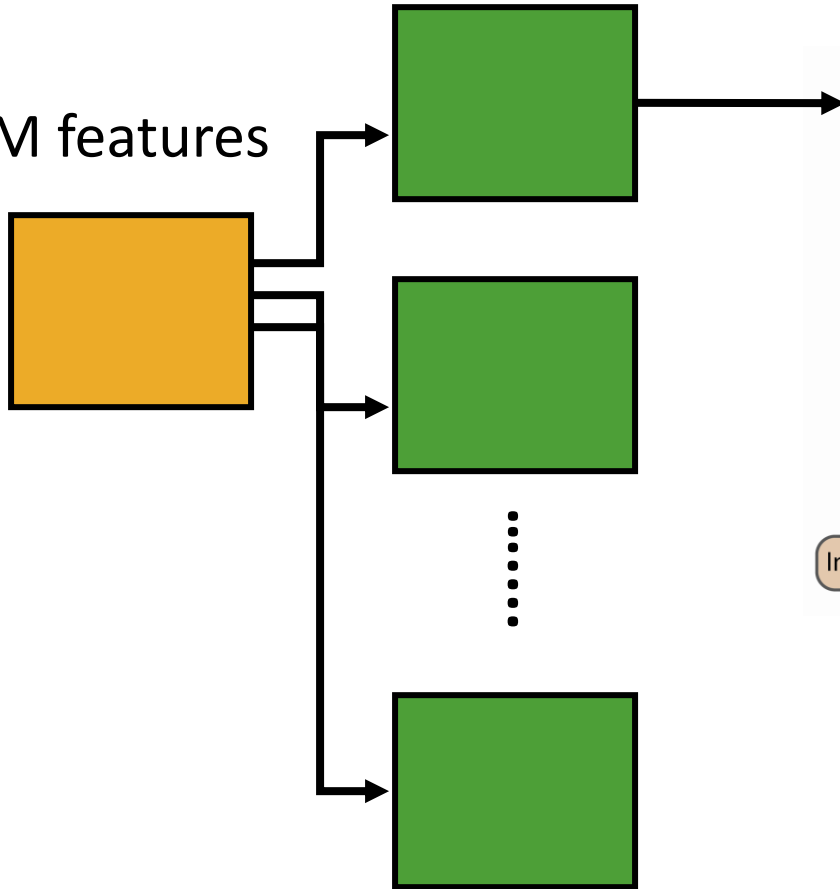
Bagging

10

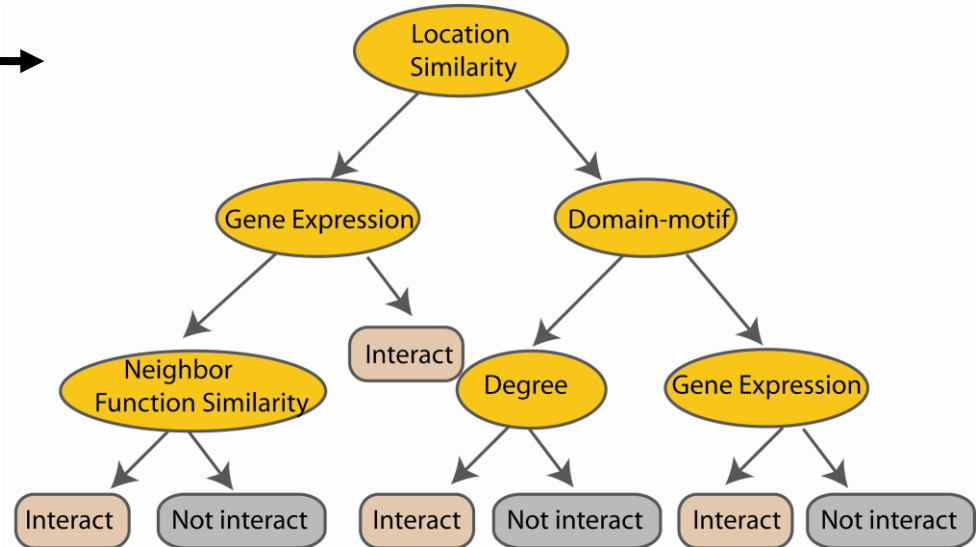
Create bootstrap samples from the training data

M features

N examples



Construct a decision tree



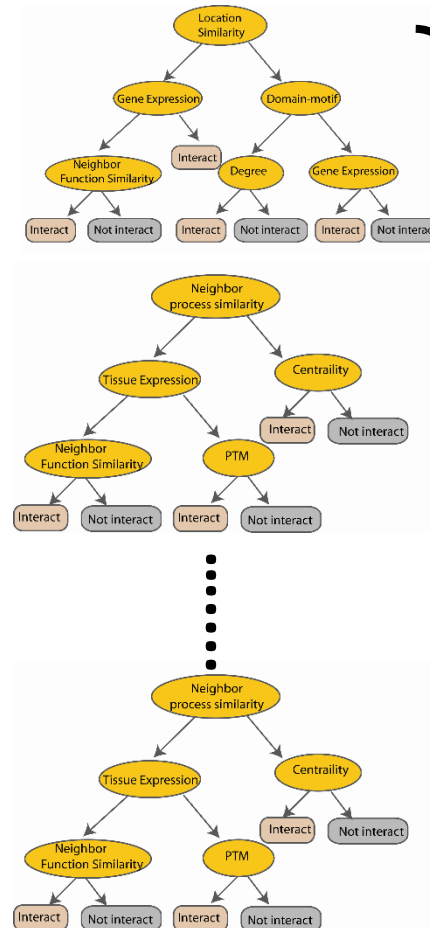
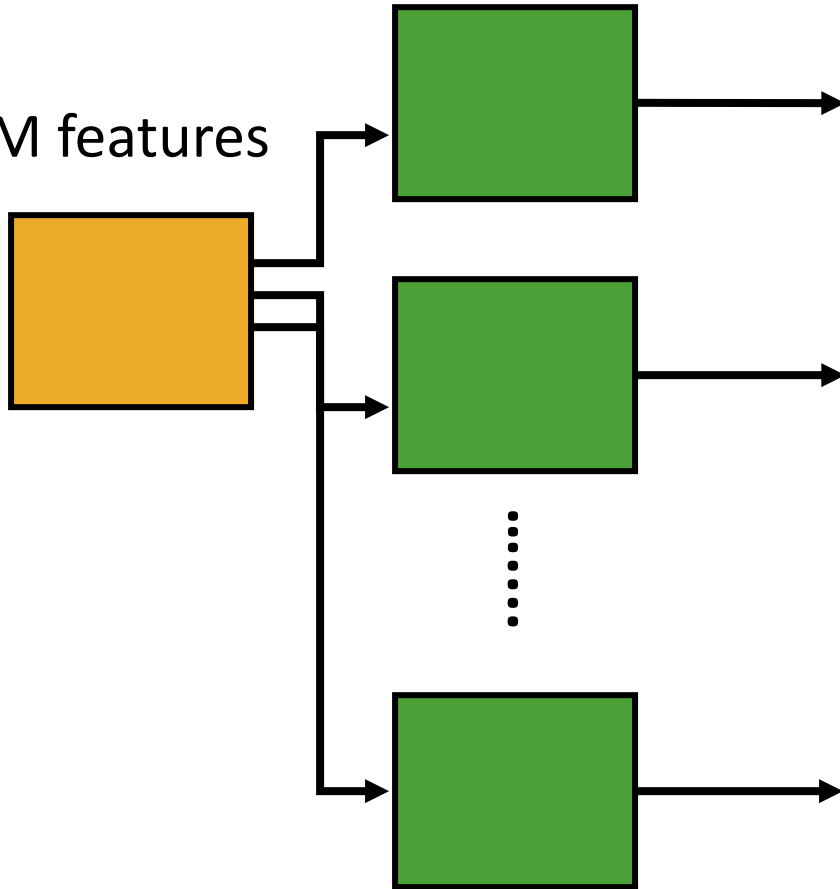
Bagging

11

Create bootstrap samples from the training data

M features

N examples



Take the majority vote

Bagging

12

- ❑ Reduces overfitting (variance)
- ❑ Normally uses **one type of classifier**
 - ❑ Decision trees are popular
- ❑ Easy to parallelize
- ❑ Bagging results in **improved accuracy** over prediction using a single tree
- ❑ Unfortunately, **difficult to interpret the resulting model.**
 - ❑ Bagging improves prediction accuracy at the expense of interpretability.

Bagging ... Issues

13

- Each tree is **identically distributed**
- the **expectation of the average of B** such trees is the same as the expectation of any one of them
- the **bias of bagged trees** is the same as that of the individual trees

Bagging generate correlated trees.

RANDOM FOREST

Random Forest

15

- **Random forest** is a **bagging** technique and **not a boosting** technique.
- The trees in **random forests** run in parallel.
- There is **no interaction between these trees** while building the trees.

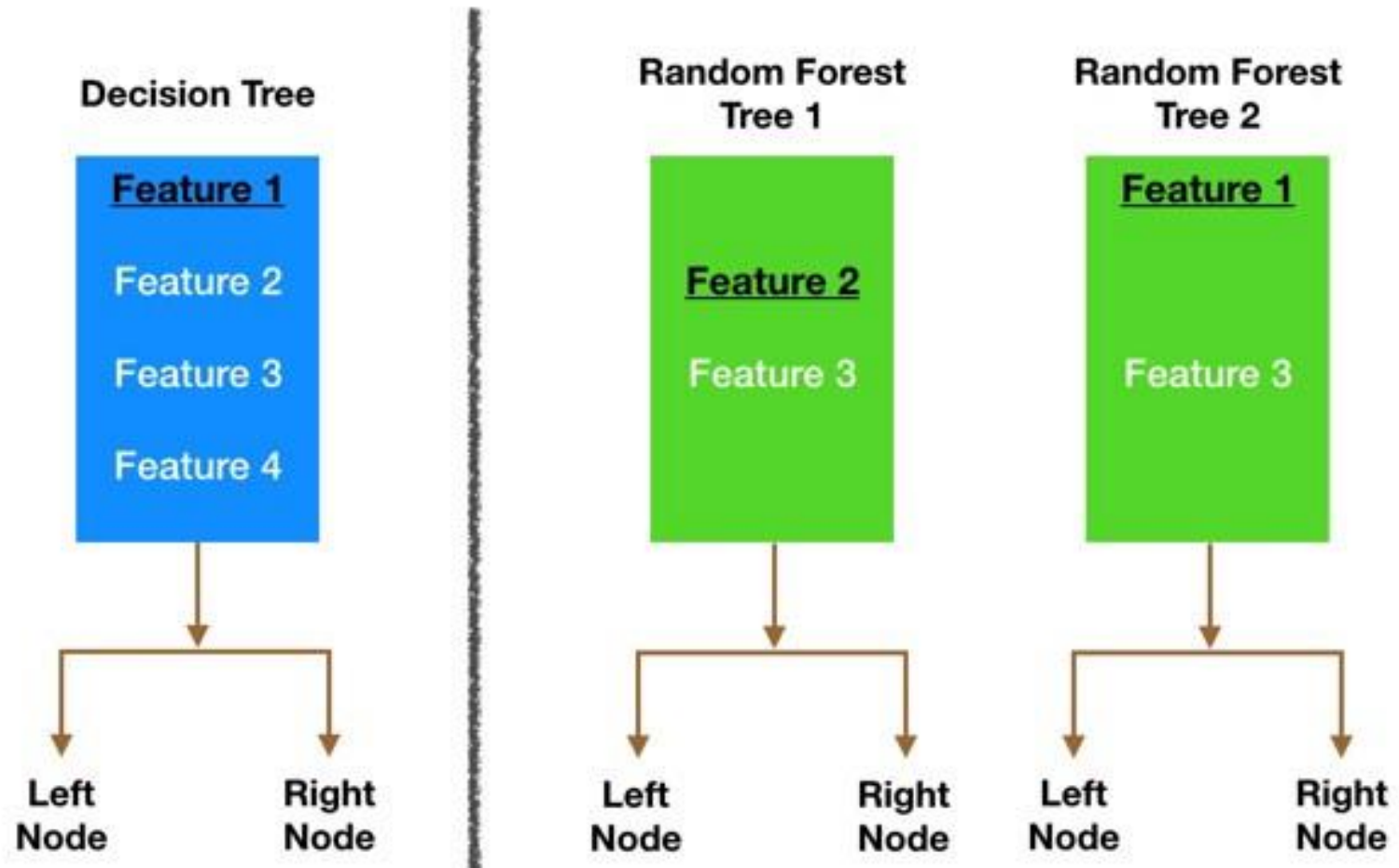
Random Forest

16

- Given n observations with p predictors.
- Input:
 - $m \ll p$ the fraction of the predictors to sample at each split (often $m = \sqrt{p}$)
 - f , the fraction of the data to use for training
 - k , the number of trees in the forest.
- Repeat k times:
 - Choose a training set by choosing $f \times n$ training cases (with replacement). This is called bagging
 - Build a decision tree as follows
 - For each node of the tree, **randomly choose m variables** and **find the best split** from among those m variables
 - Repeat until the full tree is built (no pruning)

Random Forest

17



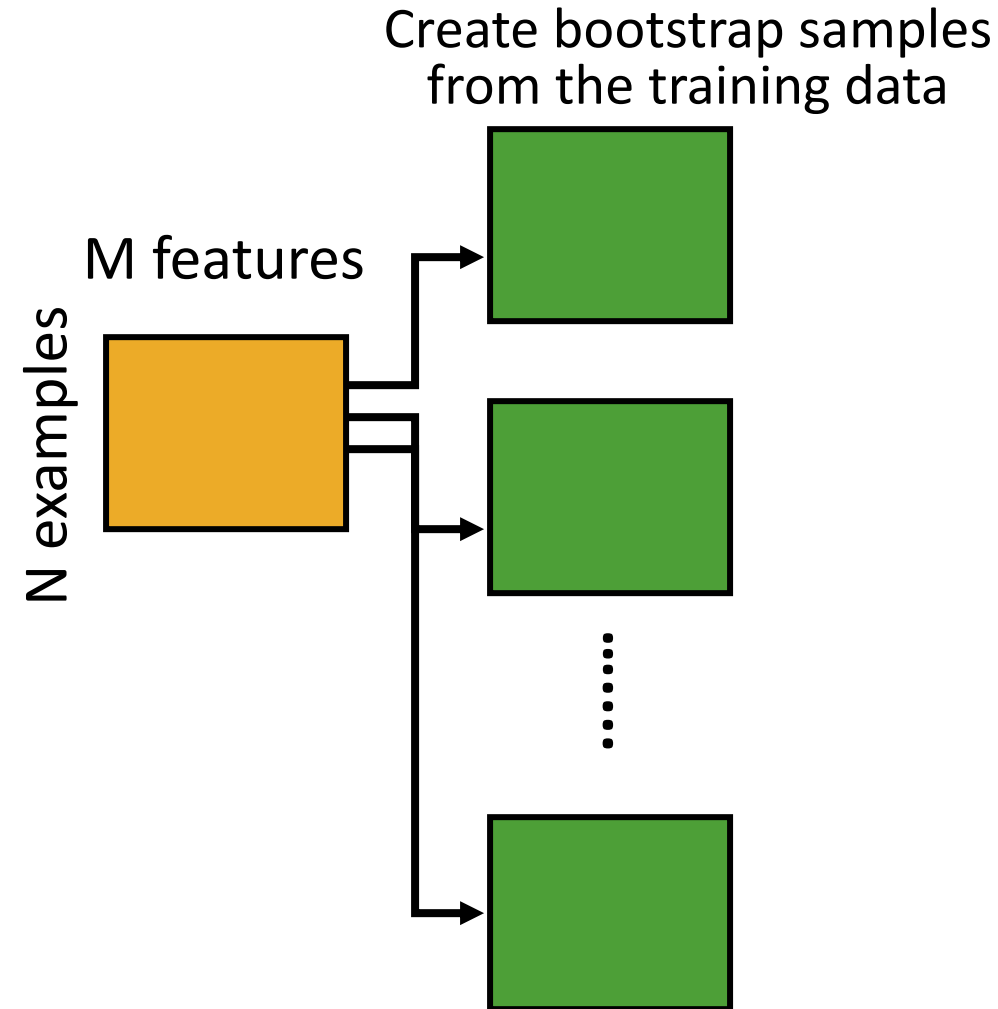
Random Forest

18

- To make a prediction at a new point x , we do:
 - **For regression:**
 - ▣ average the results
 - **For classification:**
 - ▣ majority vote

Random Forest

19



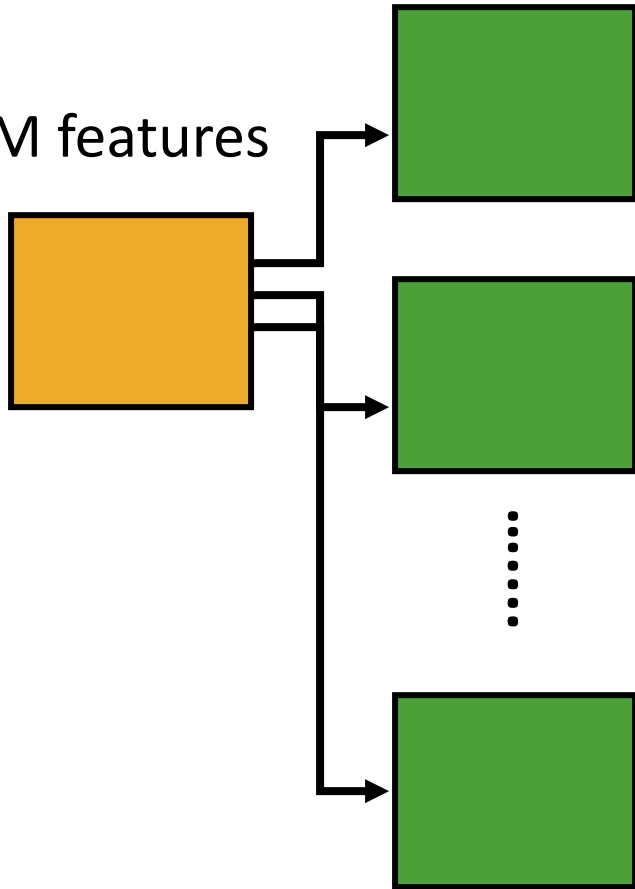
Random Forest

20

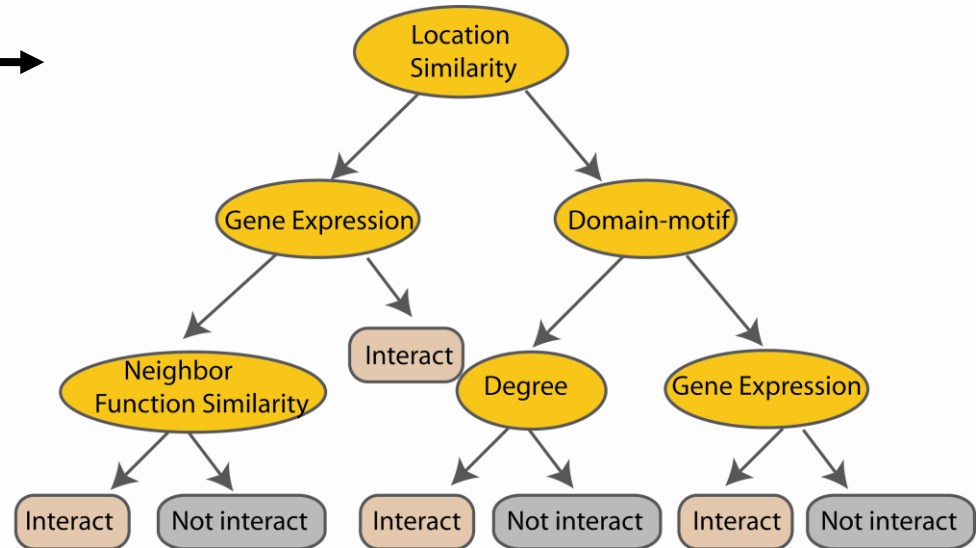
Create bootstrap samples from the training data

M features

N examples



Construct a decision tree



At each node in choosing the **split feature** choose only among $m < M$ features

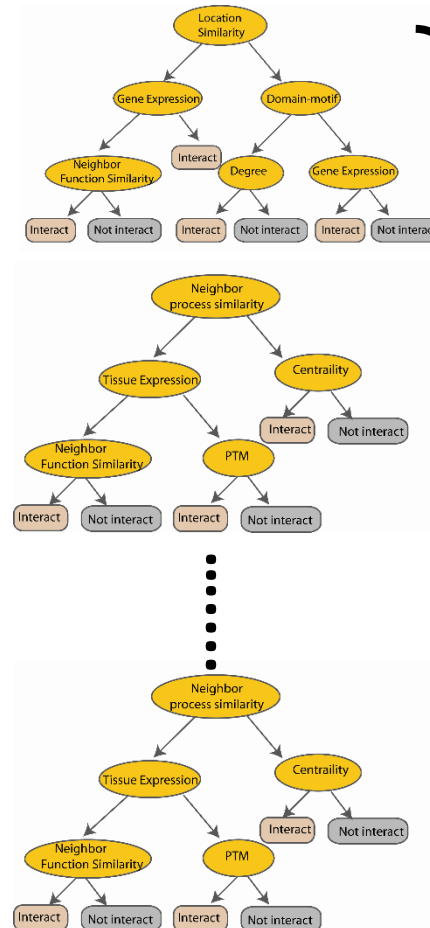
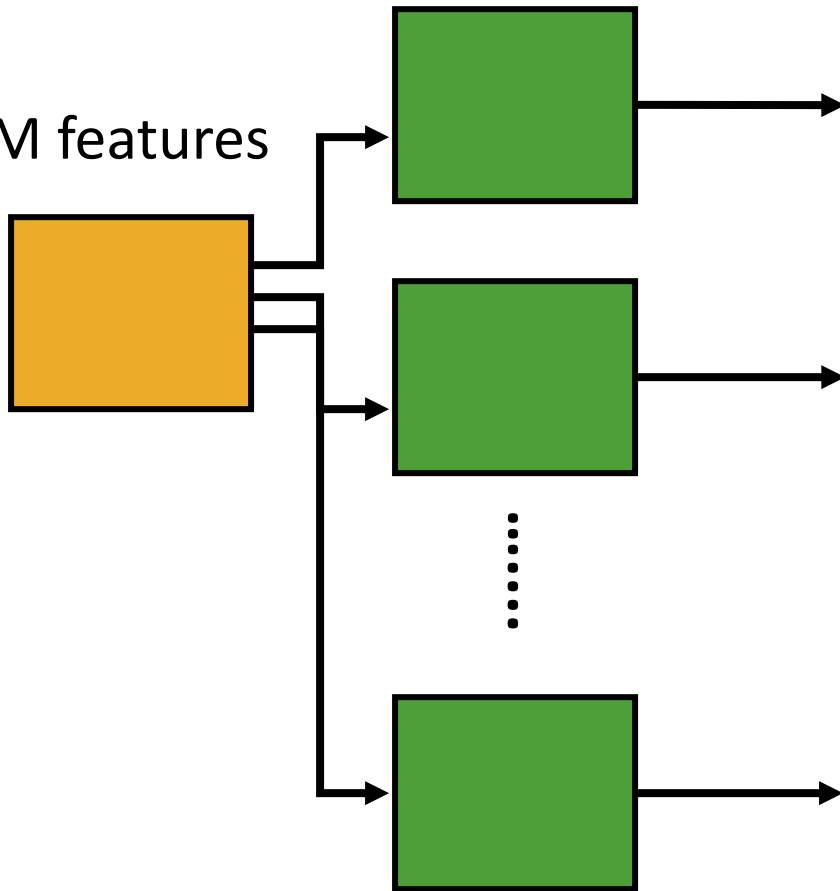
Random Forest

21

Create bootstrap samples
from the training data

M features

N examples



Take the
majority
vote

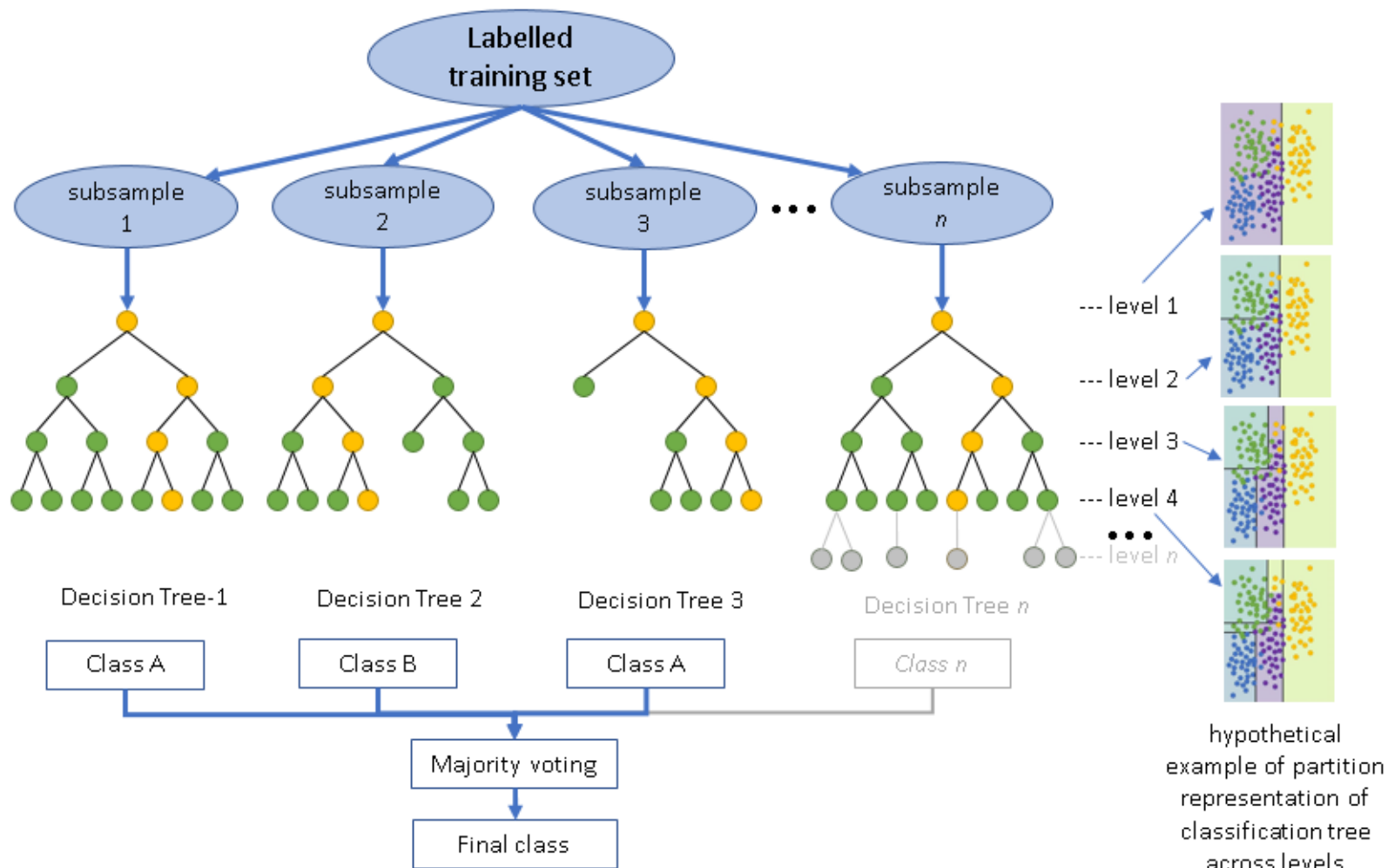
Random Forest

22

Bootstrap
sampling

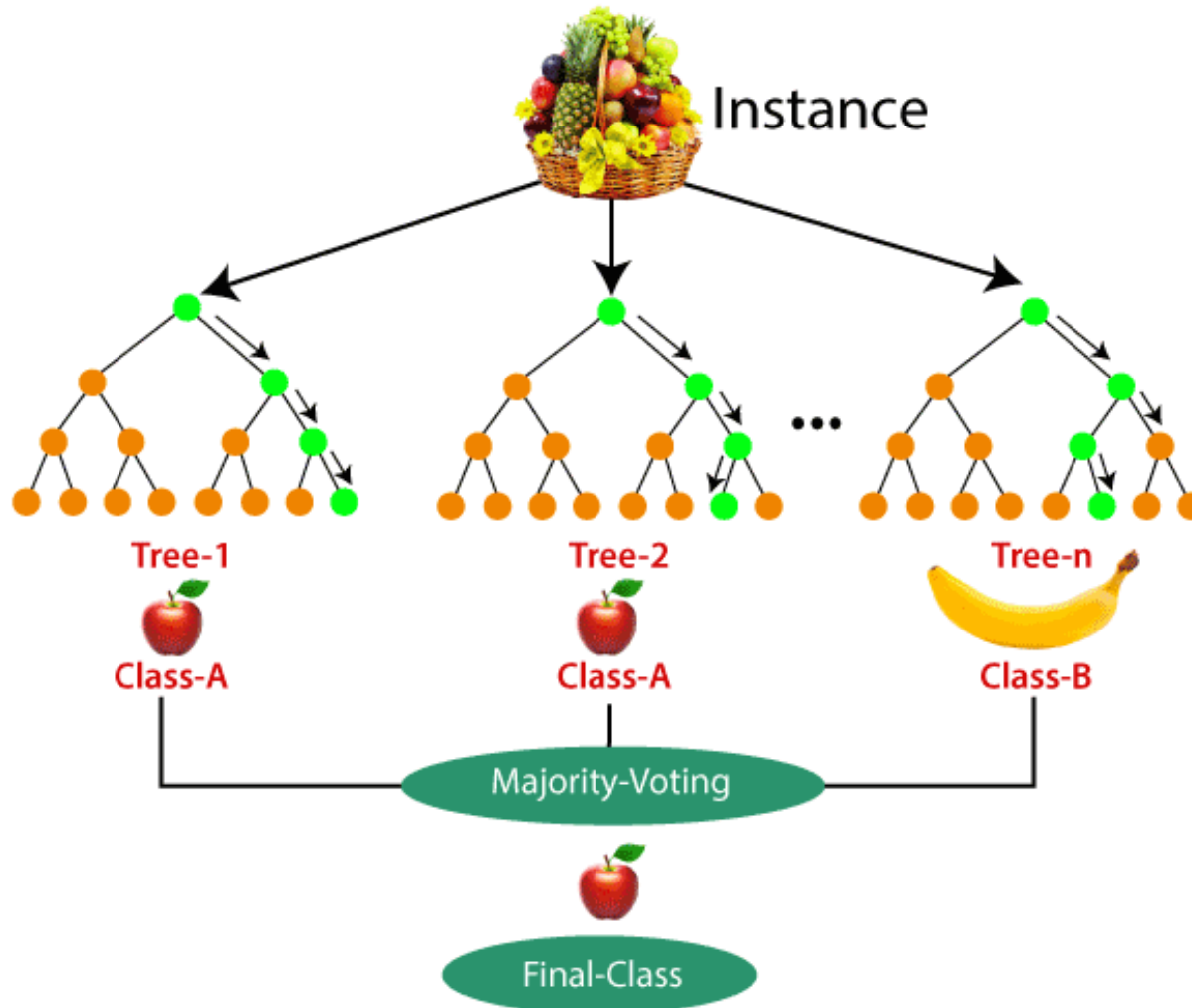
Building the trees
on a random set
of features

Bootstrap
aggregation



Random Forest

23



Random Forest

24

- It is one of the **most accurate learning algorithms** available.
 - ▣ For many datasets, it produces a **highly accurate classifier**.
- It runs **efficiently on large databases**.
- It can **handle thousands of input variables** without variable deletion.
- It gives **estimates of what variables that are important** in the classification.

Random Forest

25

- It generates an internal **unbiased estimate of the generalization error** as the forest building progresses.
- It has an **effective method for estimating missing data**,
 - ▣ maintains accuracy when a large proportion of the data are missing.

Random Forest ... Issues

26

- Random forests have been observed to **overfit for some datasets** with noisy classification/regression tasks.
- For data including *categorical variables with different number of levels*, random forests are **biased in favor of those attributes with more levels**.

ENSEMBLE LEARNING



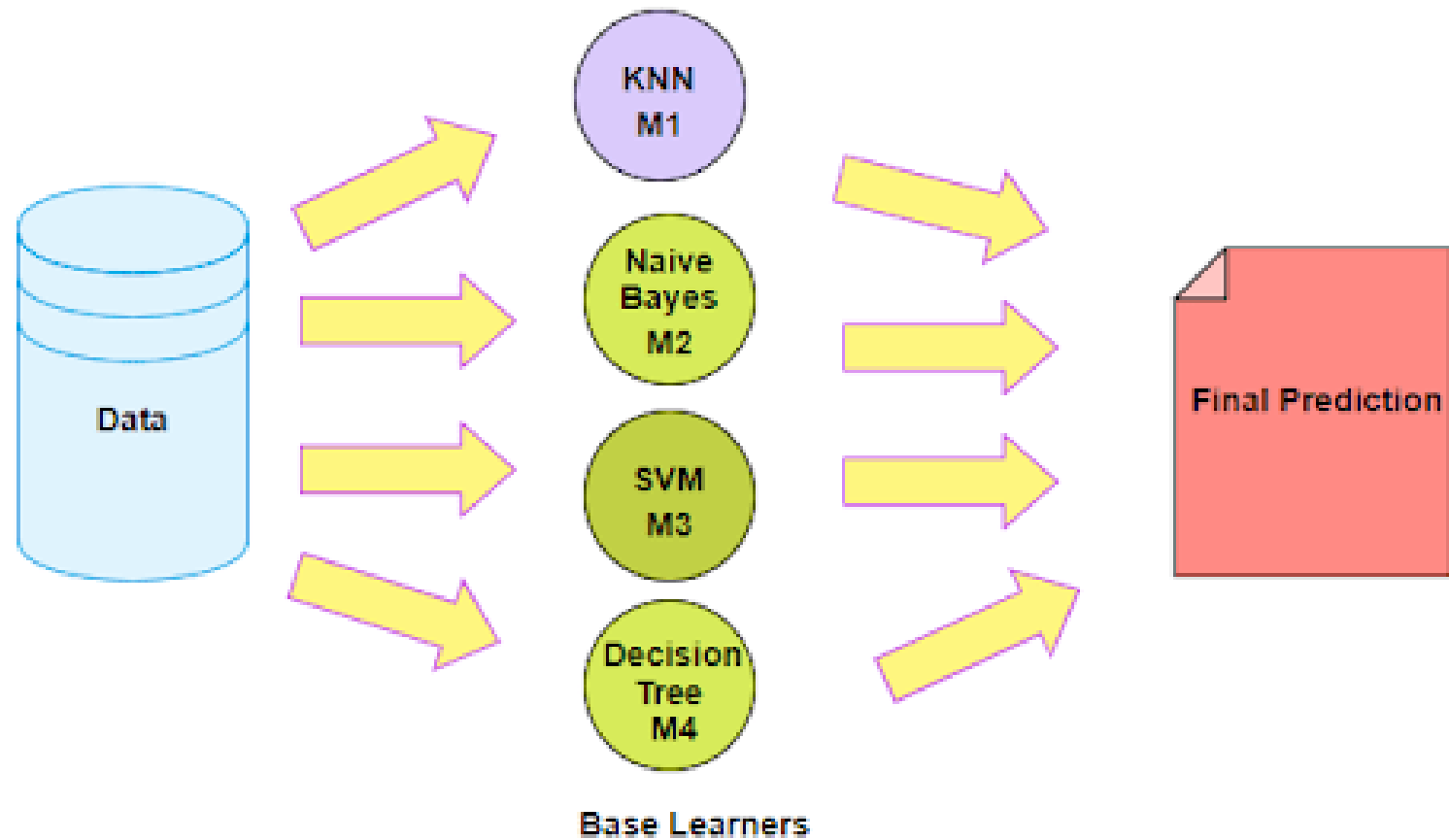
Ensemble Learning

28

- An Ensemble method is a technique that **combines**
 - ▣ the predictions from **multiple** machine learning **algorithms** together to make more accurate predictions than any individual model.
- A model comprised of many models is called an **Ensemble model**.

Ensemble Learning

29



Acknowledgement

30

Tom Mitchel, Russel & Norvig, Andrew Ng, Alpydin & Ch. Eick.