

## Decision Tree: Issues

CSE 4334 / 5334 Data Mining  
Spring 2019

**Won Hwa Kim**

(Slides courtesy of Pang-Ning Tan, Michael Steinbach and Vipin Kumar, and Jiawei Han, Micheline Kamber and Jian Pei)



## Practical Issues of Classification



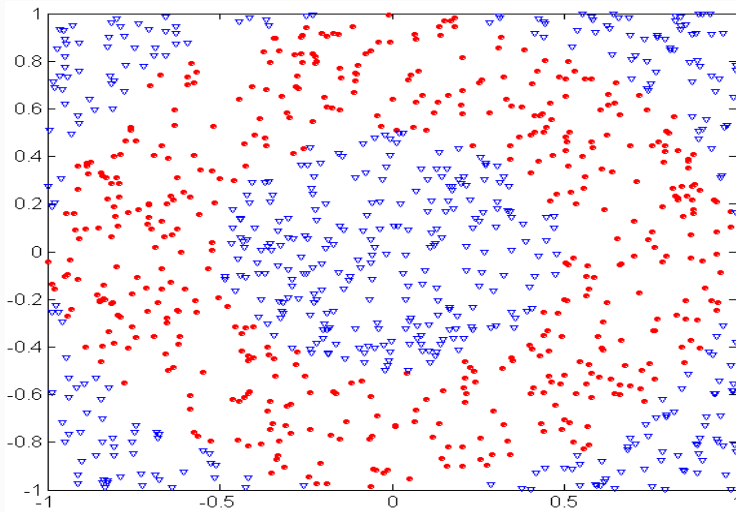
Underfitting and Overfitting

Missing Values

Costs of Classification

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Underfitting and Overfitting (Example)



500 circular and 500 triangular data points.

Circular points:

$$0.5 \leq \sqrt{x_1^2 + x_2^2} \leq 1$$

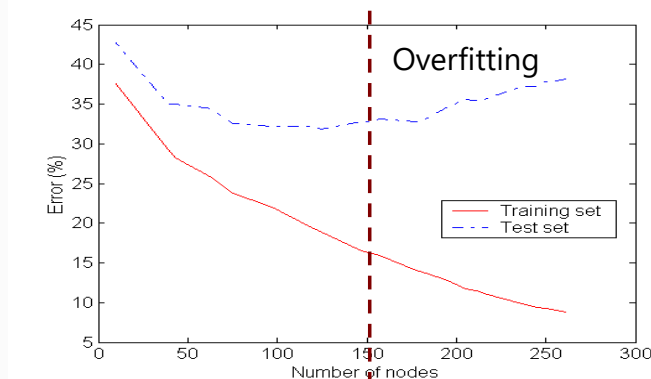
Triangular points:

$$\sqrt{x_1^2 + x_2^2} < 0.5 \text{ or}$$

$$\sqrt{x_1^2 + x_2^2} > 1$$

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Underfitting and Overfitting

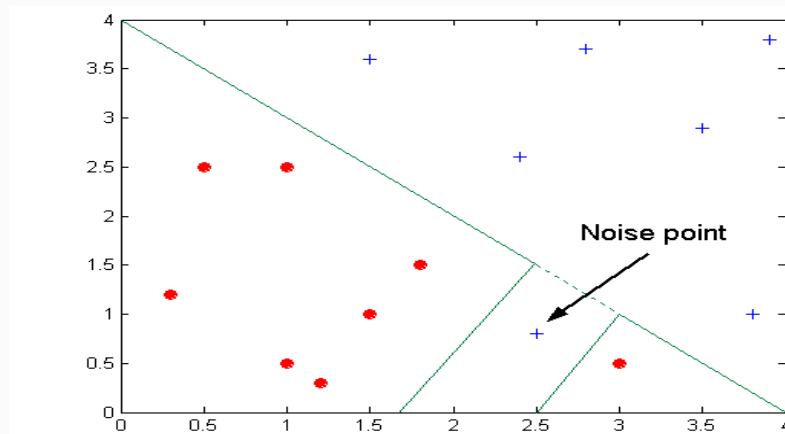


Underfitting: when model is too simple, both training and test errors are large

Overfitting: when model is too complex, test error increases even though training error decreases

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

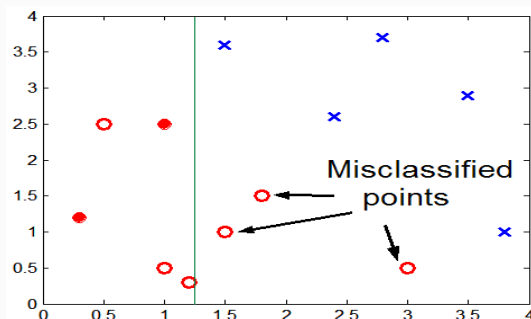
## Overfitting due to Noise



Decision boundary is distorted by noise point

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Overfitting due to Insufficient Examples



Lack of data points in the lower half of the diagram makes it difficult to predict correctly the class labels of that region

- Insufficient number of training records in the region causes the decision tree to predict the test examples using other training records that are irrelevant to the classification task

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Notes on Overfitting



Overfitting results in decision trees that are more complex than necessary

Training error no longer provides a good estimate of how well the tree will perform on previously unseen records

Need new ways for estimating errors

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Estimating Generalization Errors



Re-substitution errors: error on training ( $\sum e(t)$ )

Generalization errors: error on testing ( $\sum e'(t)$ )

Methods for estimating generalization errors:

- o Optimistic approach:  $e'(t) = e(t)$
- o Pessimistic approach:
  - o For each leaf node:  $e'(t) = (e(t) + 0.5)$
  - o Total errors:  $e'(T) = e(T) + N \times 0.5$  (N: number of leaf nodes)
  - o For a tree with 30 leaf nodes and 10 errors on training (out of 1000 instances):
    - Training error =  $10/1000 = 1\%$
    - Generalization error =  $(10 + 30 \times 0.5)/1000 = 2.5\%$
- o Reduced error pruning (REP):
  - o uses validation data set to estimate generalization error

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Occam's Razor



Given two models of similar generalization errors, one should prefer the **simpler model** over the more complex model

For complex models, there is a greater chance that it was fitted accidentally by errors in data

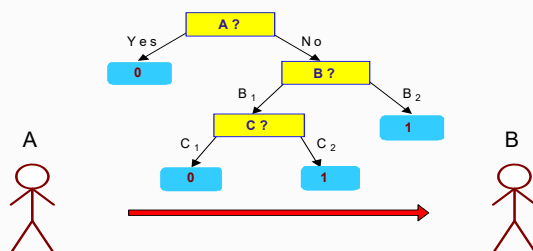
Therefore, one should include model complexity when evaluating a model

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Minimum Description Length (MDL)



X	y
$x_1$	1
$x_2$	0
$x_3$	0
$x_4$	1
...	...
$x_n$	1



X	y
$x_1$	?
$x_2$	?
$x_3$	?
$x_4$	?
...	...
$x_n$	?

$$\text{Cost}(\text{Model}, \text{Data}) = \text{Cost}(\text{Data} | \text{Model}) + \text{Cost}(\text{Model})$$

- o Cost is the number of bits needed for encoding.
- o Search for the least costly model.

$\text{Cost}(\text{Data} | \text{Model})$  encodes the misclassification errors.

$\text{Cost}(\text{Model})$  uses node encoding (number of children) plus splitting condition encoding.

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## How to Address Overfitting



### Pre-Pruning (Early Stopping Rule)

- Stop the algorithm before it becomes a fully-grown tree
- Typical stopping conditions for a node:
  - Stop if all instances belong to the same class
  - Stop if all the attribute values are the same
- More restrictive conditions:
  - Stop if number of instances is less than some user-specified threshold
  - Stop if class distribution of instances are independent of the available features (e.g., using  $\chi^2$  test)
  - Stop if expanding the current node does not improve impurity measures (e.g., Gini or information gain).

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## How to Address Overfitting...



### Post-pruning

- Grow decision tree to its entirety
- Trim the nodes of the decision tree in a bottom-up fashion
- If generalization error improves after trimming, replace sub-tree by a leaf node.
- Class label of leaf node is determined from majority class of instances in the sub-tree
- Can use MDL for post-pruning

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Example of Post-Pruning



Class = Yes	20
Class = No	10
Error = 10/30	

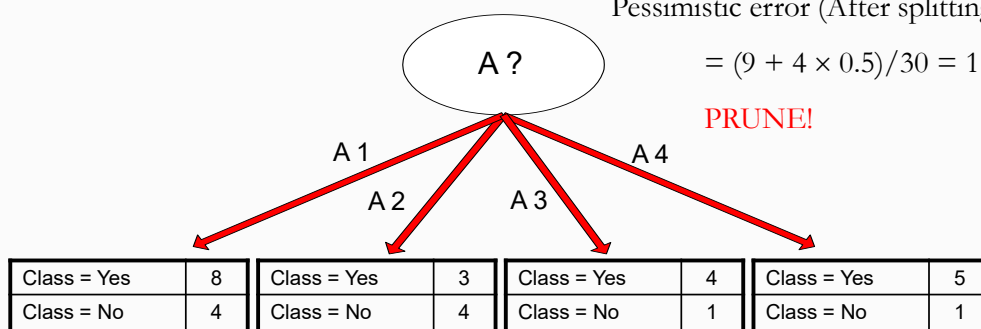
Training Error (Before splitting) = 10/30

Pessimistic error =  $(10 + 0.5)/30 = 10.5/30$

Training Error (After splitting) = 9/30

Pessimistic error (After splitting)  
 $= (9 + 4 \times 0.5)/30 = 11/30$

**PRUNE!**



Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

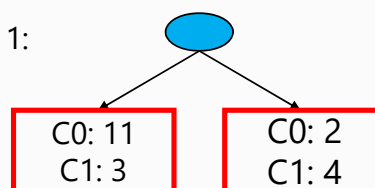
## Examples of Post-pruning



### o Optimistic error?

Don't prune for both cases

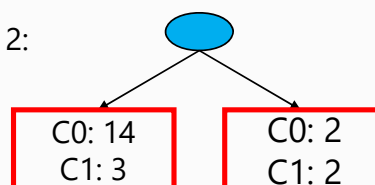
Case 1:



### o Pessimistic error?

Don't prune case 1, prune case 2

Case 2:



### o Reduced error pruning?

Depends on validation set

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Handling Missing Attribute Values

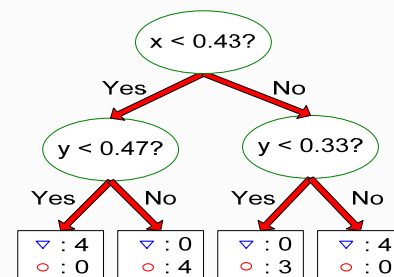
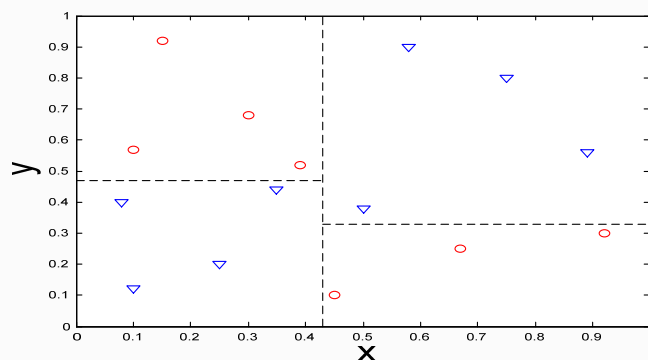


Missing values affect decision tree construction in three different ways:

- Affects how impurity measures are computed
- Affects how to distribute instance with missing value to child nodes
- Affects how a test instance with missing value is classified

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.

## Decision Boundary

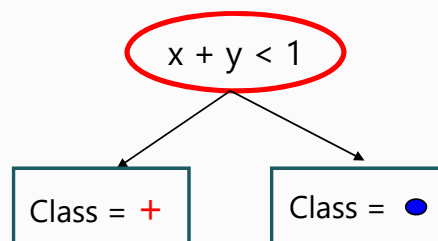
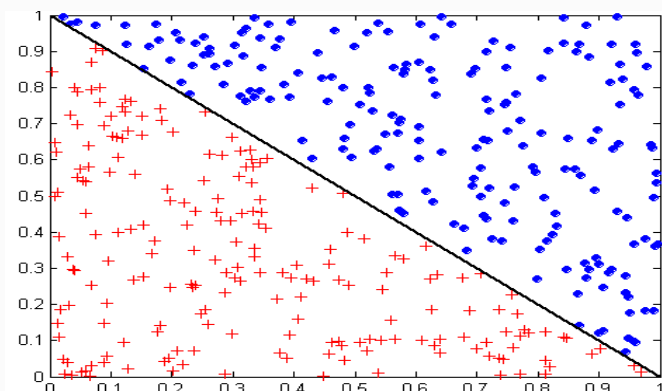


- Border line between two neighboring regions of different classes is known as decision boundary
- Decision boundary is parallel to axes because test condition involves a single attribute at-a-time

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.



## Oblique Decision Trees



- Test condition may involve multiple attributes
- More expressive representation
- Finding optimal test condition is computationally expensive

Copyright ©2007-2017 The University of Texas at Arlington. All Rights Reserved.