

An artificial example

- We'll create a training dataset

Five inputs, all bits, are generated in all 32 possible combinations

Output y = copy of e , Except a random 25% of the records have y set to the opposite of e

32 records

a	b	c	d	e	y
0	0	0	0	0	0
0	0	0	0	1	0
0	0	0	1	0	0
0	0	0	1	1	1
0	0	1	0	0	1
:	:	:	:	:	:
1	1	1	1	1	1

© Andrew W. Moore

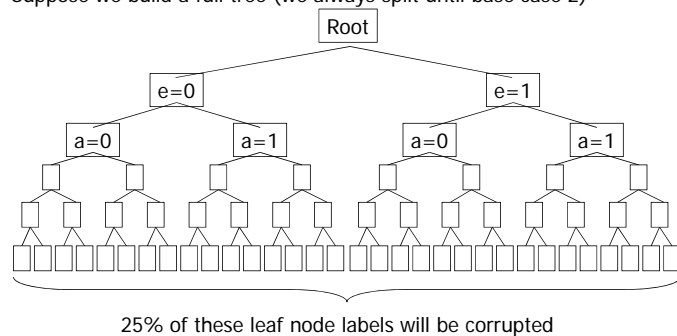
In our artificial example

- Suppose someone generates a test set according to the same method.
- The test set is identical, except that some of the y 's will be different.
- Some y 's that were corrupted in the training set will be uncorrupted in the testing set.
- Some y 's that were uncorrupted in the training set will be corrupted in the test set.

© Andrew W. Moore

Building a tree with the artificial training set

- Suppose we build a full tree (we always split until base case 2)



© Andrew W. Moore

Training set error for our artificial tree

All the leaf nodes contain exactly one record and so...

- We would have a training set error of zero

© Andrew W. Moore

Testing the tree with the test set

	1/4 of the tree nodes are corrupted	3/4 are fine
1/4 of the test set records are corrupted	1/16 of the test set will be correctly predicted for the wrong reasons	3/16 of the test set will be wrongly predicted because the test record is corrupted
3/4 are fine	3/16 of the test predictions will be wrong because the tree node is corrupted	9/16 of the test predictions will be fine

In total, we expect to be wrong on 3/8 of the test set predictions

© Andrew W. Moore

What's this example shown us?

- This explains the discrepancy between training and test set error
- But more importantly... it indicates there's something we should do about it if we want to predict well on future data.

© Andrew W. Moore

Suppose we had less data

- Let's not look at the irrelevant bits

These bits are hidden

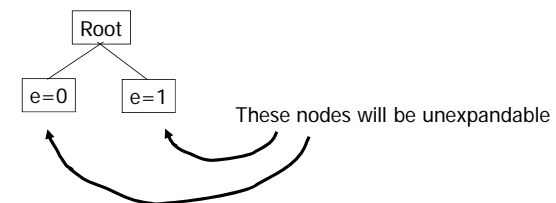
Output y = copy of e , except a random 25% of the records have y set to the opposite of e

	a	b	c	d	e	y
32 records	0	0	0	0	0	0
	0	0	0	0	1	0
	0	0	0	1	0	0
	0	0	0	1	1	1
	0	0	1	0	0	1
	:	:	:	:	:	:
	1	1	1	1	1	1

What decision tree would we learn now?

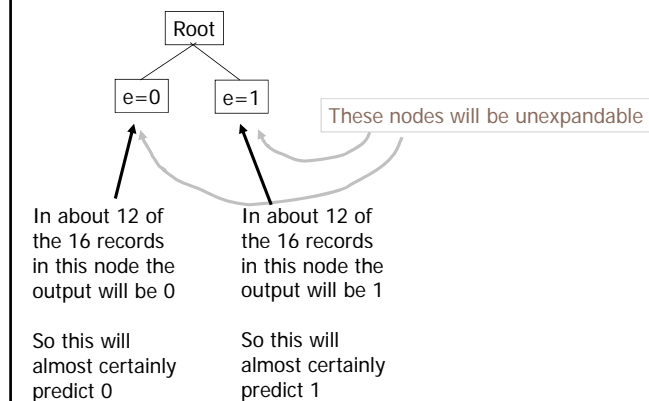
© Andrew W. Moore

Without access to the irrelevant bits...



© Andrew W. Moore

Without access to the irrelevant bits...



© Andrew W. Moore

Without access to the irrelevant bits...

	almost certainly none of the tree nodes are corrupted	almost certainly all are fine
1/4 of the test set records are corrupted	n/a	1/4 of the test set will be wrongly predicted because the test record is corrupted
3/4 are fine	n/a	3/4 of the test predictions will be fine

In total, we expect to be wrong on only 1/4 of the test set predictions

© Andrew W. Moore

Overfitting

- Definition: If your machine learning algorithm fits noise (i.e. pays attention to parts of the data that are irrelevant) it is **overfitting**
- Fact (theoretical and empirical): If your machine learning algorithm is overfitting then it may perform less well on test set data

© Andrew W. Moore