

Revenge of the Entropy



Note we're working with continuous splits here, which is an interesting problem we covered in the last lecture.



What if points have exactly same classes? Output whichever is the majority. If there's no majority, you could fip a coin, or just pick a "default" class.





We saw that if we make enough splits, we can pretty much get 100% accuracy (with some exceptions if points of different classes have the same attributes)

Welcome to the real world, where data is noisy and the answers aren't in the back of the book.

This is not only a noise problem, but a matter of sparse data. Suppose you had one female age 37 with brown hair and a white blood cell count of X and a family history of cancer, and she has cancer. Does that mean that any 37-year-old women brunettes with wbc count X and a family history have cancer?

Training and Test sets are a big deal. You never use all the data to train. Use some fraction of the data for training, and the other is testing.

When you're training, you're automatically going to overfit. Any data sample will have bias and noise. So you use the test data to test its "real" accuracy.

Oh hey we have 100% accuracy, awesome! Oh wait.

Well this improves our test a little bit, if we try pruning based on our training data (which still has bias)

Prune it a little more and the error is mostly from noise.

Why is testing usually below training accuracy? Because the training data is what we're actually trying to fit on.

Basic Questions

- How to choose the attribute/value to split on at each level of the tree?
- When to stop splitting? When should a node be declared a leaf?
- If a leaf node is impure, how should the class label be assigned?
- If the tree is too large, how can it be pruned?

Possible Overfitting Solutions

- Grow tree based on training data (*unpruned* tree)
- Prune the tree by removing useless nodes based on:
 - Additional test data (not used for training)
 - Statistical significance tests

This keeps us from going too far down in the tree.

You need to choose a threshold, which is hacky.

Decision Tree Pruning

- Construct the entire tree as before
- Starting at the leaves, recursively eliminate splits:

– At a leaf \mathcal{N} :

- Compute the *K* value for \mathcal{N} and its parent \mathcal{P} .
- If the *K* values is lower than the threshold *t*: – Eliminate all of the children of *P*
 - 𝒫 becomes a leaf
- Repeat until no more splits can be eliminated

When reporting accuracy, you'll still want to use the test set.

χ^2 Pruning

- The test on *K* is a version of a standard statistical test, the χ^2 ('chi-square') test.
- The value of *t* is retrieved from statistical tables. For example, K > t means that, with confidence 95%, the information gain due to the split is significant.
- If K < t, with high confidence, the information gain will be 0 over very large training samples
 - Reduces overfitting
 - Eliminates irrelevant attributes

Now here's a big problem. How do we decide what "features" to use?

For example, polynomial fitting. We could use linear regression, or quadratic, or ...

OR THIS. AWESOME.

Of course, there's a tradeoff between accuracy and simplicity. Note that on the right we have a lower accuracy. Chances are we can say that other point is an error, or noise, but it's not always that simple. Sometimes there's a big gap in accuracy for an only-slightly-simpler model. How do we decide how complicated we're willing to make it? For more on this, read the wikipedia article on regularization:

http://en.wikipedia.org/wiki/Regularization_(mathematics)