

INDUCTIVE BIAS IN DECISION TREE LEARNING

Inductive bias is the set of assumptions that, together with the training data, deductively justify the classifications assigned by the learner to future instances

Given a collection of training examples, there are typically many decision trees consistent with these examples. Which of these decision trees does ID3 choose?

ID3 search strategy

- (a) selects in favour of shorter trees over longer ones
- (b) selects trees that place the attributes with highest information gain closest to the root.

Approximate inductive bias of ID3: Shorter trees are preferred over larger trees

- Consider an algorithm that begins with the empty tree and searches *breadth first* through progressively more complex trees.
- First considering all trees of depth 1, then all trees of depth 2, etc.
- Once it finds a decision tree consistent with the training data, it returns the smallest consistent tree at that search depth (e.g., the tree with the fewest nodes).
- Let us call this breadth-first search algorithm BFS-ID3.
- BFS-ID3 finds a shortest decision tree and thus exhibits the bias "shorter trees are preferred over longer trees."

- ID3 can be viewed as an efficient approximation to BFS-ID3, using a greedy heuristic search to attempt to find the shortest tree without conducting the entire breadth-first search through the hypothesis space.
- Because ID3 uses the information gain heuristic and a hill climbing strategy, it exhibits a more complex bias than BFS-ID3.
- In particular, it does not always find the shortest consistent tree, and it is biased to favour trees that place attributes with high information gain closest to the root.

Restriction Biases and Preference Biases



Difference between the types of inductive bias exhibited by ID3 and by the CANDIDATE-ELIMINATION Algorithm.

ID3

- ID3 searches a complete hypothesis space
- It searches incompletely through this space, from simple to complex hypotheses, until its termination condition is met
- Its inductive bias is solely a consequence of the ordering of hypotheses by its search strategy. Its hypothesis space introduces no additional bias

CANDIDATE-ELIMINATION Algorithm

- The version space CANDIDATE-ELIMINATION Algorithm searches an incomplete hypothesis space
- It searches this space completely, finding every hypothesis consistent with the training data.
- Its inductive bias is solely a consequence of the expressive power of its hypothesis representation. Its search strategy introduces no additional bias

Restriction Biases and Preference Biases



- The inductive bias of ID3 is *a preference* for certain hypotheses over others (e.g., preference for shorter hypotheses over larger hypotheses), with no hard restriction on the hypotheses that can be eventually enumerated. This form of bias is called *a preference bias* or *a search bias*.
- The bias of the CANDIDATE ELIMINATION algorithm is in the form of *a categorical* restriction on the set of hypotheses considered. This form of bias is typically called *a restriction bias* or *a language bias*.

Which type of inductive bias is preferred in order to generalize beyond the training data: a preference bias or a restriction bias?

- A preference bias is more desirable than a restriction bias, because it allows the learner to work within a complete hypothesis space that is assured to contain the unknown target function.
- In contrast, a restriction bias that strictly limits the set of potential hypotheses is generally less desirable, because it introduces the possibility of excluding the unknown target function altogether.

Occam's razor

Occam's razor: is the problem-solving principle that the simplest solution tends to be the right one. When presented with competing hypotheses to solve a problem, one should select the solution with the fewest assumptions.

Occam's razor: *“Prefer the simplest hypothesis that fits the data”*.

Why Prefer Short Hypotheses?

Argument in favour:

Fewer short hypotheses than long ones:

- Short hypotheses fit the training data which are *less likely to be coincident*
- Longer hypotheses fit the training data might be *coincident*.

Many complex hypotheses that fit the current training data but fail to generalize correctly to subsequent data.

Argumentopposed:

- There are few small trees, and our priori chance of finding one consistent with an arbitrary set of data is therefore small. The difficulty here is that there are very many small set of hypotheses that one can define *but understood by few learner*.
- The size of a hypothesis is determined by the representation used *internally* by the learner. Occam's razor will produce *two different hypotheses from the same training examples when it is applied by two learners*, both justifying their contradictory conclusions by Occam's razor. On this basis we might be tempted to reject Occam's razor altogether.

ISSUES IN DECISION TREE LEARNING

1. Avoiding Overfitting the Data

Reduced error

pruning Rule post-pruning

2. Incorporating Continuous-Valued Attributes

3. Alternative Measures for Selecting Attributes

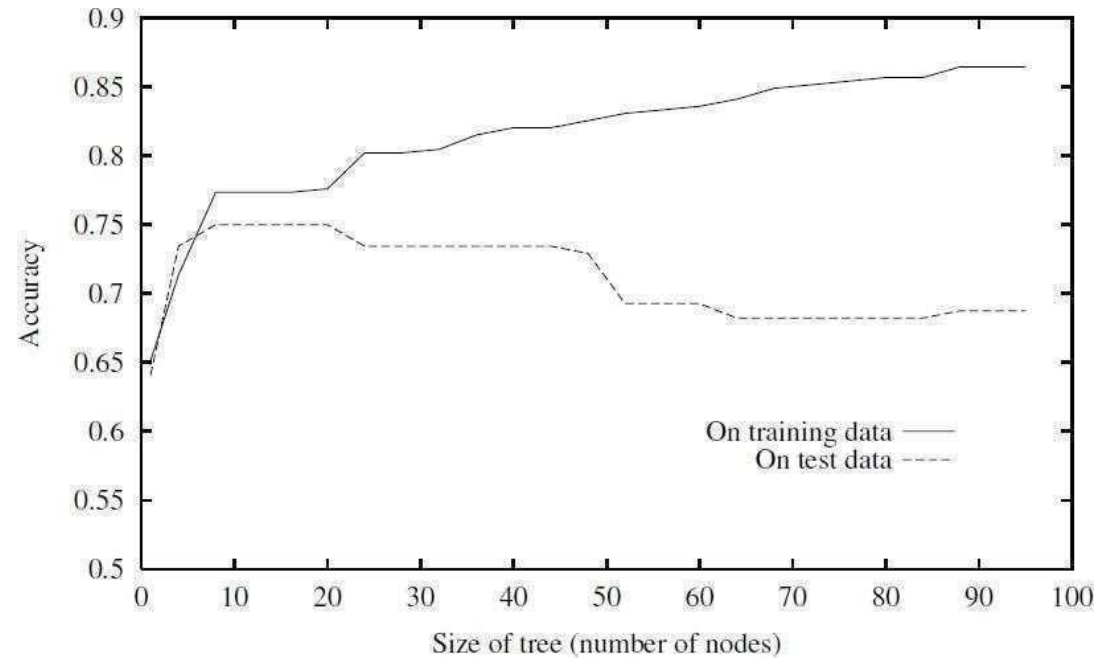
4. Handling Training Examples with Missing Attribute Values

5. Handling Attributes with Differing Costs

1. Avoiding Overfitting the Data

- The ID3 algorithm grows each branch of the tree just deeply enough to perfectly classify the training examples but it can lead to difficulties when there is noise in the data, or when the number of training examples is too small to produce a representative sample of the true target function. This algorithm can produce trees that *overfit* the training examples.
- **Definition-Overfit:** Given a hypothesis space H , a hypothesis $h \in H$ is said to overfit the training data if there exists some alternative hypothesis $h' \in H$, such that h has a smaller error than h' over the training examples, but h' has a smaller error than h over the entire distribution of instances.

- The below figure illustrates the impact of overfitting in a typical application of decision tree learning.



- The **horizontal axis** of this plot indicates the total number of nodes in the decision tree, as the tree is being constructed. The **vertical axis** indicates the accuracy of predictions made by the tree.
- The **solid line** shows the accuracy of the decision tree over the training examples. The **broken line** shows accuracy measured over an independent set of test examples.
- The **accuracy of the tree over the training examples increases monotonically** as the tree is grown. The **accuracy measured over the independent test examples first increases, then decreases.**

How can it be possible for a tree to fit the training examples better than h' , but for it to perform more poorly on subsequent examples?

1. Overfitting can occur when the training examples contain random errors or noise
2. When small numbers of examples are associated with leaf nodes.

Noisy Training Example

Example 15: $\langle \text{Sunny}, \text{Hot}, \text{Normal}, \text{Strong}, - \rangle$

- Example is noisy because the correct label is +
- Previously constructed tree misclassifies it

