

□ Criterion used to determine the correct final tree size

- Use a separate set of examples, distinct from the training examples, to evaluate the utility of post-pruning nodes from the tree
- Use all the available data for training, but apply *a statistical test* to estimate whether expanding (or pruning) a particular node is likely to produce an improvement beyond the training set
- Use a measure of the complexity for encoding the training examples and the decision tree, halting growth of the tree when this encoding size is minimized. This approach is the Minimum Description Length (or called

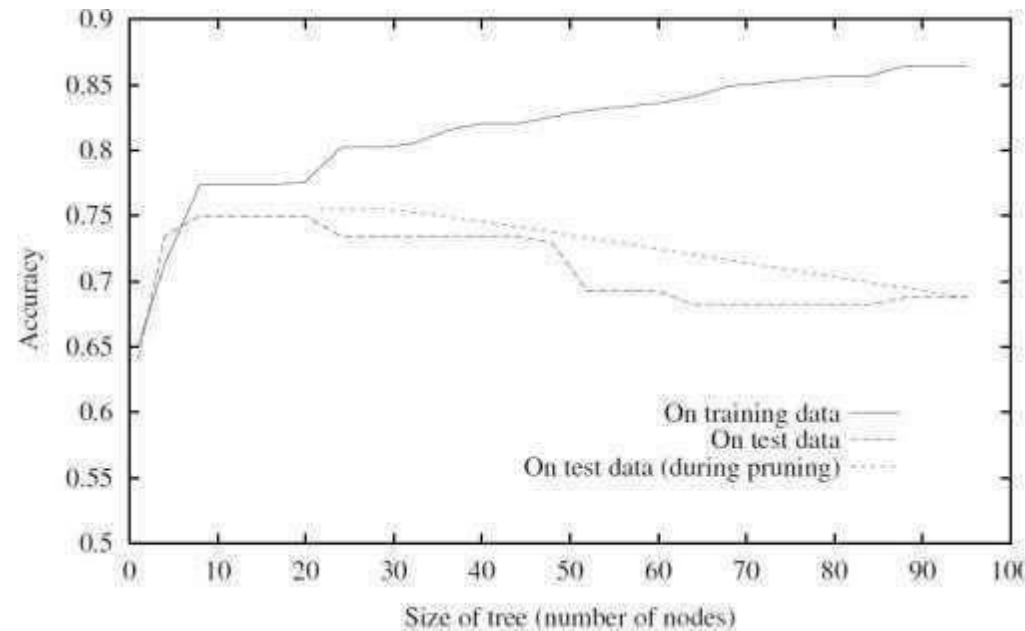
□ *MDL*–

Minimize: size(tree) + size(misclassifications(tree))

Reduced-Error Pruning

- **Reduced-error pruning**, is to consider each of the decision nodes in the tree to be candidates for pruning
- **Pruning** a decision node consists of removing the subtree rooted at that node, making it a leaf node, and assigning it the most common classification of the training examples affiliated with that node
- Nodes are removed only if the resulting pruned tree performs no worse than the original over the validation set.
- Reduced error pruning has the effect that any leaf node added due to coincidental regularities in the training set is likely to be pruned because these same coincidences are unlikely to occur in the validation set

The impact of reduced-error pruning on the accuracy of the decision tree is illustrated in below figure



- The additional line in figure shows *accuracy over the test examples* as the tree is pruned. When pruning begins, the tree is at its maximum size and lowest accuracy over the test set. As pruning proceeds, the number of nodes is reduced and accuracy over the test set increases.
- The available data has been split into three subsets: the training examples, the validation examples used for pruning the tree, and a set of test examples used to provide an unbiased estimate of accuracy over future unseen examples. The plot shows accuracy over the training and test sets.

Summary

- Statistical theory provides a basis for estimating the true error ($error_D(h)$) of hypothesis h , based on its observed error ($error_{sh}$) over a sample S of data.
- In general, the problem of estimating confidence intervals is approached by identifying the parameter to be estimated ($error_D(h)$) and an estimator ($error_{sh}$) for this quantity.
- Because the estimator is a random variable it can be characterised by its probability distribution that governs its value.

unbiased estimator, the observed value of the estimator is likely to vary from one experiment to another.

The variance of the distribution governing the estimate is likely to or characterises how widely this

Confidence intervals can then be calculated by determining the interval that contains the desired probability mass under this distribution.

A cause of estimation error is the variance in the estimate. Even with an

Prodigy and Explanation Based Learning

Prodigy defines a set of target concepts to learn, e.g., which operator given the current state takes you to the goal state?

An example of a rule learned by Prodigy in the block-stacking problem is:

IF One subgoal to be solved is $On(x,y)$ AND
 One subgoal to be solved is $On(y,z)$
THEN Solve the subgoal $On(y,z)$ before $On(x,y)$

Prodigy and Explanation Based Learning

The rationale behind the rule is that it would avoid a conflict when stacking blocks.

Prodigy learns by first encountering a conflict, then explaining the reason for the conflict and creating a rule like the one above.

Experiments show an improvement in efficiency by a factor of two to four.

Problems with EBL

- ✓ The number of control rules that must be learned is very large.
- ✓ If the control rules are many, much time will be spent looking for the best rule.
Utility analysis is used to determine what rules to keep and what rules to forget.

Prodigy:

328 possible rules → 69 pass test → 19 were retained

Problems with EBL

- ✓ Another problem with EBL is that it is sometimes intractable to create an explanation for the target concept.

For example, in chess, learning a concept like:
“states for which operator A leads to a solution”
The search here grows exponentially.

Summary

- ❖ Different from inductive learning, analytical learning looks for a hypothesis that fit the background knowledge and covers the training examples.
- ❖ Explanation based learning is one kind of analytical learning that divides into three steps:
 - a. Explain the target value for the current example
 - b. Analyze the explanation (generalize)
 - c. Refine the hypothesis

Summary

- ❖ Prolog-EBG constructs intermediate features after analyzing examples.
- ❖ Explanation based learning can be used to find search control rules.
- ❖ In all cases we depend on a perfect domain theory.



Machine Learning

Chapter 12. Combining Inductive and Analytical Learning

Tom M. Mitchell



Inductive and Analytical Learning

Inductive learning

- Hypothesis fits data
- Statistical inference
- Requires little prior knowledge
- Syntactic inductive bias

Analytical learning

- Hypothesis fits domain the
- Deductive inference
- Learns from scarce data
- Bias is domain theory



What We Would Like

Inductive learning

Analytical learning



General purpose learning method:

- No domain theory → learn as well as inductive methods
- Perfect domain theory → learn as well as Prolog-EBG
- Accomodate arbitrary and unknown errors in domain theory
- Accomodate arbitrary and unknown errors in training data



Domain theory:

- Cup ← Stable, Lifiable, Open Vessel
- Stable ← BottomIsFlat
- Lifiable ← Graspable, Light
- Graspable ← HasHandle
- Open Vessel ← HasConcavity, ConcavityPointsUp

Training examples:

	Cups				Non-Cups					
BottomIsFlat	✓	✓	✓	✓	✓	✓	✓			✓
ConcavityPoints Up	✓	✓	✓	✓	✓		✓	✓		
Expensive	✓		✓				✓		✓	
Fragile	✓	✓			✓	✓		✓		✓
HandleOnTop					✓		✓			
HandleOnSide	✓			✓					✓	
HasConcavity	✓	✓	✓	✓	✓		✓	✓	✓	✓
HasHandle	✓			✓	✓		✓		✓	✓
Light	✓	✓	✓	✓	✓	✓	✓		✓	
MadeOfCeramic	✓				✓		✓	✓		
MadeOfPaper				✓					✓	
MadeOfStyrofoam		✓	✓			✓				✓

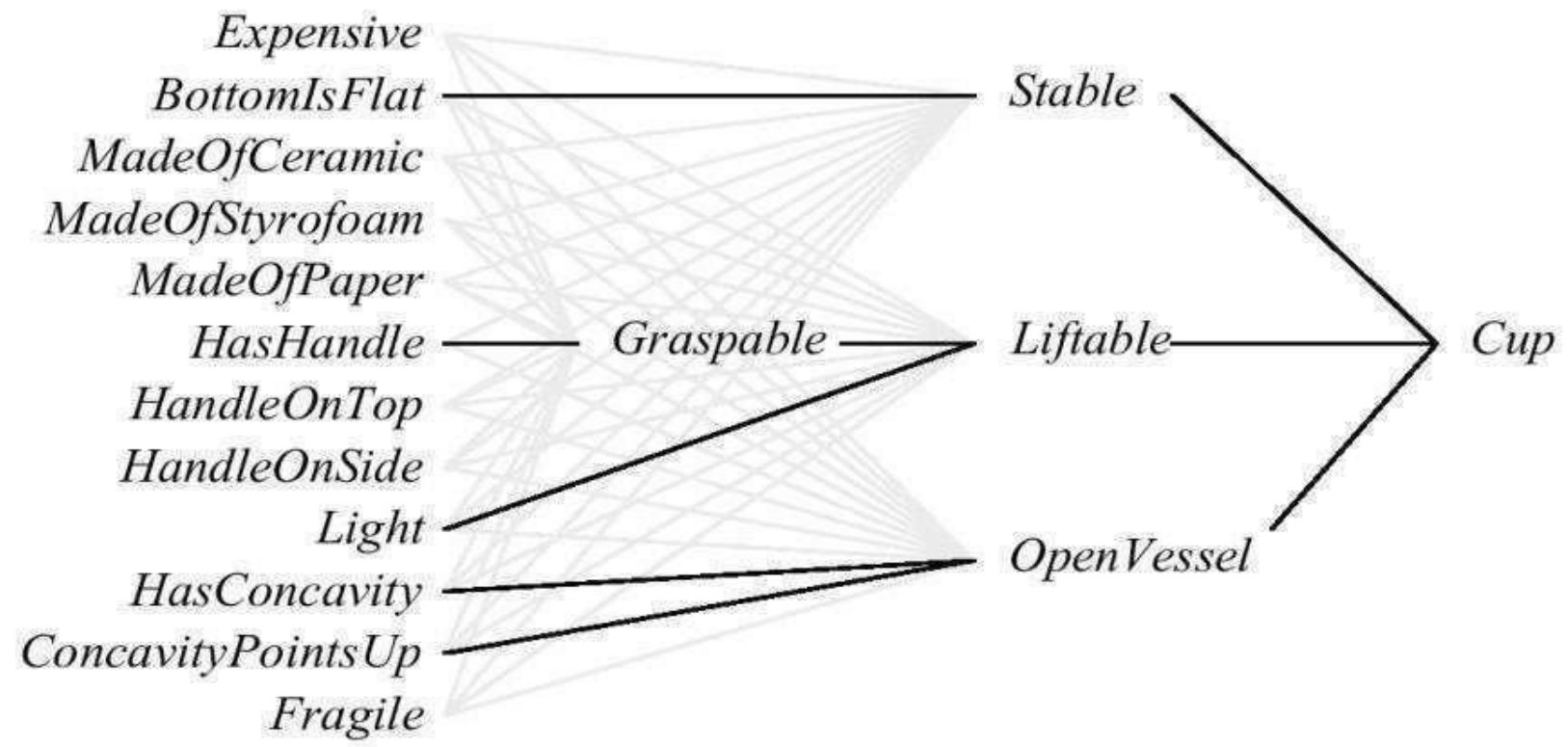


KBANN

- KBANN (data D , domain theory B)
 1. Create a feedforward network h equivalent to B
 2. Use BACKPROP to tune h to $t D$



Neural Net Equivalent to Domain Theory





Creating Network Equivalent to Domain Theory

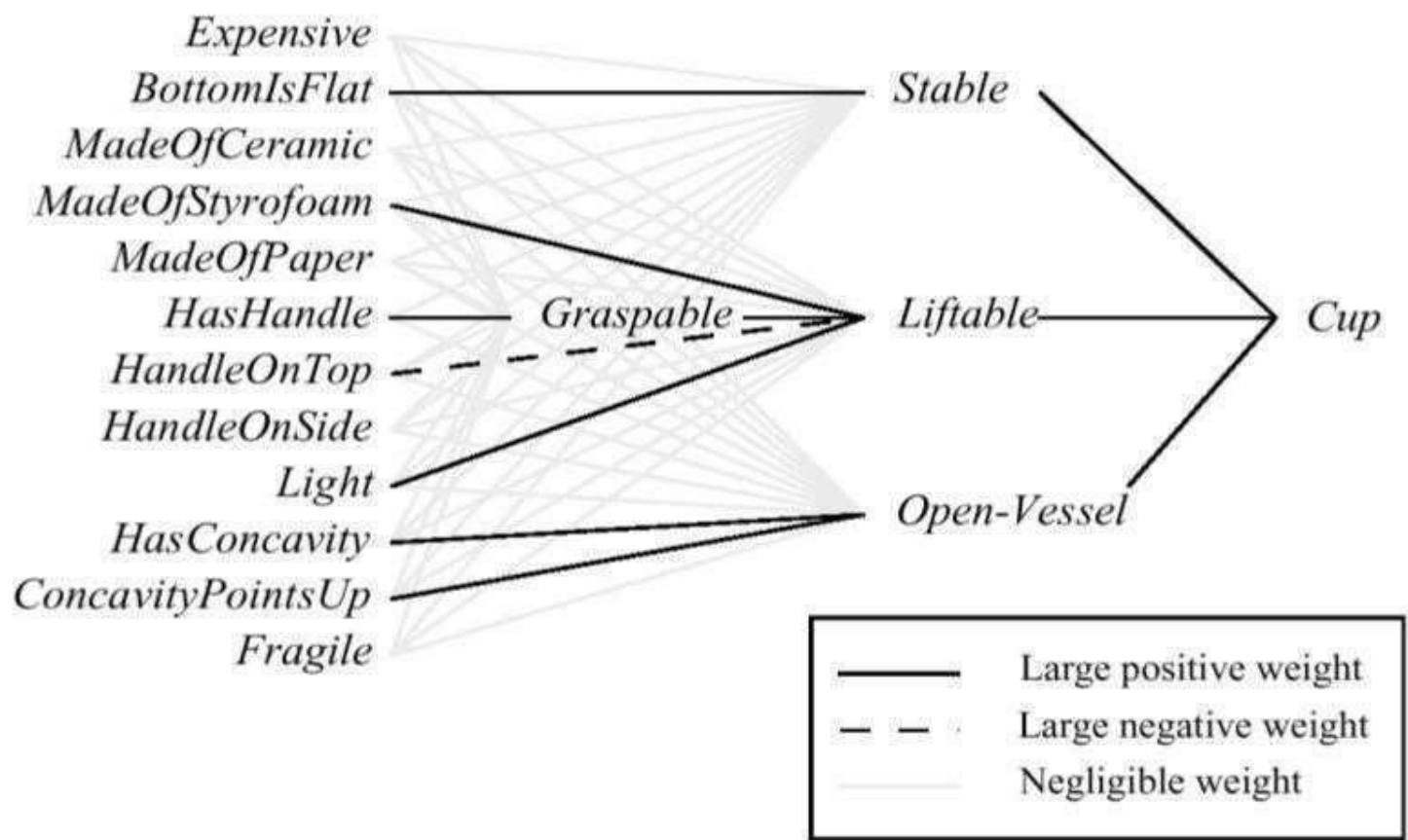
Create one unit per horn clause rule (i.e., an AND unit)

- Connect unit inputs to corresponding clause antecedents
- For each non-negated antecedent, corresponding input weight $w \leftarrow W$, where W is some constant
- For each negated antecedent, input weight $w \leftarrow -W$
- Threshold weight $w_0 \leftarrow -(n-.5)W$, where n is number of non-negated antecedents

Finally, add many additional connections with near-zero weights

Liftable \leftarrow *Graspable*, \neg *Heavy*

Result of refining the network





KBANN Results

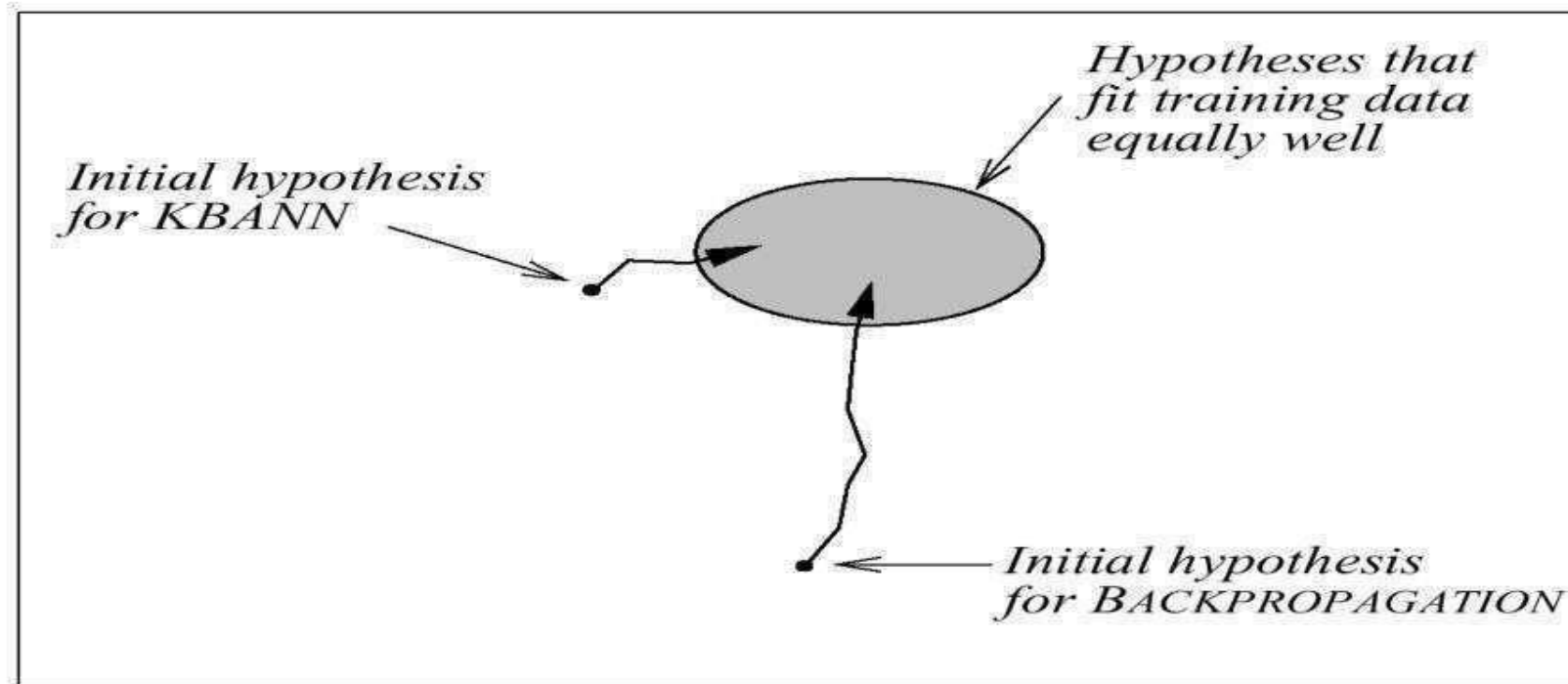
Classifying promoter regions in DNA leave one out testing:

- Backpropagation : error rate 8/106
- KBANN: 4/106

Similar improvements on other classification, control tasks.

Hypothesis space search in KBANN

Hypothesis Space





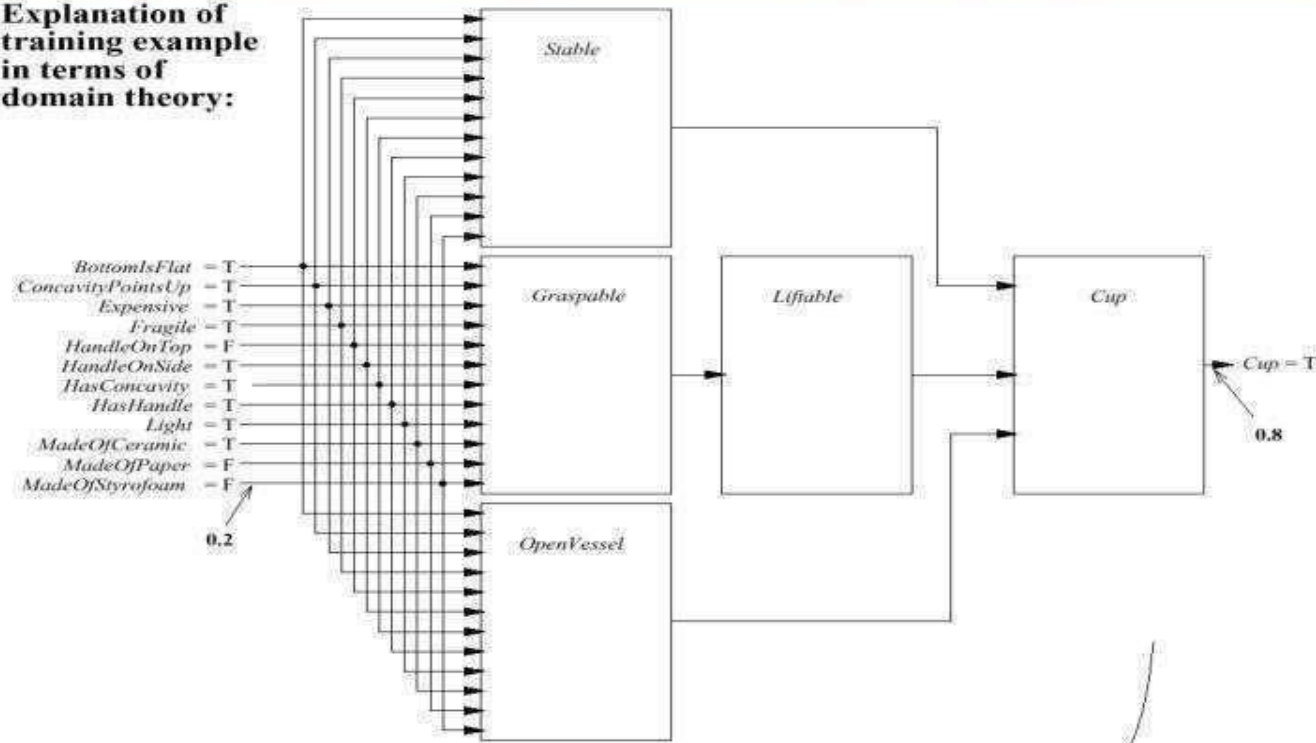
EBNN

Key idea:

- Previously learned approximate domain theory
- Domain theory represented by collection of neural networks
- Learn target function as another neural network



Explanation of training example in terms of domain theory:



Target network:





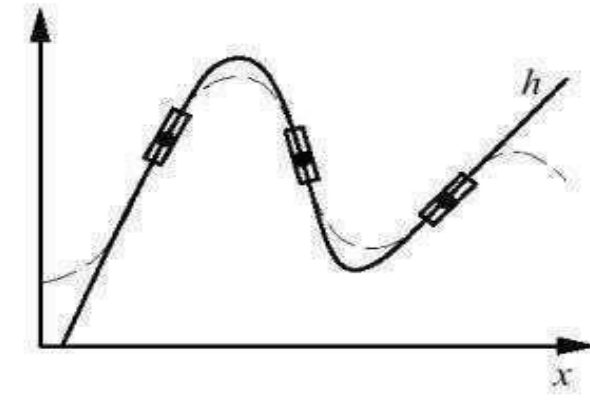
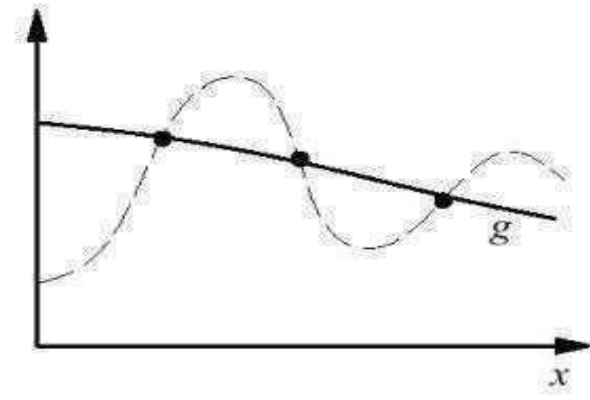
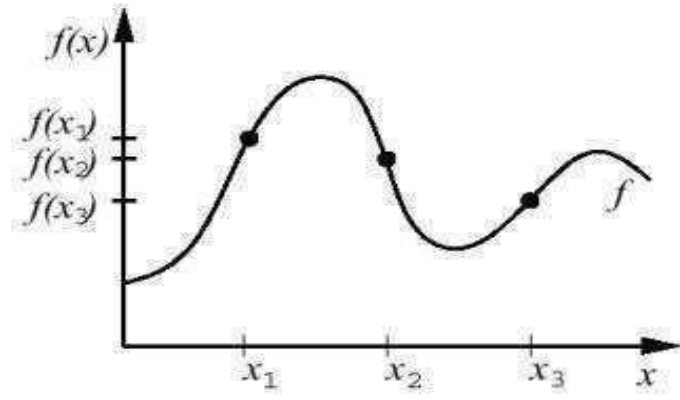
Modified Objective for Gradient Descent

$$E = \sum_i \left[(f(x_i) - \hat{f}(x_i))^2 + \mu_i \sum_j \left(\frac{\partial A(x)}{\partial x^j} - \frac{\partial \hat{f}(x)}{\partial x^j} \right)_{(x=x_i)}^2 \right]$$

where

$$\mu_i \equiv 1 - \frac{|A(x_i) - f(x_i)|}{c}$$

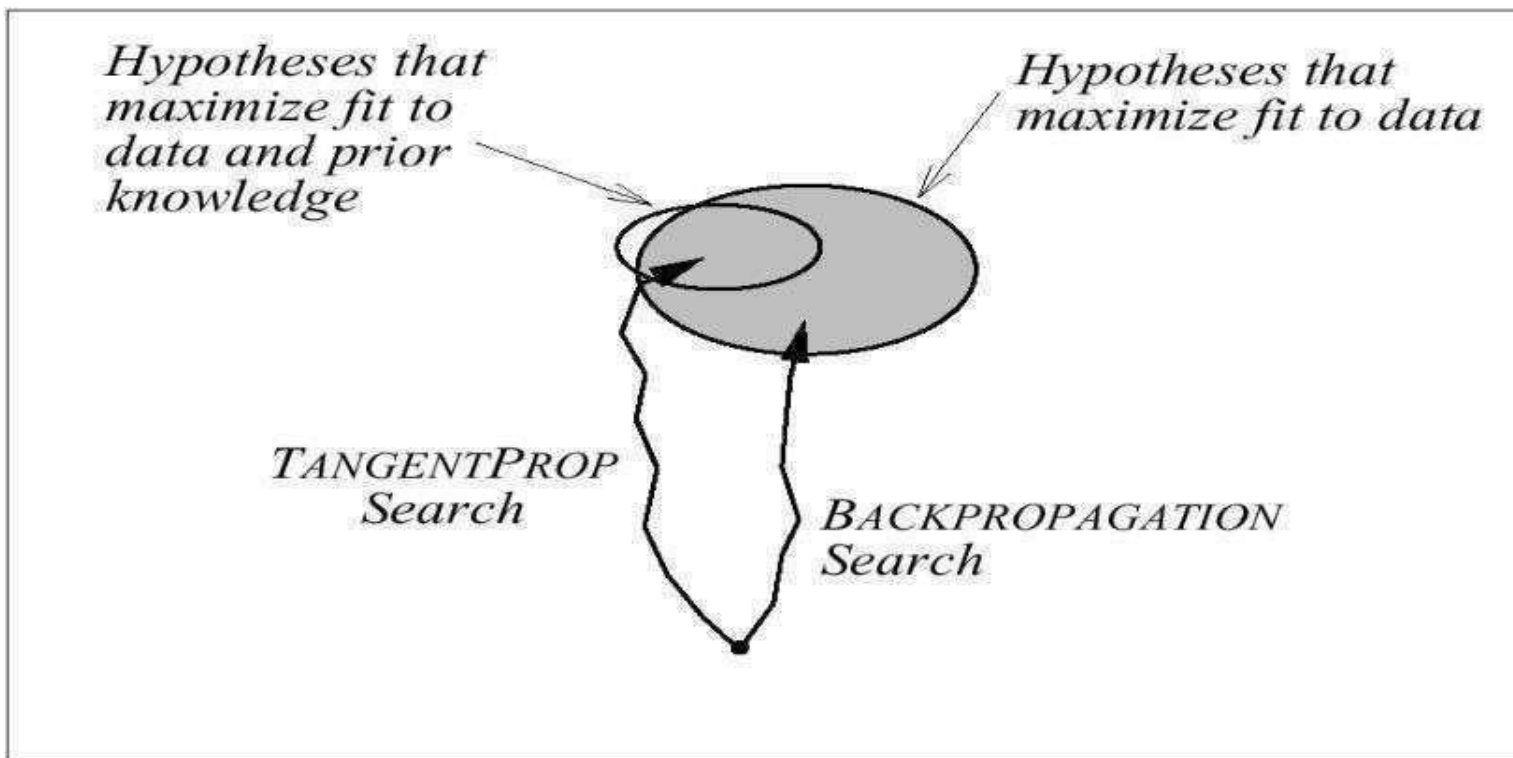
- $f(x)$ is target function
- $\hat{f}(x)$ is neural net approximation to $f(x)$
- $A(x)$ is domain theory approximation to $f(x)$



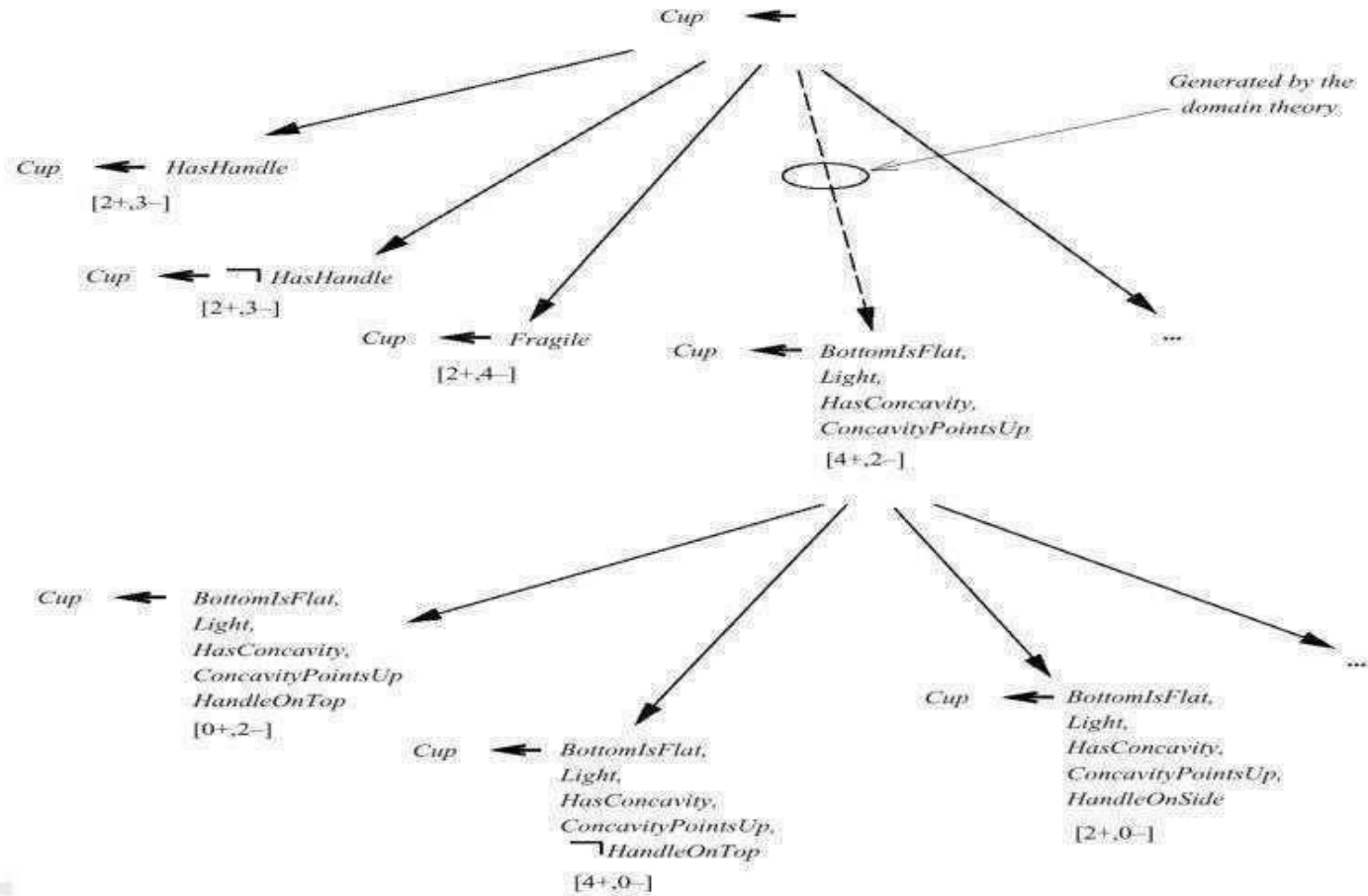


Hypothesis Space Search in EBNN

Hypothesis Space



Search in FOCL



THANKYOU