

# 23IT503:MachineLearning

**Mr.V.Anil Kumar  
Assistant Professor  
Department of IT & DS**



**NARSIMHA REDDY ENGINEERING COLLEGE**  
**UGC AUTONOMOUS INSTITUTION**

Maisammaguda (V), Kompally - 500100, Secunderabad, Telangana State, India

**UGC - Autonomous Institute**  
Accredited by **NBA** & NAAC with 'A' Grade  
Approved by **AICTE**  
Permanently affiliated to **JNTUH**

# References

## TextBooks:

1. Tom M. Mitchell, Machine Learning, India Edition 2013, McGraw Hill Education.

## Reference Books:

1. Trevor Hastie, Robert Tibshirani, Jerome Friedman, *The Elements of Statistical Learning*, 2nd edition, Springer series in statistics.
2. Ethem Alpaydın, *Introduction to machine learning*, second edition, MIT press.

# Prerequisites

For Machine Learning Course we recommend that students meet the following prerequisites:

- Basic programming skills (in Python)
- Algorithm design
- Basics of probability & statistics

# Content

- Unit-1      Introduction, Concept Learning, Decision Tree
- Unit-2      LearningArtificialNeuralNetworks-1,ArtificialNeural Networks-2, Evaluating Hypothesis,
- Unit-3      BayesianLearning,Computationallearningtheory, Instance Based Learning,
- Unit-4      GeneticAlgorithms,LearningSetsofRules, Reinforcement Learning
- Unit-5      AnalyticalLearning-1,AnalyticalLearning-2,CombiningInductiveand Analytical Learning

# UNIT-1

# MachineLearning

## Introduction

Ever since computers were invented, we have wondered whether they might be made to learn. If we could understand how to program them to learn - to improve automatically with experience - the impact would be dramatic.

- Imagine computers learning from medical records which treatments are most effective for new diseases
- Houses learning from experience to optimize energy costs based on the of their occupants.
- Personal software assistants learning the evolving interests of their to highlight especially relevant stories from the online morning newspaper

# Examples of Successful Applications of Machine Learning

- Learning to recognize spoken words
- Learning to drive an autonomous vehicle
- Learning to classify new astronomical structures
- Learning to play world-class backgammon

# Why is Machine Learning Important?

- Some tasks cannot be defined well, except by examples (e.g., recognizing people).
- Relationships and correlations can be hidden within large amounts of data. Machine Learning/Data Mining maybe able to find these relationships.
- Human designers often produce machines that do not work as well as desired in the environments in which they are used.
- The amount of knowledge available about certain tasks might be too large for explicit encoding by humans (e.g., medical diagnostic).
- Environments change over time.
- New knowledge about tasks is constantly being discovered by humans. It may be difficult to continuously re-design systems “by hand”.

# Areas of Influence for Machine Learning

- **Statistics:** How best to use samples drawn from unknown probability distributions to help decide from which distributions some new sample is drawn?
- **Brain Models:** Non-linear elements with weighted inputs (Artificial Neural Networks) have been suggested as simple models of biological neurons.
- **Adaptive Control Theory:** How to deal with controlling a process having unknown parameters that must be estimated during operation?
- **Psychology:** How to model human performance on various learning tasks?
- **Artificial Intelligence:** How to write algorithms to acquire the knowledge humans are able to acquire, at least, as well as humans?
- **Evolutionary Models:** How to model certain aspects of biological evolution to improve the performance of computer programs?

# MachineLearning: A Definition

A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E.

# Why“Learn”?

Learning is used when:

- Human expertise does not exist (navigating on Mars)
- Humans are unable to explain their expertise (speech recognition)
- Solution changes in time (routing on a computer network)
- Solution needs to be adapted to particular cases (user biometrics)

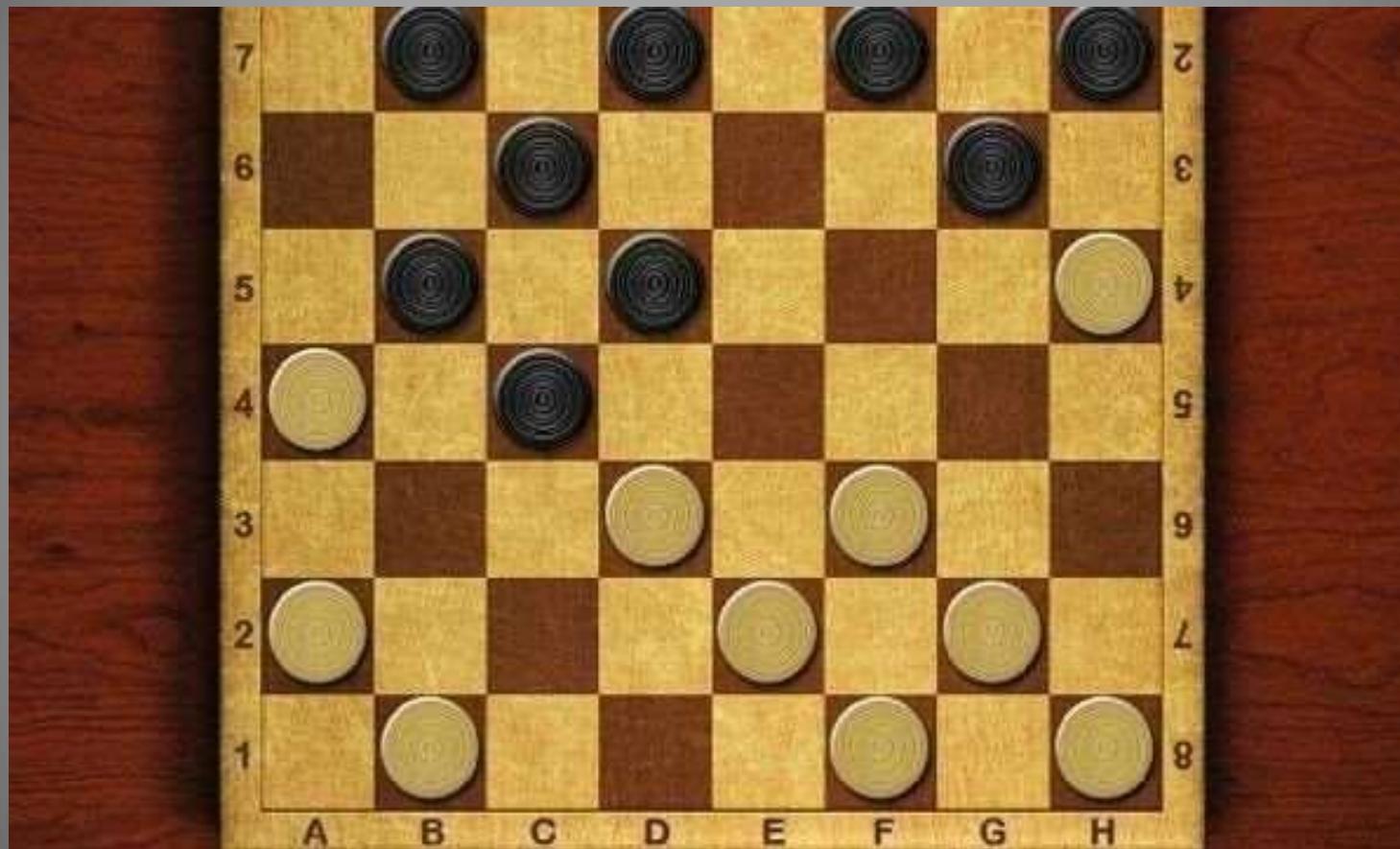
# Well-Posed Learning Problem

**Definition:** A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance on tasks in T, as measured by P, improves with experience E.

To have a well-defined learning problem, three features need to be identified:

1. The class of tasks
2. The measure of performance to be improved
3. The source of experience

# CheckersGame



Computer Science and Engineering, NRCM

# GameBasics

- Checkers is played by two players. Each player begins the game with 12 colored discs.(One set of pieces is black and the other red.) Each player places his or her pieces on the 12 dark squares closest to him or her. Black moves first. Players then alternate moves.
- The board consists of 64 squares, alternating between 32 dark and 32 light squares.
- It is positioned so that each player has a light square on the right side corner closest to him or her.
- A player wins the game when the opponent cannot make a move. In most cases, this is because all of the opponent's pieces have been captured, but it could also be because all of his pieces are blocked in.

# Rules of the Game

- Moves are allowed only on the dark squares, so pieces always move diagonally. Single pieces are always limited to forward moves (toward the opponent).
- A piece making a non-capturing move (not involving a jump) may move only one square.
- A piece making a capturing move (a jump) leaps over one of the opponent's pieces, landing in a straight diagonal line on the other side. Only one piece may be captured in a single jump; however, multiple jumps are allowed during a single turn.
- When a piece is captured, it is removed from the board.
- If a player is able to make a capture, there is no option; the jump must be made.
- If more than one capture is available, the player is free to choose whichever he or she prefers.

# Rules of the Game Cont.

- When a piece reaches the furthest row from the player who controls that piece, it is crowned and becomes a king. One of the pieces which had been captured is placed on top of the king so that it is twice as high as a single piece.
- Kings are limited to moving diagonally but may move both forward and backward. (Remember that single pieces, i.e. non-kings, are always limited to forward moves.)
- Kings may combine jumps in several directions, forward and backward, on the same turn. Single pieces may shift direction diagonally during a multiple capture return, but must always jump forward (toward the opponent).

# Well-Defined Learning Problem

*A checkers learning problem:*

- Task T: playing checkers
- Performance measure P: percent of games won against opponents
- Training experience E: playing practice games against itself

*A handwriting recognition learning problem:*

- Task T: recognizing and classifying handwritten words within images
- Performance measure P: percent of words correctly classified
- Training experience E: a database of handwritten words with given classifications

## *A robot driving learning problem:*

- Task T: driving on public four-lane highways using vision sensors
- Performance measure P: average distance travelled before an error (as judged by human overseer)
- Training experience E: a sequence of images and steering commands recorded while observing a human driver

# Designing a Learning System

1. Choosing the Training Experience
2. Choosing the Target Function
3. Choosing a Representation for the Target Function
4. Choosing a Function Approximation Algorithm
  1. Estimating training values
  2. Adjusting the weights
5. The Final Design

# 1. Choosing the Training Experience

- The first design choice is to choose the type of training experience from which the system will learn.
- The type of training experience available can have a significant impact on success or failure of the learner.

There are three attributes which impact on success or failure of the learner

1. Whether the training experience provides direct or indirect feedback regarding the choices made by the performance system.
2. The degree to which the learner controls the sequence of training examples
3. How well it represents the distribution of examples over which the final system performance must be measured.

# 1. Whether the training experience provides direct or indirect feedback regarding the choices made by the performance system ■

**For example, in checkers game:**

- In learning to play checkers, the system might learn from direct training examples consisting of individual ***Checkers board states and the correct move for each.***
- Indirect training examples consisting of the ***move sequences and final outcomes of various games played.***
- The information about the correctness of specific moves early in the game must be inferred indirectly from the fact that the game was eventually won or lost.
- Here the learner faces an additional problem of **credit assignment**, or determining the degree to which each move in the sequence deserves credit or blame for the final outcome.
- Credit assignment can be a particularly difficult problem because the game can be lost even when early moves are optimal, if these are followed later by poor moves.
- Hence, learning from direct training feedback is typically easier than learning from indirect feedback.

## *2. A second important attribute of the training experience is the degree to which the learner controls the sequence of training examples*

**For example, in checkers game:**

- The learner might depend on the teacher to select informative board states and to provide the correct move for each.
- Alternatively, the learner might itself propose board states that it finds particularly confusing and ask the teacher for the correct move.
- The learner may have complete control over both the board states and (indirect) training classifications, as it does when it learns by playing against itself with no teacher present.
- Notice in this last case the learner may choose between experimenting with novel board states that it has not yet considered, or honing its skill by playing minor variations of lines of play it currently finds most promising.

### 3. A third attribute of the training experience is how well it represents *the distribution of examples over which the final system performance must be measured* ■

Learning is most reliable when the training examples follow a distribution similar to that of future test examples.

#### **For example, in checkers game:**

- In checkers learning scenario, the performance metric  $P$  is the percent of games the system wins in the world tournament.
- If its training experience  $E$  consists only of games played against itself, there is a danger that this training experience might not be fully representative of the distribution of situations over which it will later be tested. For example, the learner might never encounter certain crucial board states that are very likely to be played by the human checkers champion.
- It is necessary to learn from a distribution of examples that is somewhat different from those on which the final system will be evaluated. Such situations are problematic because mastery of one distribution of examples will not necessarily lead to strong performance over some other distribution.

### 3. Choosing the Target Function

The next design choice is to determine exactly what type of knowledge will be learned and how this will be used by the performance program.

- Let's begin with a checkers-playing program that can generate the legal moves from any board state.
- The program needs only to learn how to choose the best move from among these legal moves. This learning task is representative of a large class of tasks for which the legal moves that define some large search space are known *a priori*, but for which the best search strategy is not known.

Given this setting where we must learn to choose among the legal moves, the most obvious choice for the type of information to be learned is a program, or function, that chooses the best move for any given board state.

1. Let ***ChooseMove*** be the target function and the notation is

***ChooseMove:B***  $\rightarrow$  ***M***

Which indicate that this function accepts as input any board from the set of legal board states ***B*** and produces as output some move from the set of legal moves ***M***.

***ChooseMove*** is a choice for the target function in checkers example, but this function will turn out to be very difficult to learn given the kind of indirect training experience available to our system

2. An alternative target function is an *evaluation function* that assigns a ***numerical score*** to any given board state

Let the target function  $V$  and the notation

$$V: B \longrightarrow R$$

Which denotes that  $V$  maps any legal board state from the set  $B$  to some real value

We intend for this target function  $V$  to assign higher scores to better board states. If the system can successfully learn such a target function  $V$ , it can easily use it to select the best move from any current board position.

Let us define the target value  $V(b)$  for an arbitrary board state  $b$  in  $B$ , as follows:

1. If  $b$  is a final board state that is won, then  $V(b) = 100$
2. If  $b$  is a final board state that is lost, then  $V(b) = -100$
3. If  $b$  is a final board state that is drawn, then  $V(b) = 0$
4. If  $b$  is not a final state in the game, then  $V(b) = V(b')$ ,

Where  $b'$  is the best final board state that can be achieved starting from  $b$  and playing optimally until the end of the game

### 3. Choosing a Representation for the Target Function

Let us choose a simple representation - for any given board state, the function will be calculated as a linear combination of the following board features:

X<sub>1</sub>: the number of black pieces on the board  
X<sub>2</sub>: the number of red pieces on the board  
X<sub>3</sub>: the number of black kings on the board  
X<sub>4</sub>: the number of red kings on the board

X<sub>5</sub>: the number of black pieces threatened by red (i.e., which can be captured on red's next turn)

X<sub>6</sub>: the number of red pieces threatened by black

$$\hat{V}(b) = w_0 + w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + w_5x_5 + w_6x_6$$

Where,

- $w_0$  through  $w_6$  are numerical coefficients, or weights, to be chosen by the learning algorithm.
- Learned values for the weights  $w_1$  through  $w_6$  will determine the relative importance of the various board features in determining the value of the board
- The weight  $w_0$  will provide an additive constant to the board value

## Partial design of a checkers learning program:

- Task T: playing checkers
- Performance measure P: percent of games won in the world tournament
- Training experience E: games played against itself
- Target function:  $V: \text{Board} \rightarrow \mathbb{R}$
- Target function representation

$$\hat{V}(b) = w_0 + w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + w_5x_5 + w_6x_6$$

The first three items above correspond to the specification of the learning task, whereas the final two items constitute design choices for the implementation of the learning program.

# 4. Choosing a Function Approximation Algorithm

- In order to learn the target function we require a set of training examples, each describing a specific board state  $b$  and the training value  $V_{\text{train}}(b)$  for  $b$ .
- Each training example is an ordered pair of the form  $(b, V_{\text{train}}(b))$ .
- For instance, the following training example describes a board state  $b$  in which black has won the game (note  $x_2 = 0$  indicates that red has no remaining pieces) and for which the target function value  $V_{\text{train}}(b)$  is therefore +100.

$$((x_1=3, x_2=0, x_3=1, x_4=0, x_5=0, x_6=0), +100)$$

# Function Approximation Procedure

1. Derive training examples from the indirect training experience available to the learner
2. Adjusts the weights  $w_i$  to best fit these training examples

[Type here]

Where,

$\hat{v}$  is the learner's current approximation to  $V$

$\text{Successor}(b)$  denotes the next board state following for which it is again the program's turn to move

## **Rule for estimating training values**

$$V_{\text{train}}(b) \leftarrow \hat{v}(\text{Successor}(b))$$

Specify the learning algorithm for choosing the weights  $w_i$  to best fit the set of training examples  $\{(b, V_{train}(b))\}$

A first step is to define what we mean by the best fit to the training data.

- One common approach is to define the best hypothesis, or set of weights, as that which minimizes the squared error  $E$  between the training values and the values predicted by the hypothesis.

$$E \equiv \sum_{(b, V_{train}(b)) \in \text{training examples}} (V_{train}(b) - \hat{V}(b))^2$$

- Several algorithms are known for finding weights of a linear function that minimize  $E$ .

One such algorithm is called the **least mean squares**, or **LMS training rule**. For each observed training example it adjusts the weights as small amount in the direction that reduces the error on this training example

**LMS weight update rule**:- For each training example

( $b, V_{train}(b)$ ) Use the current weights to calculate  $\hat{v}(b)$

For each weight  $w_i$ , update it as

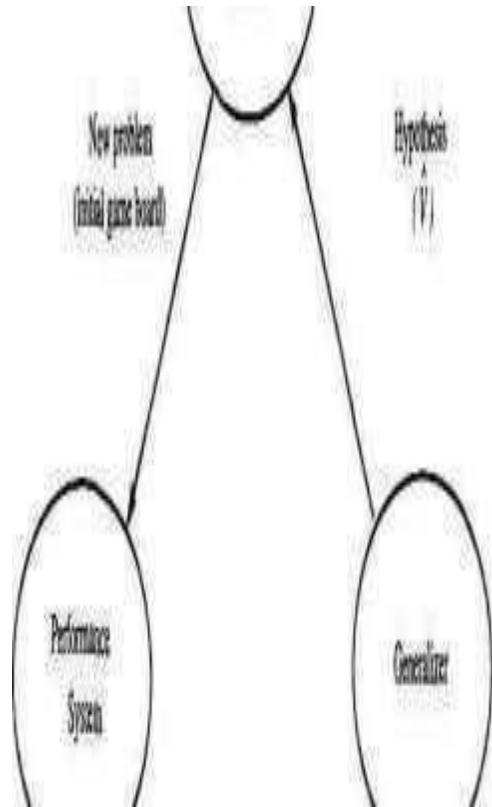
$$w_i \leftarrow w_i + \eta (V_{train}(b) - v(b)) x_i$$

# g., 0.1) that moderates the size of the weight update. Working of weight update rule

- When the error ( $V_{train}(b) - v(b)$ ) is zero, no weights are changed.
- When  $(V_{train}(b) - v(b))$  is positive (i.e., when  $v(b)$  is too low), then each weight is increased in proportion to the value of its corresponding feature. This will raise the value of  $v(b)$ , reducing the error.
- If the value of some feature  $x_i$  is zero, then its weight is not altered regardless of the error, so that the only weights updated are those whose features actually occur on the training example board.

The final design of checkers learning system can be described by four distinct program modules that represent the central components in many learning systems

## 5. The Final Design



**1. The Performance System** is the module that must solve the given performance task by using the learned target function(s).

It takes an instance of a new problem (new game) as input and produces a trace of its solution (game history) as output.

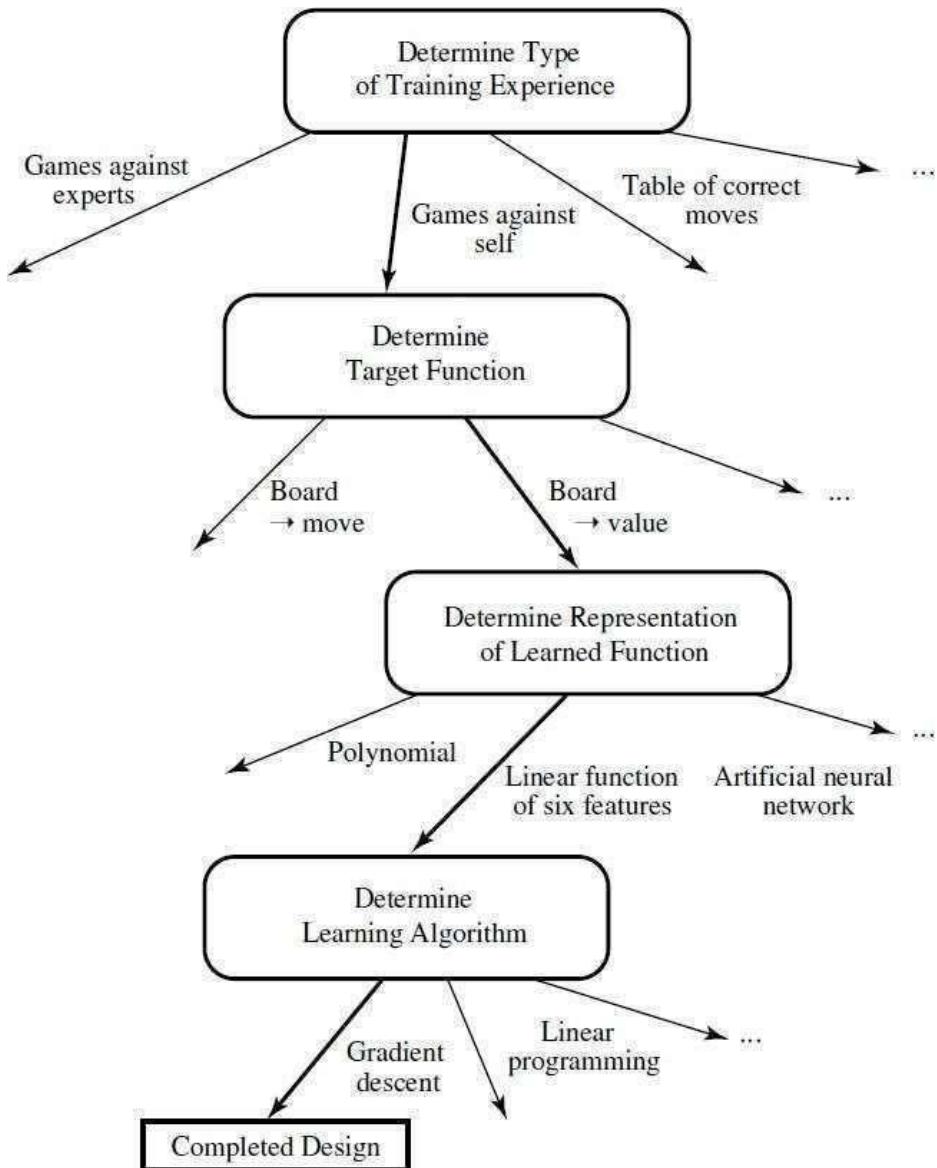
In checkers game, the strategy used by the Performance System to select its next move at each step is determined by the learned  $\hat{V}$  evaluation function. Therefore, we expect its performance to improve as this evaluation function becomes increasingly accurate.

**2. The Critic** takes as input the history or trace of the game and produces as output a set of training examples of the target function. As shown in the diagram, each training example in this case corresponds to some game state in the trace, along with an estimate  $V_{\text{train}}$  of the target function value for this example.

**3. The Generalizer** takes as input the training examples and produces an output hypothesis that is its estimate of the target function.

It generalizes from the specific training examples, hypothesizing a general function that covers these examples and other cases beyond the training examples. In our example, the Generalizer corresponds to the LMS algorithm, and the output hypothesis is the function  $\hat{V}$  described by the learned weights  $w_0, \dots, w_6$ .

**4. The Experiment Generator** takes as input the current hypothesis and outputs a new problem (i.e., initial board state) for the Performance System to explore. Its role is to pick new practice problems that will maximize the learning rate of the overall system. In our example, the Experiment Generator always proposes the same initial game board to begin a new game.



# Perspectives of Machine Learning

Perspective of machine learning involves searching very large space of possible hypothesis to determine one that

Best fits the observed data and any prior knowledge held by learner.

# Issues in Machine Learning

- What algorithms exist for learning general target functions from specific training examples? In what settings will particular algorithms converge to the desired function, given sufficient training data? Which algorithms perform best for which types of problems and representations?
- How much training data is sufficient? What general bounds can be found to relate the confidence in learned hypotheses to the amount of training experience and the character of the learner's hypothesis space?
- When and how can prior knowledge held by the learner guide the process of generalizing from examples? Can prior knowledge be helpful even when it is only approximately correct?

- What is the best strategy for choosing a useful next training experience, and how does the choice of this strategy alter the complexity of the learning problem?
- What is the best way to reduce the learning task to one or more function approximation problems? Put another way, what specific functions should the system attempt to learn? Can this process itself be automated?
- How can the learner automatically alter its representation to improve its ability to represent and learn the target function?

# Concept Learning

- Learning involves acquiring general concepts from specific training examples.  
Example: People continually learn general concepts or categories such as "bird," "car," "situations in which I should study more in order to pass the exam," etc.
- Each such concept can be viewed as describing some subset of objects or events defined over a larger set
- Alternatively, each concept can be thought to be a Boolean-valued function defined over this larger set. (Example: A function defined over all animals, whose value is true for birds and false for other animals).

**Concept learning**-Inferring a Boolean-valued function from training examples of its input and output

# A Concept Learning Task

Consider the example task of learning the target concept  
 "Day on which my friend Aldo enjoys his favorite water sport."

Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes

Table - Describes a set of example days, each represented by a set of attributes

DEPARTMENT OF COMPUTER SCIENCE, CANARA COLLEGE

## What hypothesis representation is provided to the learner?

Let's consider a simple representation in which each hypothesis consists of a conjunction of constraints on the instance attributes.

Let each hypothesis be a vector of six constraints, specifying the values of the six attributes Sky, AirTemp, Humidity, Wind, Water, and Forecast.

For each attribute, the hypothesis will be either

- Indicate by a "?" that any value is acceptable for this attribute,
- Specify a single required value (e.g., Warm) for the attribute, or
- Indicate by a " $\Phi$ " that no value is acceptable

The hypothesis that PERSON enjoys his favorite sport only on cold days with high humidity (independent of the values of the other attributes) is represented by the expression

$$(\text{?,Cold,High,?,?,?})$$

The most general hypothesis - that every day is a positive example - is represented by

$$(\text{?,?,?,?,?,?})$$

The most specific possible hypothesis - that day is a positive example - is not represented by

# Notation

The set of items over which the concept is defined is called the set of *instances*, which we denote by  $X$ .

**Example:**  $X$  is the set of all possible days, each represented by the attributes: Sky, AirTemp, Humidity, Wind, Water, and Forecast

The concept or function to be learned is called the *target concept*, which we denote by  $c$ .  $c$  can be any Boolean valued function defined over the instances  $X$ :  $c: X \rightarrow \{0, 1\}$

**Example:** The target concept corresponds to the value of the attribute *EnjoySport* (i.e.,  $c(x) = 1$  if *EnjoySport* = Yes, and  $c(x) = 0$  if *EnjoySport* = No).

- Instances for which  $c(x) = 1$  are called positive examples, or members of the target concept.
- Instances for which  $c(x) = 0$  are called negative examples, or non-members of the target concept.
- The ordered pair  $(x, c(x))$  to describe the training example consisting of the instance  $x$  and its target concept value  $c(x)$ .
- $D$  to denote the set of available training examples
- The symbol  $H$  to denote the set of all possible hypotheses that the learner may consider regarding the identity of the target concept. Each hypothesis  $h$  in  $H$  represents a Boolean-valued function defined over  $X$

$$h: X \longrightarrow \{0, 1\}$$

- The goal of the learner is to find a hypothesis such that  $h(x) = c(x)$  for all  $x$  in  $X$ .

Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes

- Given:

- Instances  $X$ : Possible days, each described by the attributes
  - Sky (with possible values *Sunny*, *Cloudy*, and *Rainy*),
  - *AirTemp* (with values *Warm* and *Cold*),
  - *Humidity* (with values *Normal* and *High*),
  - *Wind* (with values *Strong* and *Weak*),
  - *Water* (with values *Warm* and *Cool*), and
  - *Forecast* (with values *Same* and *Change*).
- Hypotheses  $H$ : Each hypothesis is described by a conjunction of constraints on the attributes *Sky*, *AirTemp*, *Humidity*, *Wind*, *Water*, and *Forecast*. The constraints may be “?” (any value is acceptable), “ $\emptyset$ ” (no value is acceptable), or a specific value.
- Target concept  $c$ :  $\text{EnjoySport} : X \rightarrow \{0, 1\}$
- Training examples  $D$ : Positive and negative examples of the target function (see Table 2.1).

- Determine:

- A hypothesis  $h$  in  $H$  such that  $h(x) = c(x)$  for all  $x$  in  $X$ .

---

TABLE The *EnjoySport* concept learning task.

Any hypothesis found to approximate the target function well over a sufficiently large set of training examples will also approximate the target function well over other unobserved examples.

## The Inductive Learning Hypothesis

# Concept learning as Search

- Concept learning can be viewed as the task of searching through a large space of hypotheses implicitly defined by the hypothesis representation.
- The goal of this search is to find the hypothesis that best fits the training examples.

**Example,** the instances  $X$  and hypotheses  $H$  in the *Enjoy Sport* learning task. The attribute *Sky* has three possible values, and *AirTemp*, *Humidity*, *Wind*, *WaterForecast* each have two possible values, the instances space  $X$  contains exactly

- $3 \cdot 2 \cdot 2 \cdot 2 \cdot 2 = 96$  Distinct instances
- $5 \cdot 4 \cdot 4 \cdot 4 \cdot 4 = 5120$  syntactically distinct hypotheses within  $H$ .

Every hypothesis containing one or more " $\Phi$ " symbols represents the empty set of instances; that is, it classifies every instance as **negative**.

- $1 + (4 \cdot 3 \cdot 3 \cdot 3 \cdot 3 \cdot 3) = 973$ . Semantically distinct hypotheses

# General-to-Specific Ordering of Hypotheses

- Consider the two hypotheses

$$h_1 = (\text{Sunny}, ?, ?, \text{Strong}, ?, ?)$$

$$h_2 = (\text{Sunny}, ?, ?, ?, ?, ?, ?)$$

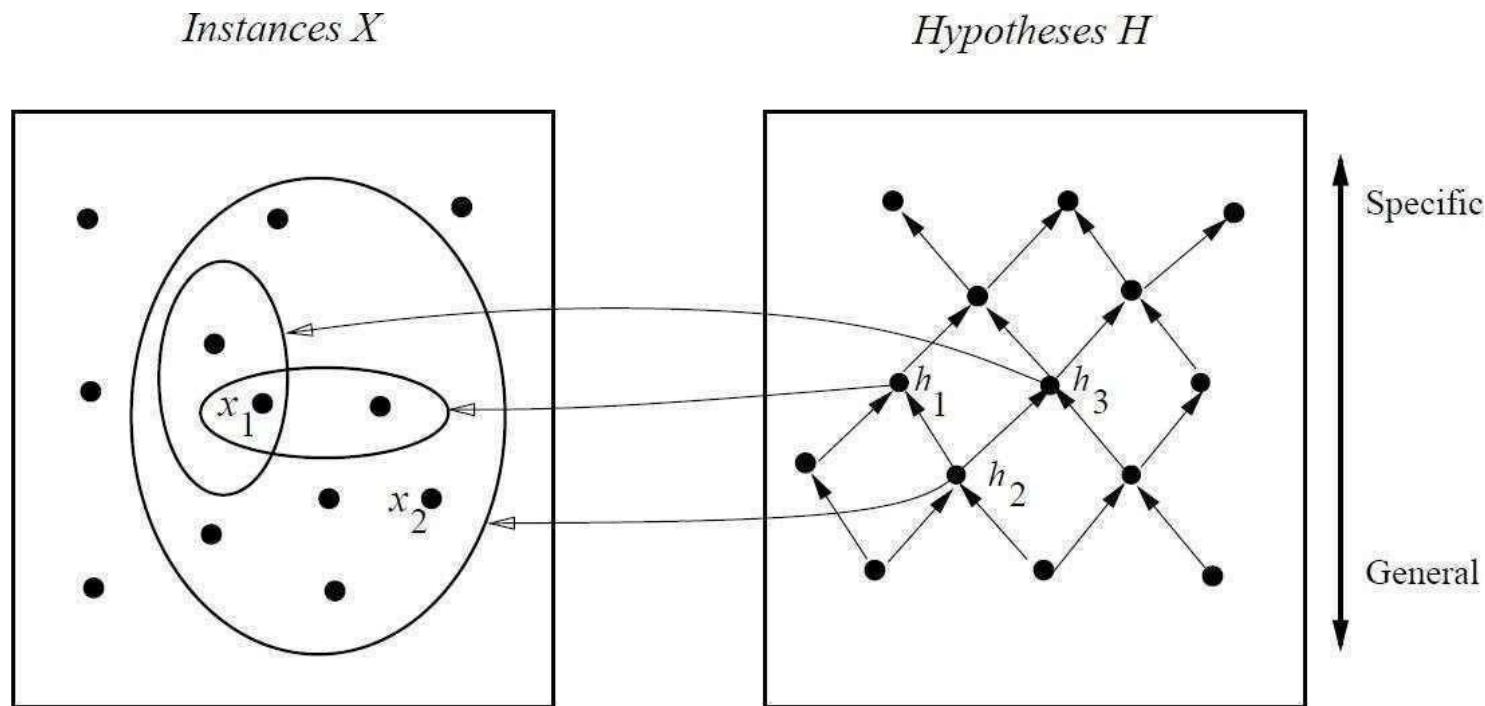
- Consider the sets of instances that are classified positive by  $h_1$  and by  $h_2$ .
- $h_2$  imposes fewer constraints on the instance, it classifies more instances as positive. So, any instance classified positive by  $h_1$  will also be classified positive by  $h_2$ . Therefore,  $h_2$  is more general than  $h_1$ .

# General-to-Specific Ordering of Hypotheses

- Given hypotheses  $h_j$  and  $h_k$ ,  $h_j$  is more-general-than-or-equal to  $h_k$  if and only if any instance that satisfies  $h_k$  also satisfies  $h_j$ .

**Definition:** Let  $h_j$  and  $h_k$  be Boolean-valued functions defined over  $X$ . Then  $h_j$  is more-general-than-or-equal to  $h_k$  (written  $h_j \geq h_k$ ) if and only if

$$(\forall x \in X)[(h_k(x) = 1) \rightarrow (h_j(x) = 1)]$$



$x_1 = \langle \text{Sunny}, \text{Warm}, \text{High}, \text{Strong}, \text{Cool}, \text{Same} \rangle$   
 $x_2 = \langle \text{Sunny}, \text{Warm}, \text{High}, \text{Light}, \text{Warm}, \text{Same} \rangle$

$h_1 = \langle \text{Sunny}, ?, ?, \text{Strong}, ?, ? \rangle$   
 $h_2 = \langle \text{Sunny}, ?, ?, ?, ?, ? \rangle$   
 $h_3 = \langle \text{Sunny}, ?, ?, ?, \text{Cool}, ? \rangle$

- In the figure, the box on the left represents the set  $X$  of all instances, the box on the right the set  $H$  of all hypotheses.
- Each hypothesis corresponds to some subset of  $X$  – the subset of instances that it classifies positive.
- The arrows connecting hypotheses represent the **more - general than** relation, with the arrow pointing toward the less general hypothesis.
- Note the subset of instances characterized by  $h_2$  subsumes the subset characterized by  $h_1$ , hence  $h_2$  is more-general than  $h_1$ .

# FIND-S: Finding a Maximally Specific Hypothesis

## FIND-S Algorithm

1. Initialize  $h$  to the most specific hypothesis in  $H$
2. For each positive training instance  $x$ 
  - For each attribute constraint  $a_i$  in  $h$ 
    - If the constraint  $a_i$  is satisfied by  $x$ 
      - Then do nothing
    - Else replace  $a_i$  in  $h$  by the next more general constraint that is satisfied by  $x$
  - 3. Output hypothesis  $h$

Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes

The first step of FIND-S to initialize  $h$  to the most specific hypothesis in  $H$   

$$h-(\emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset)$$

$x_1 = \langle \text{Sunny} \text{ Warm} \text{ Normal} \text{ Strong} \text{ Warm} \text{ Same} \rangle, +$ 

Observing the first training example, it is clear that our hypothesis is too specific. In particular, none of the " $\emptyset$ " constraints in  $h$  are satisfied by this example, so each is replaced by the next more general constraint that fits the example

 $h_1 = \langle \text{Sunny} \text{ Warm} \text{ Normal} \text{ Strong} \text{ Warm} \text{ Same} \rangle$ 

This is still very specific; it asserts that all instances are negative except for the single positive training example

 $x_2 = \langle \text{Sunny}, \text{ Warm}, \text{ High}, \text{ Strong}, \text{ Warm}, \text{ Same} \rangle, +$ 

The second training example forces the algorithm to further generalize  $h$ , this time substituting a "?" in place of any attribute value in  $h$  that is not satisfied by the new example

 $h_2 = \langle \text{Sunny} \text{ Warm? Strong} \text{ Warm} \text{ Same} \rangle$

$x3 = \langle Rainy, Cold, High, Strong, Warm, Change \rangle, -$

Upon encountering the third training the algorithm makes no change to h. The FIND-S algorithm simply ignores every negative example.

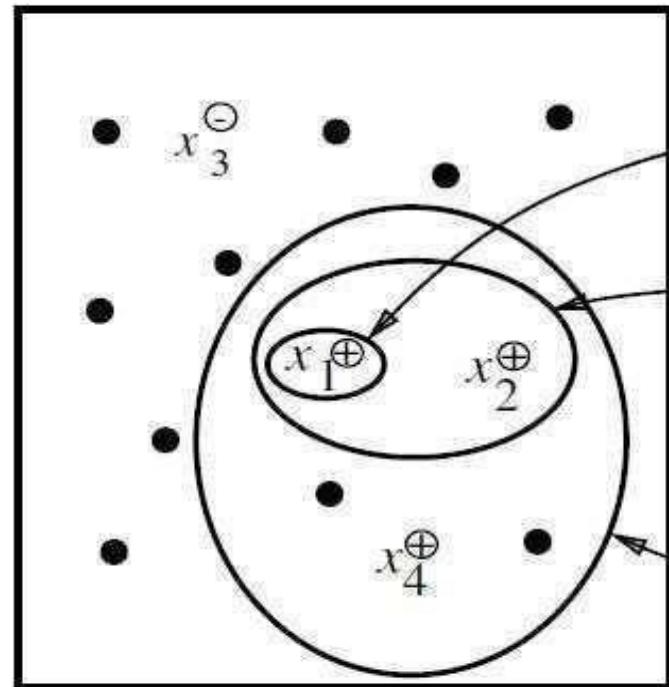
$h3 = \langle Sunny, Warm?, Strong, Warm, Same \rangle$

$x4 = \langle Sunny, Warm, High, Strong, Cool, Change \rangle, +$

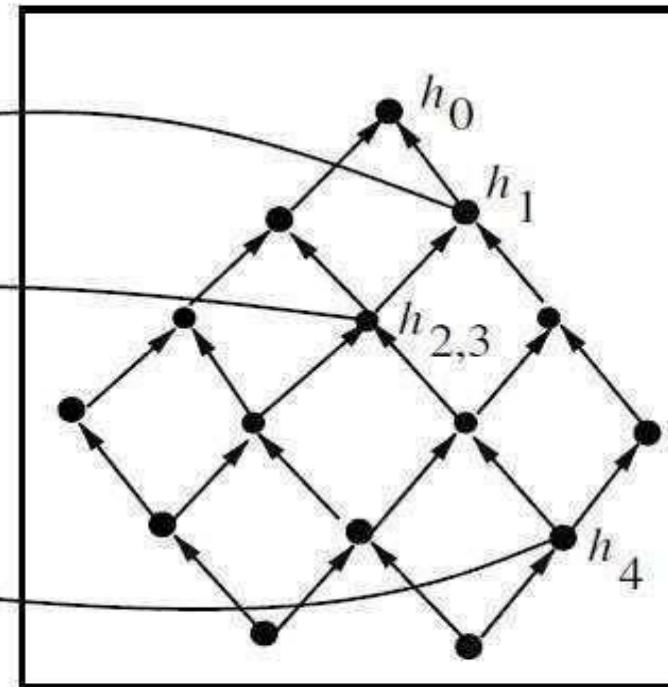
The fourth example leads to a further generalization of h

$h4 = \langle Sunny, Warm?, Strong?? \rangle$

*Instances X*



*Hypotheses H*



Specific

General

$$h_0 = \langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle$$

$$h_1 = \langle \text{Sunny Warm Normal Strong Warm Same}, + \rangle$$

$$h_2 = \langle \text{Sunny Warm ? Strong Warm Same} \rangle$$

$$h_3 = \langle \text{Sunny Warm ? Strong Warm Same} \rangle$$

$$h_4 = \langle \text{Sunny Warm ? Strong ? ?} \rangle$$

$x_1 = \langle \text{Sunny Warm Normal Strong Warm Same} \rangle, +$

$x_2 = \langle \text{Sunny Warm High Strong Warm Same} \rangle, +$

$x_3 = \langle \text{Rainy Cold High Strong Warm Change} \rangle, -$

$x_4 = \langle \text{Sunny Warm High Strong Cool Change} \rangle, +$

## The key property of the FIND-S algorithm is

- FIND-S is guaranteed to output the most specific hypothesis within  $H$  that is inconsistent with the positive training examples
- FIND-S algorithm's final hypothesis will also be consistent with the negative examples provided the correct target concept is contained in  $H$ , and provided the training examples are correct.

# Unanswered by FIND-S

1. Has the learner converged to the correct target concept?
2. Why prefer the most specific hypothesis?
3. Are the training examples consistent?
4. What if there are several maximally specific consistent hypotheses?

# Version Space and CANDIDATE-ELIMINATION Algorithm

The key idea in the CANDIDATE-ELIMINATION algorithm is to output a description of the set of all ***hypotheses consistent with their training examples***.

## Representation

- ***Definition:*** A hypothesis  $h$  is **consistent** with a set of training examples  $D$  if and only if  $h(x) = c(x)$  for each example  $(x, c(x))$  in  $D$ .

$$\text{Consistent}(h, D) \equiv (\forall \langle x, c(x) \rangle \in D) h(x) = c(x)$$

Note difference between definitions of *consistent* and *satisfies*

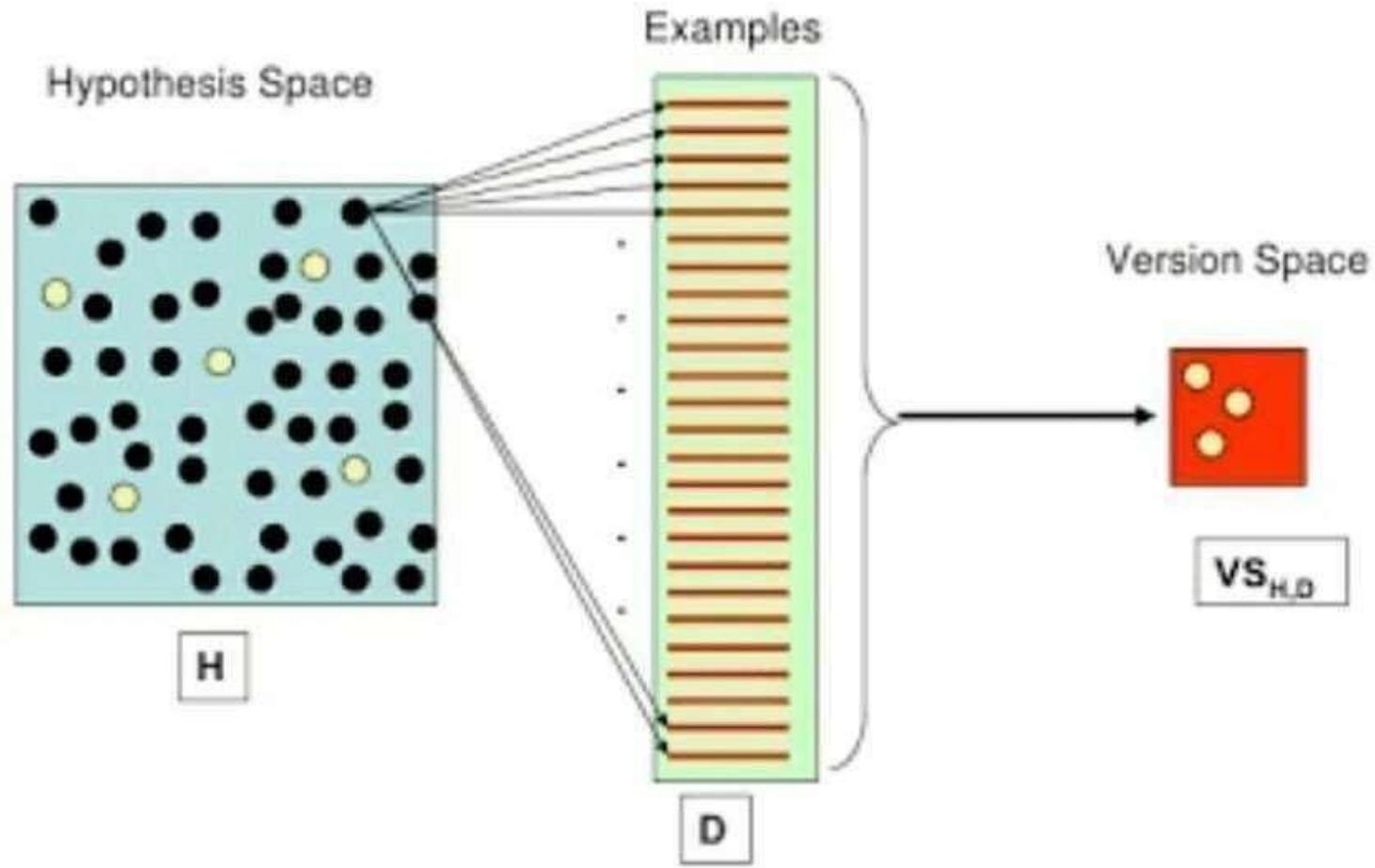
- An example is said to **satisfy** hypothesis  $h$  when  $h(x) = 1$ , regardless of whether  $x$  is a positive or negative example of the target concept.
- an example is said to **consistent** with hypothesis  $h$  iff  $h(x) = c(x)$

# Version Space

A representation of the set of all hypotheses which are consistent with D

**Definition:** The **version space**, denoted  $VS_{H,D}$  with respect to hypothesis space  $H$  and training examples D, is the subset of hypotheses from  $H$  consistent with the training examples in D

$$VS_{H,D} \equiv \{h \in H | Consistent(h, D)\}$$



# The LIST-THEN-ELIMINATE      Algorithm

The LIST-THEN-ELIMINATE algorithm first initializes the version space to contain all hypotheses in  $H$  and then eliminates any hypothesis found inconsistent with any training example.

- 
1. *VersionSpace* is a list containing every hypothesis in  $H$
  2. For each training example,  $(x, c(x))$   
remove from *VersionSpace* any hypothesis  $h$  for which  $h(x) \neq c(x)$
  3. Output the list of hypotheses in *VersionSpace*
- 

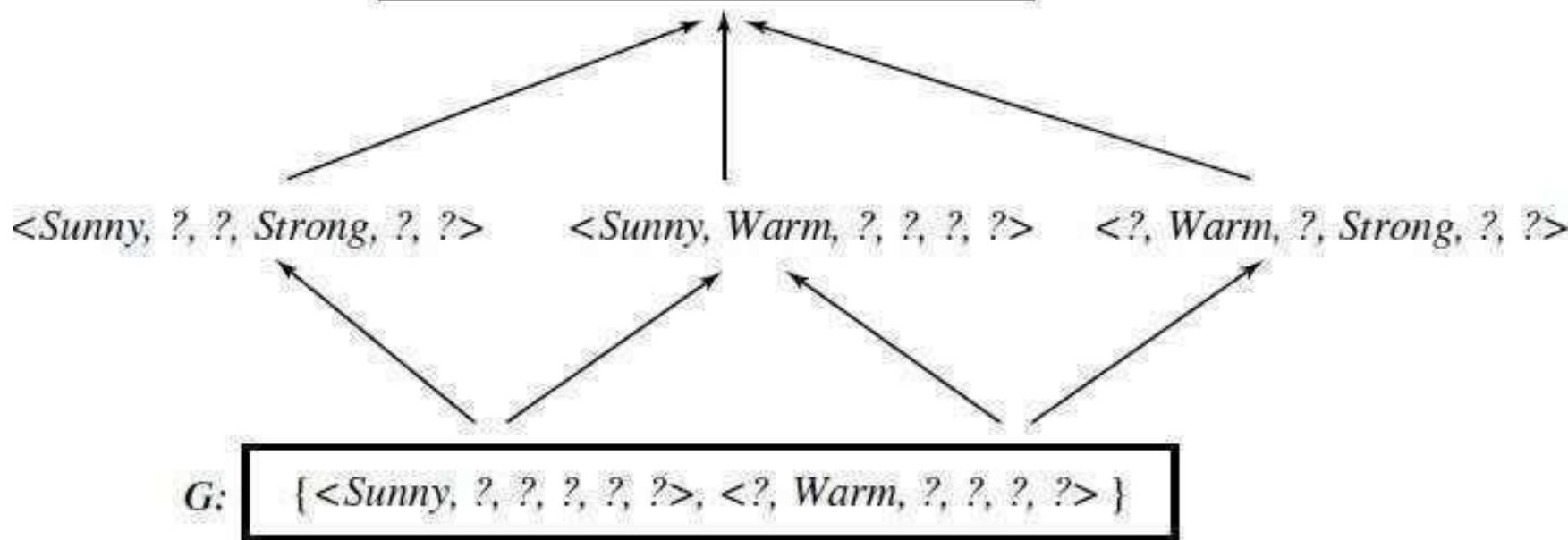
The LIST-THEN-ELIMINATE Algorithm

- **List-Then-Eliminate** works in principle, so long as version space is finite.
- However, since it requires exhaustive enumeration of all hypotheses in practice it is not feasible.

# A More Compact Representation for Version Spaces

- The version space is represented by its most general and least general members.
- These members form general and specific boundary sets that delimit the version space within the partially ordered hypothesis space.

S: {<Sunny, Warm, ?, Strong, ?, ?> }



- Aversionspacewithitsgeneral and specific boundarysets
- The version space your includes allsix hypotheses shown here, butcanberepresentedmoresimpl yby SandG.
- Arrowsindicateinstanceof the*more-general-than* relation. Thisistheversion spaceforthe *EnjoySport*conceptlearning problemandtrainingexamples describedinbelow

Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes

Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes

$$G \equiv \{g \in H | Consistent(g, D) \wedge (\neg \exists g' \in H)[(g' >_g g) \wedge Consistent(g', D)]\}$$

**Definition:** The specific boundary  $S$ , with respect to hypothesis space  $H$  and training data  $D$ , is the set of minimally general (i.e., maximally specific) members of  $H$  consistent with  $D$ .

$$S \equiv \{s \in H | Consistent(s, D) \wedge (\neg \exists s' \in H)[(s >_g s') \wedge Consistent(s', D)]\}$$

# Version Space representation theorem

**Theorem:** Let  $X$  be an arbitrary set of instances and Let  $H$  be a set of Boolean-valued hypotheses defined over  $X$ . Let  $c : X \rightarrow \{0, 1\}$  be an arbitrary target concept defined over  $X$ , and let  $D$  be an arbitrary set of training examples  $\{(x, c(x))\}$ . For all  $X, H, c$ , and  $D$  such that  $S$  and  $G$  are well defined,

$$VS_{H,D} = \{ h \in H | (\exists s \in S)(\exists g \in G)(g \geq_g h \geq_g s) \}$$

$$VSH,D = \{h \in H \mid (\exists s \in S)(\exists g \in G)(g \geq_g h \geq_g s)\}$$

ToProve:

1. Every  $h$  satisfying the right hand side of the above expression is in  $VSH,D$
2. Every member of  $VSH,D$  satisfies the right-hand side of the expression

Sketch of proof:

1. let  $g, h, s$  be arbitrary members of  $G, H, S$  respectively with  $g \geq_g h \geq_g s$

By the definition of  $S$ ,  $s$  must be satisfied by all positive examples in  $D$ . Because  $h \geq_g s$ ,  $h$  must also be satisfied by all positive examples in  $D$ .

By the definition of  $G$ ,  $g$  cannot be satisfied by any negative example in  $D$ , and because  $g \geq_g h$   $h$  cannot be satisfied by any negative example in  $D$ . Because  $h$  is satisfied by all positive examples in  $D$  and by no negative examples in  $D$ ,  $h$  is consistent with  $D$ , and therefore  $h$  is a member of  $VSH,D$

2. It can be proven by assuming some  $h$  in  $VSH,D$  that does not satisfy the right-hand side of the expression, then showing that this leads to an inconsistency

# The CANDIDATE-ELIMINATION Learning Algorithm

The CANDIDATE-ELIMINATION algorithm computes the *version space* containing all hypotheses from  $H$  that are consistent with an observed sequence of training examples.

Initialize  $G$  to the set of maximally general hypotheses in  $H$   
 Initialize  $S$  to the set of maximally specific hypotheses in  $H$   
 Foreach training example  $d$ , do

- If  $d$  is a positive example
  - Remove from  $G$  any hypothesis inconsistent with  $d$
  - For each hypothesis in  $S$  that is not consistent with  $d$ 
    - Removes from  $S$
    - Add to  $S$  all minimal generalizations of  $s$  such that
      - $h$  is consistent with  $d$ , and some member of  $G$  is more general than  $h$
    - Remove from  $S$  any hypothesis that is more general than another hypothesis in  $S$
- If  $d$  is a negative example
  - Remove from  $S$  any hypothesis inconsistent with  $d$
  - For each hypothesis in  $G$  that is not consistent with  $d$ 
    - Remove  $g$  from  $G$
    - Add to  $G$  all minimal specializations of  $g$  such that
      - $h$  is consistent with  $d$ , and some member of  $S$  is more specific than  $h$
    - Remove from  $G$  any hypothesis that is less general than another hypothesis in  $G$

# An Illustrative Example

The boundary sets are first initialized to  $G_0$  and  $S_0$ , the most general and most specific hypotheses in  $H$ .

**$S_0$**

$$\langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle$$

**$G_0$**

$$\langle ?, ?, ?, ?, ?, ? \rangle$$

For training exampled,

$\langle \text{Sunny}, \text{Warm}, \text{Normal}, \text{Strong}, \text{Warm}, \text{Same} \rangle +$

**S<sub>0</sub>**

$\langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle$



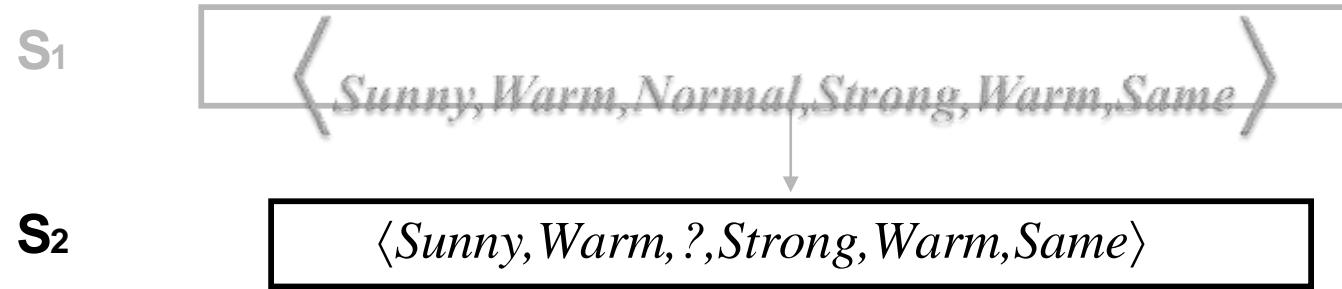
**S<sub>1</sub>**

$\langle \text{Sunny}, \text{Warm}, \text{Normal}, \text{Strong}, \text{Warm}, \text{Same} \rangle$

$\langle ?, ?, ?, ?, ?, ? \rangle$

For training example,

$\langle \text{Sunny}, \text{Warm}, \text{High}, \text{Strong}, \text{Warm}, \text{Same} \rangle +$



$\langle ?, ?, ?, ?, ?, ? \rangle$

For training example,

$\langle \text{Rainy}, \text{Cold}, \text{High}, \text{Strong}, \text{Warm}, \text{Change} \rangle -$

$S_2, S_3$

$\langle \text{Sunny}, \text{Warm}, ?, \text{Strong}, \text{Warm}, \text{Same} \rangle$

**G<sub>3</sub>**

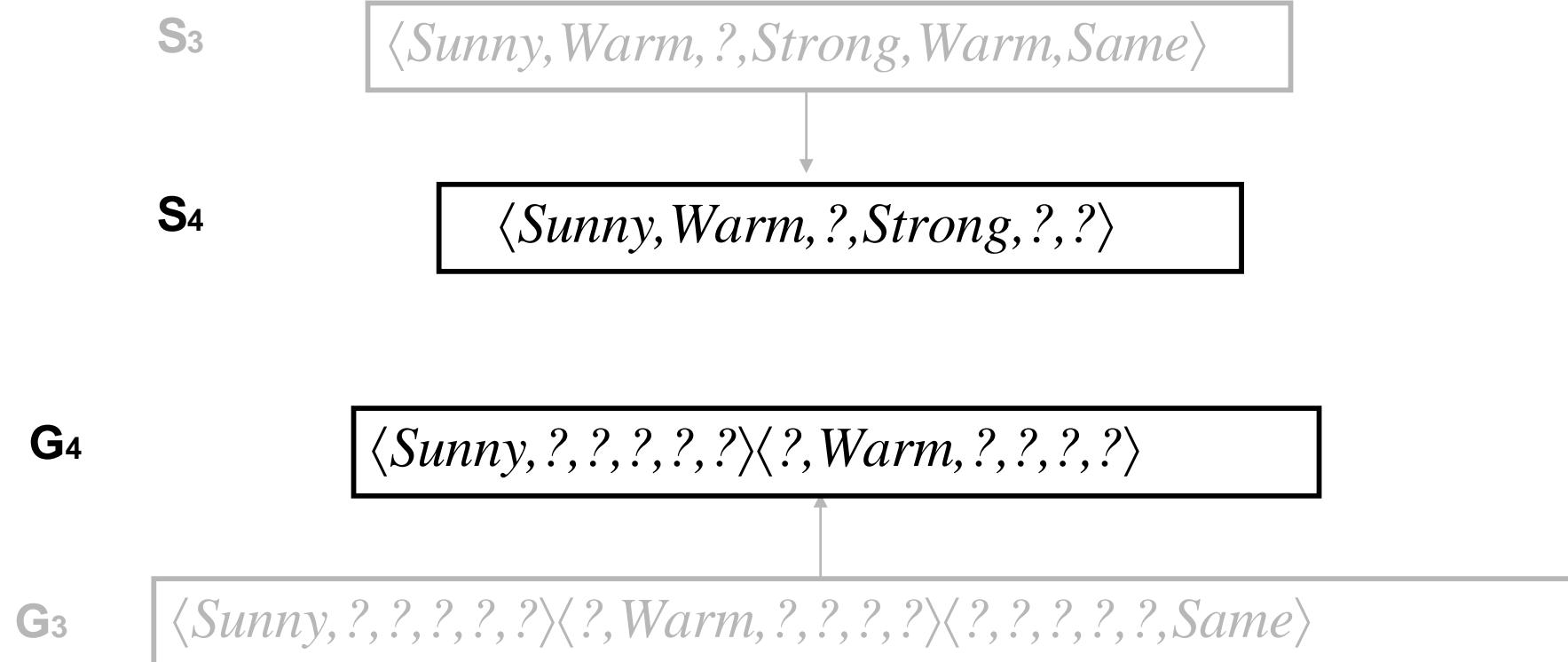
$\langle \text{Sunny}, ?, ?, ?, ?, ?, ? \rangle \langle ?, \text{Warm}, ?, ?, ?, ?, ? \rangle \langle ?, ?, ?, ?, ?, ?, \text{Same} \rangle$

**G<sub>2</sub>**

$\langle ?, ?, ?, ?, ?, ?, ? \rangle$

For training exampled,

$\langle \text{Sunny}, \text{Warm}, \text{High}, \text{Strong}, \text{CoolChange} \rangle +$



**S<sub>4</sub>**
 $\{ <\text{Sunny}, \text{Warm}, ?, \text{Strong}, ?, ?> \}$ 
 $<\text{Sunny}, ?, ?, \text{Strong}, ?, ?>$ 
 $<\text{Sunny}, \text{Warm}, ?, ?, ?, ?>$ 
 $<?, \text{Warm}, ?, \text{Strong}, ?, ?>$ 
**G<sub>4</sub>**
 $\{ <\text{Sunny}, ?, ?, ?, ?, ?>, <?, \text{Warm}, ?, ?, ?, ?> \}$ 

The final version space for the *Enjoy Sport* concept learning problem and

# Inductive Bias

The fundamental questions for inductive inference

- What if the target concept is not contained in the hypothesis space?
- Can we avoid this difficulty by using a hypothesis space that includes every possible hypothesis?
- How does the size of this hypothesis space influence the ability of the algorithm to generalize to unobserved instances?
- How does the size of the hypothesis space influence the number of training examples that must be observed?

# Effect of incomplete hypothesis space

Preceding algorithms work if target function is in  $H$

Will generally not work if target function not in  $H$

Consider following examples which represent target function  
“sky=sunny or sky=cloudy”:

$\langle \text{Sunny Warm Normal Strong Cool Change} \rangle$	Y
$\langle \text{Cloudy Warm Normal Strong Cool Change} \rangle$	Y
$\langle \text{Rainy Warm Normal Strong Cool Change} \rangle$	N

If apply Candidate Elimination algorithm as before, end up with empty Version Space After first two training example

$$S = \langle ? \text{Warm Normal Strong Cool Change} \rangle$$

New hypothesis is overly general and it covers the third negative training example! Our

# An Unbiased Learner

## Incomplete hypothesis space

- If  $c$  not in  $H$ , then consider generalizing representation of  $H$  to contain  $c$
- The size of the instances space  $X$  of days described by the six available attributes is 96. The number of distinct subsets that can be defined over a set  $X$  containing  $|X|$  elements (i.e., the size of the powerset of  $X$ ) is  $2^{|X|}$
- Recall that there are 96 instances in *Enjoy Sport*; hence there are  $2^{96}$  possible hypotheses in full space  $H$
- Can do this by using full propositional calculus with AND, OR, NOT
- Hence  $H$  defined only by conjunctions of attributes is biased (containing only 973  $h'$ s)

- Let us reformulate the *Enjoy sport* learning task in an unbiased way by defining a new hypothesis space  $H'$  that can represent every subset of instances; that is, let  $H'$  correspond to the powerset of  $X$ .
- One way to define such an  $H'$  is to allow arbitrary disjunctions, conjunctions, and negations of our earlier hypotheses.

For instance, the target concept "*Sky=Sunny* or *Sky=Cloudy*" could then be described as

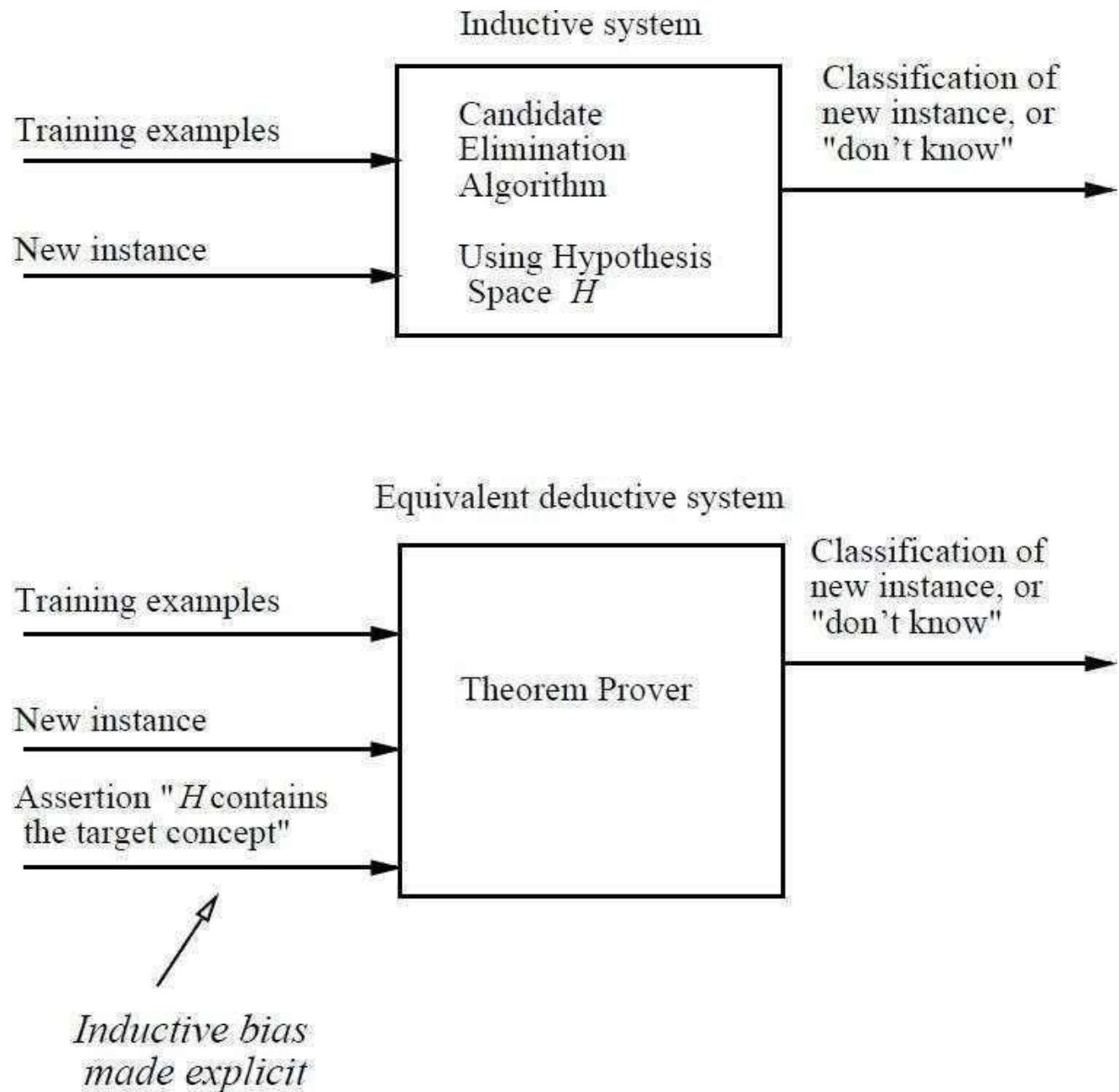
$$(\textit{Sunny}, ?, ?, ?, ?, ?) \vee (\textit{Cloudy}, ?, ?, ?, ?, ?)$$

## ***Definition:***

Consider a concept learning algorithm  $L$  for the set of instances  $X$ .

- Let  $c$  be an arbitrary concept defined over  $X$
- Let  $D_c = \{(x, c(x))\}$  be an arbitrary set of training examples of  $c$ .
- Let  $L(x_i, D_c)$  denote the classification assigned to the instance  $x_i$  by  $L$  after training on the data  $D_c$ .
- The inductive bias of  $L$  is any minimal set of assertions  $B$  such that for any target concept  $c$  and corresponding training examples  $D_c$

$$(\forall \langle x_i \in X \rangle [(B \wedge D_c \wedge x_i) \models L(x_i, D_c)])$$



Modelling inductive systems by equivalent deductive stems.

The input-output behavior of the CANDIDATE-ELIMINATION algorithm using a hypothesis space  $H$  is identical to that of a deductive theorem prover utilizing the assertion " $H$  contains the target concept." This assertion is therefore called the **inductive bias** of the CANDIDATE-ELIMINATION algorithm.

.

characterizing inductive systems

by their inductive bias allows modelling them by their equivalent deductive systems. This provides a way to compare inductive systems according to their policies for generalizing beyond the observed training data

# DECISION TREE LEARNING

# DECISION TREE REPRESENTATION

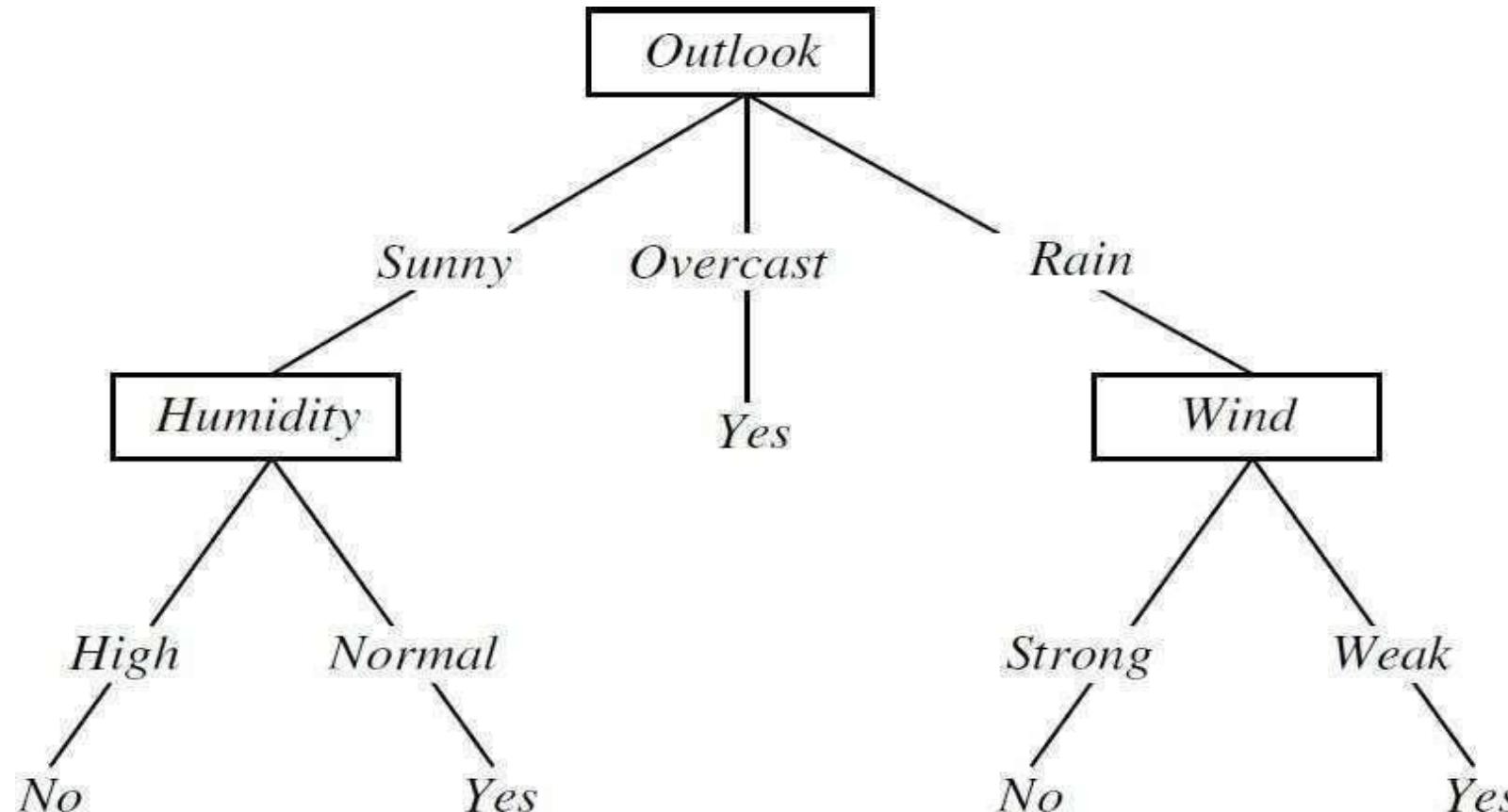


FIGURE: A decision tree for the concept **PlayTennis**. An example is classified by sorting it through the tree to the appropriate leaf node, then returning the classification associated with this leaf

- Decision trees classify instances by sorting them down the tree from the root to some leaf node, which provides the classification of the instance.
- Each node in the tree specifies a test of some attribute of the instance, and each branch descending from that node corresponds to one of the possible values for this attribute.
- An instance is classified by starting at the root node of the tree, testing the attribute specified by this node, then moving down the tree branch corresponding to the value of the attribute in the given example. This process is then repeated for the subtree rooted at the new node.

- Decision trees represent a disjunction of conjunctions of constraints on the attribute values of instances.
- Each path from the tree root to a leaf corresponds to a conjunction of attribute tests, and the tree itself to a disjunction of these conjunctions

For example,

The decision tree shown in above figure corresponds to the expression

(Outlook = Sunny  $\wedge$  Humidity = Normal)

(Outlook = Overcast)

(Outlook = Rain  $\wedge$  Wind = Weak)

# APPROPRIATE PROBLEMS FOR DECISION TREE LEARNING

Decision tree learning is generally best suited to problems with the following characteristics:

1. *Instances are represented by attribute-value pairs* – Instances are described by a fixed set of attributes and their values
2. *The target function has discrete output values* – The decision tree assigns a Boolean classification (e.g., yes or no) to each example. Decision tree methods easily extend to learning functions with more than two possible output values.
3. *Disjunctive descriptions may be required*

4. ***The training data may contain errors***—Decision tree learning methods are robust to errors, both errors in classifications of the training examples and errors in the attribute values that describe these examples.
  
5. ***The training data may contain missing attribute values***—Decision tree methods can be used even when some training examples have unknown values
  - Decision tree learning has been applied to problems such as learning to classify *medical patients by their disease, equipment malfunctions by their cause, and loan applicants by their likelihood of defaulting on payments*.
  
  - Such problems, in which the task is to classify examples into one of a discrete set of possible categories, are often referred to as ***classification problems***.

# THE BASIC DECISION TREE ALGORITHM

## LEARNING



- Most algorithms that have been developed for learning decision trees are variations on a core algorithm that employs a top-down, greedy search through the space of possible decision trees. This approach is exemplified by the ID3 algorithm and its successor C4.5

# What is the ID3 algorithm?

- ID3 stands for Iterative Dichotomiser 3
- ID3 is a precursor to the C4.5 Algorithm.
- The ID3 algorithm was invented by Ross Quinlan in 1975
- Used to generate a decision tree from a given dataset by employing a top-down, greedy search, to test each attribute at every node of the tree.
- The resulting tree is used to classify future samples.

# ID3algorithm

***ID3(Examples,Target\_attribute,Attributes)***

Examples are the training examples. Target\_attribute is the attribute whose value is to be predicted by the tree. Attributes is a list of other attributes that may be tested by the learned decision tree. Returns a decision tree that correctly classifies the given Examples.

- Create a Root node for the tree
- If all Examples are positive, Return the single-node tree Root, with label = +
- If all Examples are negative, Return the single-node tree Root, with label = -
- If Attributes is empty, Return the single-node tree Root, with label = most common value of Target\_attribute in Examples

- OtherwiseBegin
  - $A \leftarrow$  the attribute from Attributes that best\* classifies Examples
  - The decision attribute for Root  $\leftarrow A$
  - For each possible value,  $v_i$ , of A,
    - Add a new tree branch below Root, corresponding to the test  $A = v_i$
    - Let  $Examples_{v_i}$ , be the subset of Examples that have value  $v_i$  for A
    - If  $Examples_{v_i}$  is empty
      - Then below this new branch add a leaf node with label = most common value of Target\_attribute in Examples
      - Else below this new branch add the subtree  $ID3(Examples_{v_i}, \text{Targe\_tattribute}, \text{Attributes} - \{A\})$
- End
- Return Root

\*The best attribute is the one with highest information gain

# Which Attribute Is the Best Classifier?

- The central choice in the ID3 algorithm is selecting which attribute to test at each node in the tree.
- A statistical property called ***information gain*** that measures how well a given attribute separates the training examples according to their target classification.
- ID3 uses ***information gain*** measure to select among the candidate attributes at each step while growing the tree.

# ENTROPY MEASURES HOMOGENEITY OF EXAMPLES

- To define information gain, we begin by defining a measure called entropy.  
*Entropy measures the impurity of a collection of examples.*
- Given a collection S, containing positive and negative examples of some target concept, the entropy of S relative to this Boolean classification is

$$\text{Entropy } (S) \equiv -p_{\oplus} \log_2 p_{\oplus} - p_{\ominus} \log_2 p_{\ominus}$$

Where,

$p_{+}$  is the proportion of positive examples in S

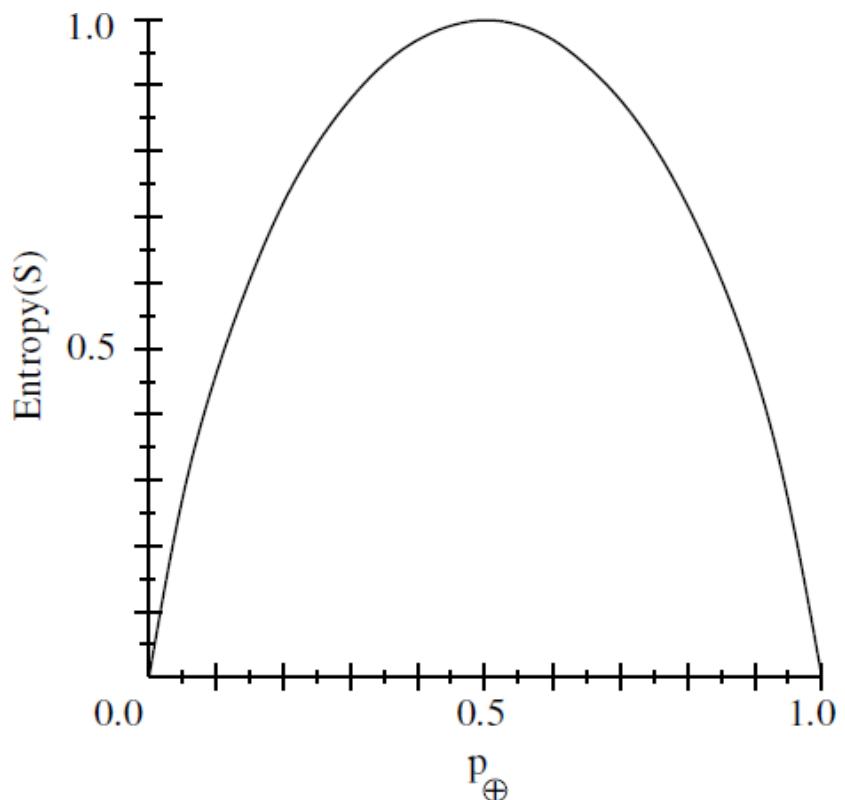
$p_{-}$  is the proportion of negative examples in S.

## Example:Entropy

- Suppose S is a collection of 14 examples of some boolean concept, including 9positive and 5 negative examples. Then the entropy of S relative to this boolean classification is

$$\begin{aligned} \text{Entropy}([9+, 5-]) &= -(9/14) \log_2(9/14) - (5/14) \log_2(5/14) \\ &= 0.940 \end{aligned}$$

- The entropy is 0 if all members of S belong to the same class
- The entropy is 1 when the collection contains an equal number of positive and negative examples
- If the collection contains unequal numbers of positive and negative examples, the entropy is between 0 and 1



**FIGURE** The entropy function relative to a boolean classification,  
as the proportion,  $p_{\oplus}$ , of positive examples varies between 0 and 1.

# INFORMATION GAIN MEASURES THE EXPECTED REDUCTION IN ENTROPY

- *Information gain*, is the expected reduction in entropy caused by partitioning the examples according to this attribute.
- The information gain,  $\text{Gain}(S, A)$  of an attribute  $A$ , relative to a collection of examples  $S$ , is defined as

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

## Example:Informationgain

Let,  $Values(Wind) = \{Weak, Strong\}$

$$S = [9+, 5-]$$

$$S_{Weak} = [6+, 2-]$$

$$S_{Strong} = [3+, 3-]$$

Informationgainofattribute  $Wind$ :

$$\begin{aligned}
 Gain(S, Wind) &= Entropy(S) - \frac{8}{14}Entropy(S_{Weak}) - \frac{6}{14}Entropy(S_{Strong}) \\
 &= 0.94 - (8/14)*0.811 - (6/14)*1.00 \\
 &= 0.048
 \end{aligned}$$

## An Illustrative Example

- To illustrate the operation of ID3, consider the learning task represented by the training examples of below table.
- Here the target attribute ***PlayTennis***, which can have values ***yes*** or ***no*** for different days.
- Consider the first step through the algorithm, in which the top most node of the decision tree is created.

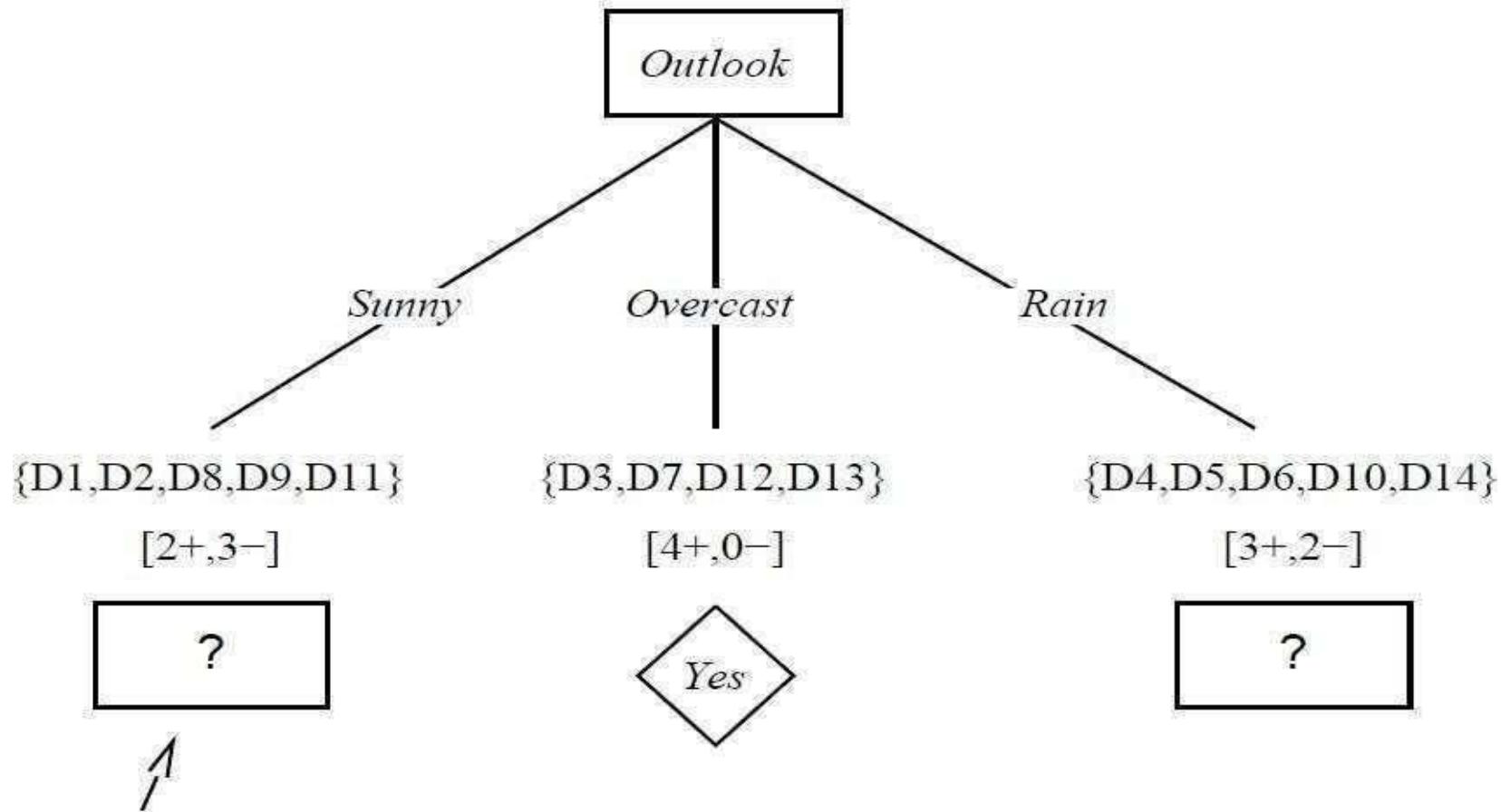
Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

The information gain values for all four attributes are

- Gain(S, Outlook) = 0.246
- Gain(S, Humidity) = 0.151
- Gain(S, Wind) = 0.048
- Gain(S, Temperature) = 0.029
- According to the information gain measure, the **Outlook** attribute provides the best prediction of the target attribute, **PlayTennis**, over the training examples. Therefore, **Outlook** is selected as the decision attribute for the root node, and branches are created below the root for each of its possible values i.e.,

{D1, D2, ..., D14}

[9+,5-]



*Which attribute should be tested here?*

$$S_{sunny} = \{D1, D2, D8, D9, D11\}$$

$$Gain(S_{sunny}, \text{Humidity}) = .970 - (3/5)0.0 - (2/5)0.0 = .970$$

$$Gain(S_{sunny}, \text{Temperature}) = .970 - (2/5)0.0 - (2/5)1.0 - (1/5)0.0 = .570$$

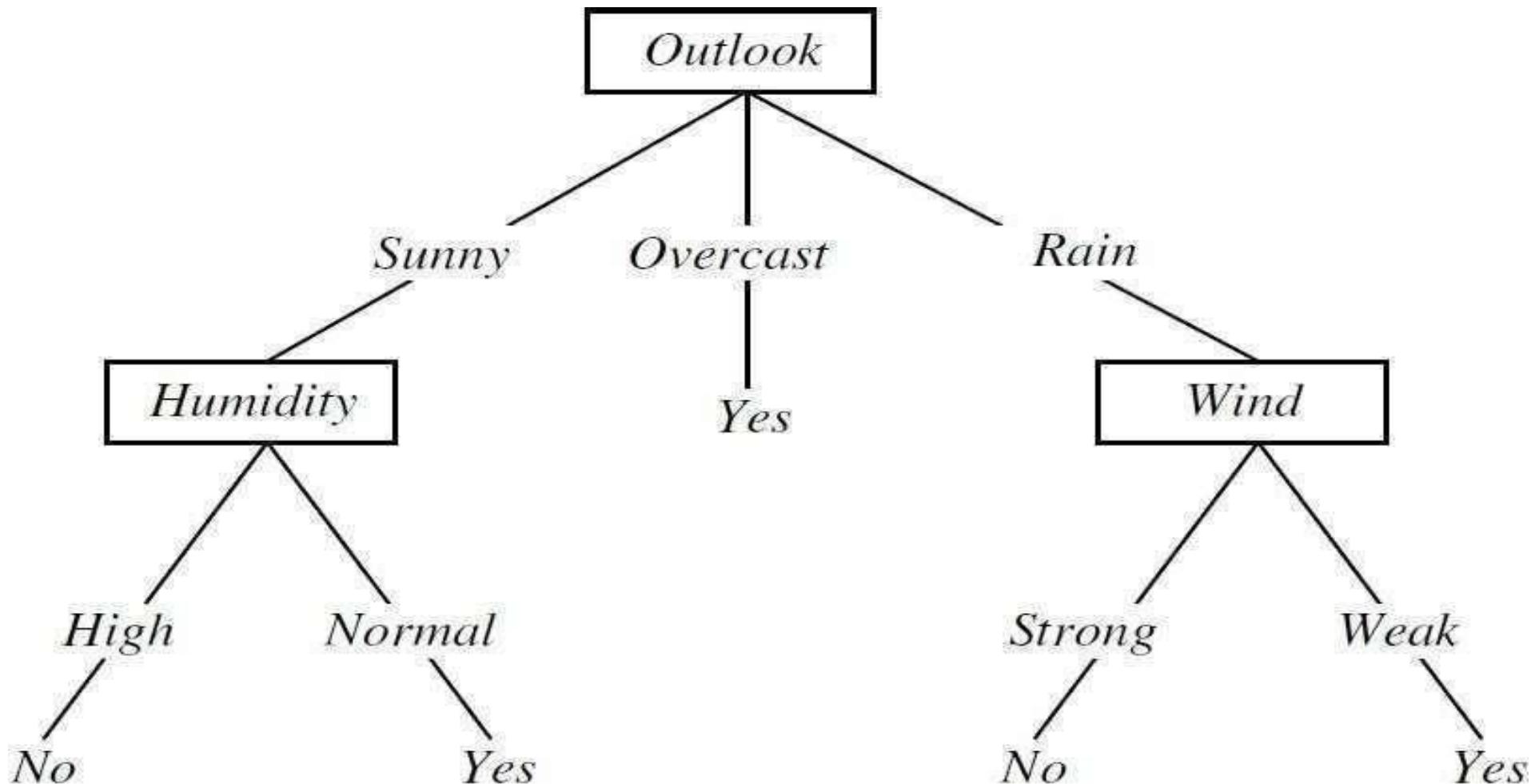
$$Gain(S_{sunny}, \text{Wind}) = .970 - (2/5)1.0 - (3/5)0.918 = .019$$

$$S_{Rain} = \{D4, D5, D6, D10, D14\}$$

$$Gain(S_{Rain}, \text{Humidity}) = 0.970 - (2/5)1.0 - (3/5)0.917 = 0.019$$

$$Gain(S_{Rain}, \text{Temperature}) = 0.970 - (0/5)0.0 - (3/5)0.918 - (2/5)1.0 = 0.019$$

$$Gain(S_{Rain}, \text{Wind}) = 0.970 - (3/5)0.0 - (2/5)0.0 = 0.970$$

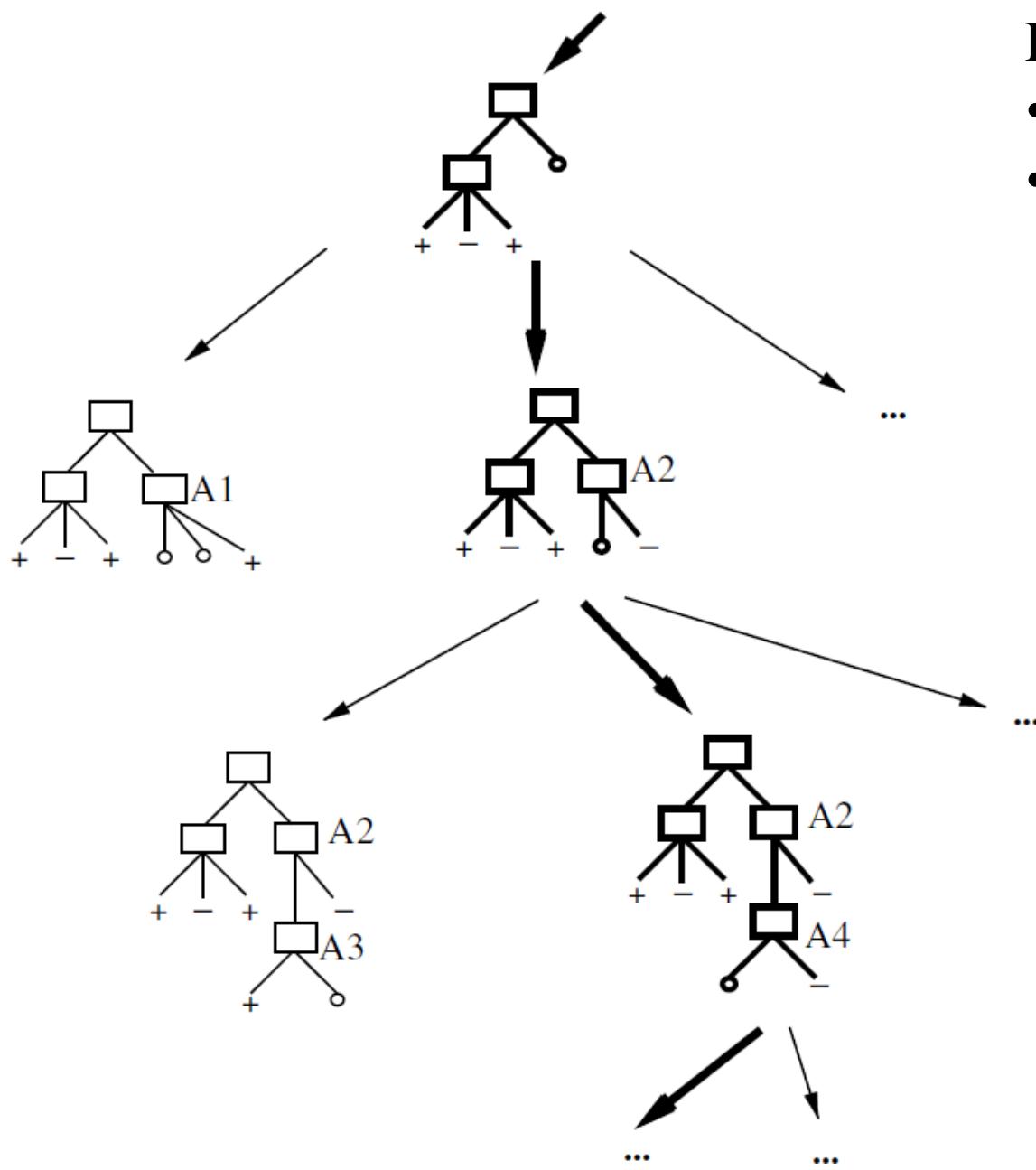


# HYPOTHESISSPACESEARCHINDECISIONTREE EARNING

- ID3 can be characterized as searching a space of hypotheses for one that fits the training examples.
- The hypothesis space searched by ID3 is the set of possible decision trees.
- ID3 performs a *simple-to complex, hill-climbing search* through this hypothesis space, beginning with the empty tree, then considering progressively more elaborate hypotheses in search of a decision tree that correctly classifies the training data

## Figure:

- Hypothesis space search by ID3.
- ID3 searches through the space of possible decision trees from simplest to increasingly complex, guided by the information gain heuristic



1. ID3's hypothesis space of all deision trees is a *complete* space of finite discrete-valued functions, relative to the available attributes. Because every finite discrete-valued function can be represented by some decision tree
  - ID3 avoids one of the major risks of methods that *search in incomplete hypothesis spaces* : that the hypothesis space might not contain the target function.

2. ID3 maintains **only a single current hypothesis** as it searches through the space of decision trees.

For example, with the earlier versions space candidate elimination method, which maintains the set of ***all*** hypotheses consistent with the available training examples.

By determining only a single hypothesis, ID3 loses the capabilities that follow from explicitly representing all consistent hypotheses.

For example, it does not have the ability to determine how many alternative decision trees are consistent with the available training data, or to pose new instance queries that optimally resolve among these competing hypotheses

**3. ID3** in its pure form performs *no backtracking in its search*. Once it selects an attribute to test at a particular level in the tree, it never backtracks to reconsider this choice.

- In the case of **ID3**, a locally optimal solution corresponds to the decision tree it selects along the single search path it explores. However, this locally optimal solution may be less desirable than trees that would have been encountered along a different branch of the search.

**4. ID3 uses all training examples at each step** in the search to make statistically based decisions regarding how to refine its current hypothesis.

- One advantage of using statistical properties of all the examples is that the resulting search is much *less sensitive to errors* in individual training examples.
- **ID3** can be easily extended to handle noisy training data by modifying its termination criterion to accept hypotheses that imperfectly fit the training data.

# INDUCTIVEBIASINDECISIONTREELEARNING

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 ***breadthfirst*** 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 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 ***apreference 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 language bias.***

## Whichtypeofinductivebiasispreferredinordertogeneralizebeyondthetrainingdata.apreference 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.

## Argument opposed:

- 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 sets of hypotheses that one can define but understood by fewer learner.
- The size of a hypothesis 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

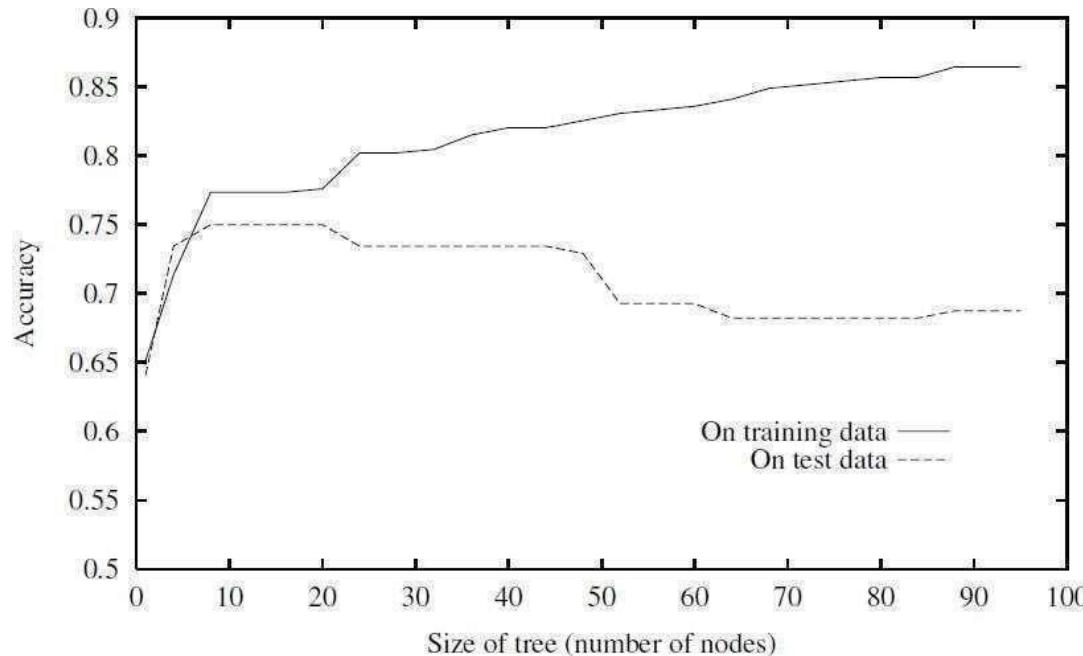
## 4. Handling Training Examples with Missing Attribute Values

## 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 lines show the accuracy of the decision tree over the training examples. The broken lines show 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.

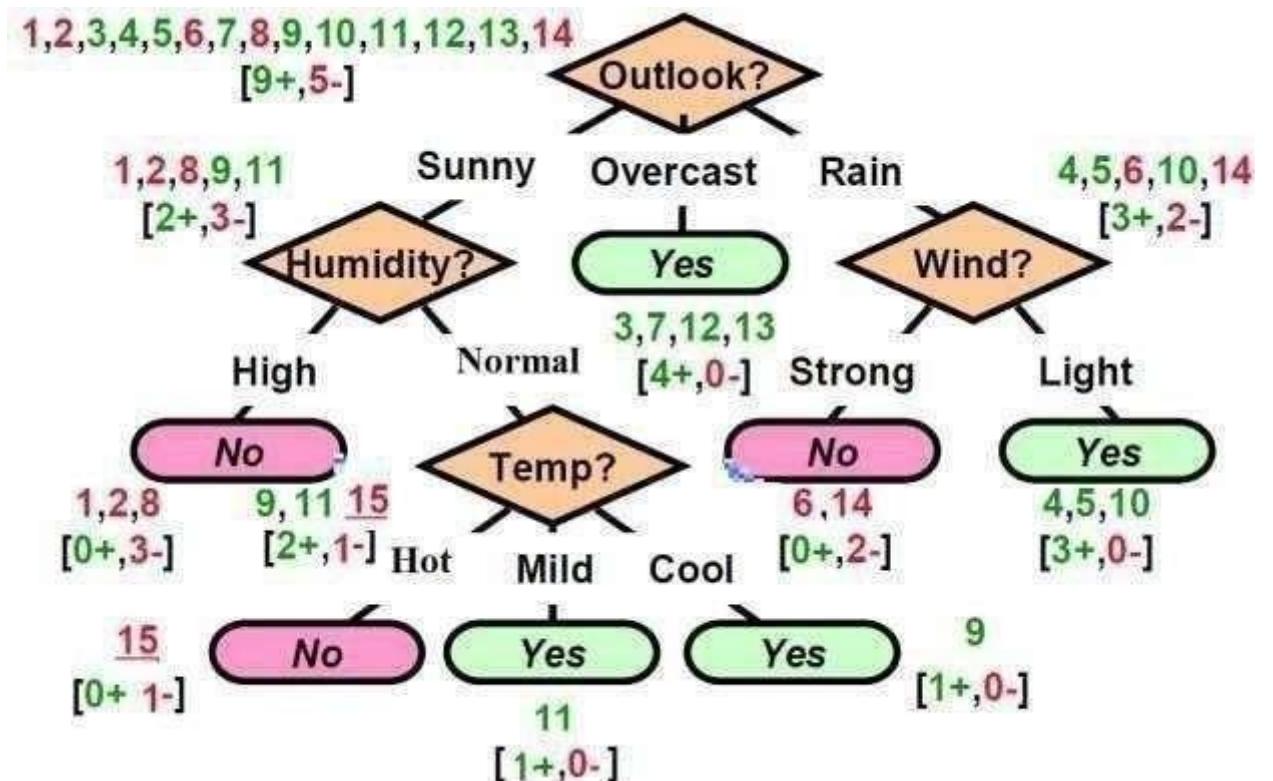
**How can it be possible for tree to fit the training examples better than 'h', but for it to perform more poorly over subsequent examples?**

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: <Sunny, Hot, Normal, Strong, ->

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



## I 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 (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

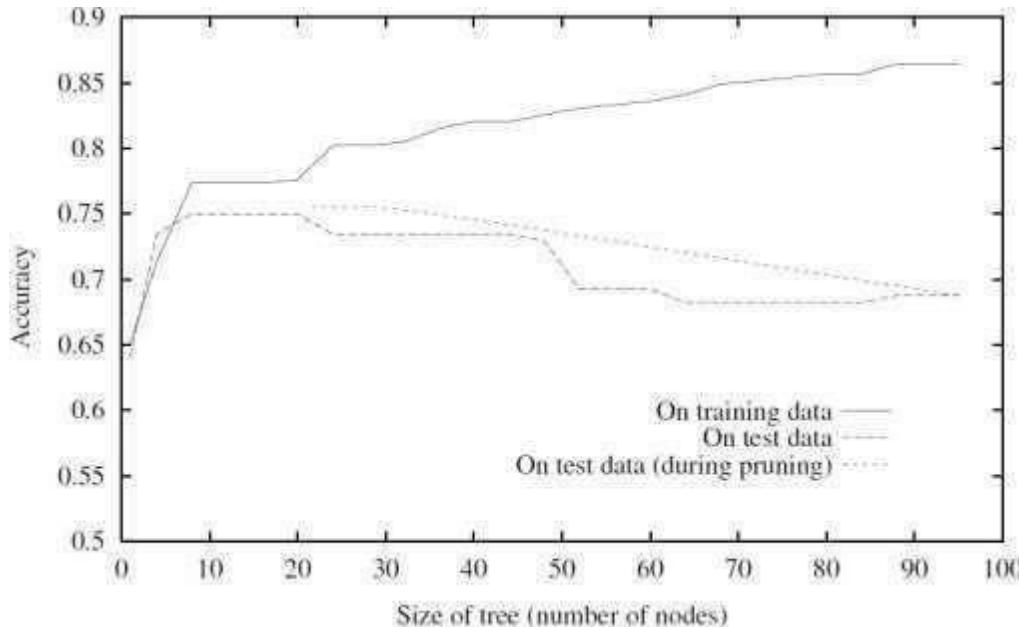
### I *MDL*–

*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 ( $error_{dh}$ ) of hypothesis, based on its observed error ( $error_{sh}$ ) over a sample of data.
- ' In general, the problem of estimating confidence intervals is approached by identifying the parameter to be estimated ( $error_{dh}$ ) 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.
- ' An 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 characterise 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  $\text{On}(x,y)$  AND  
              One subgoal to be solved is  $\text{On}(y,z)$   
**THEN**   Solve the subgoal  $\text{On}(y,z)$  before  $\text{On}(x,y)$

## 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

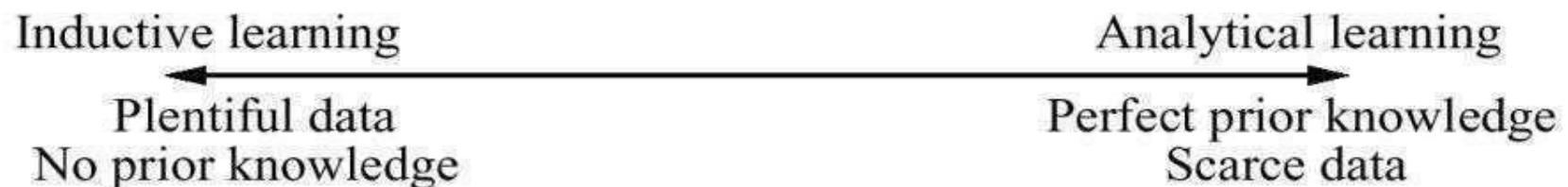
## Inductive learning

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

## Analytical learning

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

# What We Would Like



## **General purpose learning method:**

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



### Domain theory:

Cup ← Stable, Liftable, Open Vessel

Stable ← BottomIsFlat

Liftable ← 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	✓	✓		✓	✓			✓

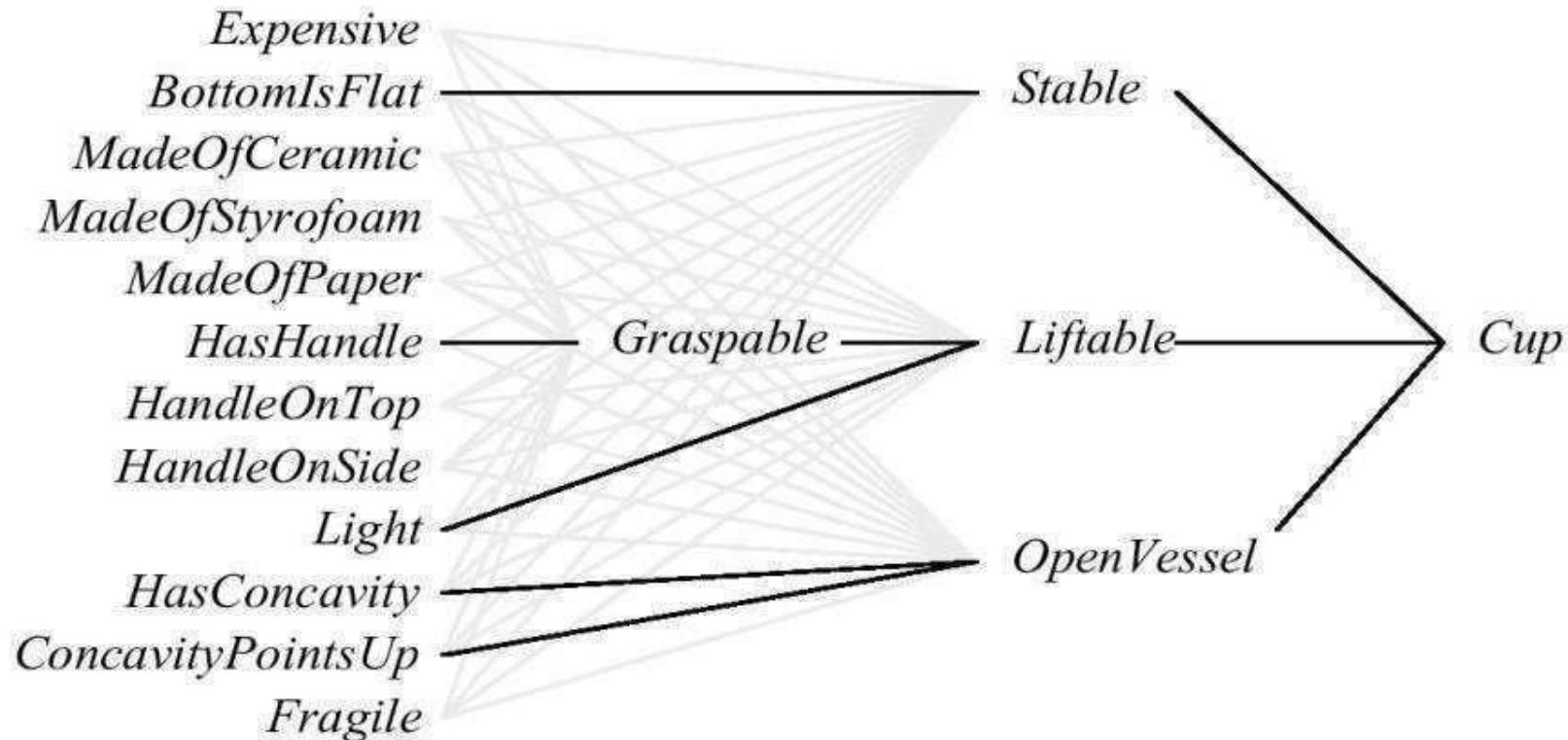


# KBANN

- 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)

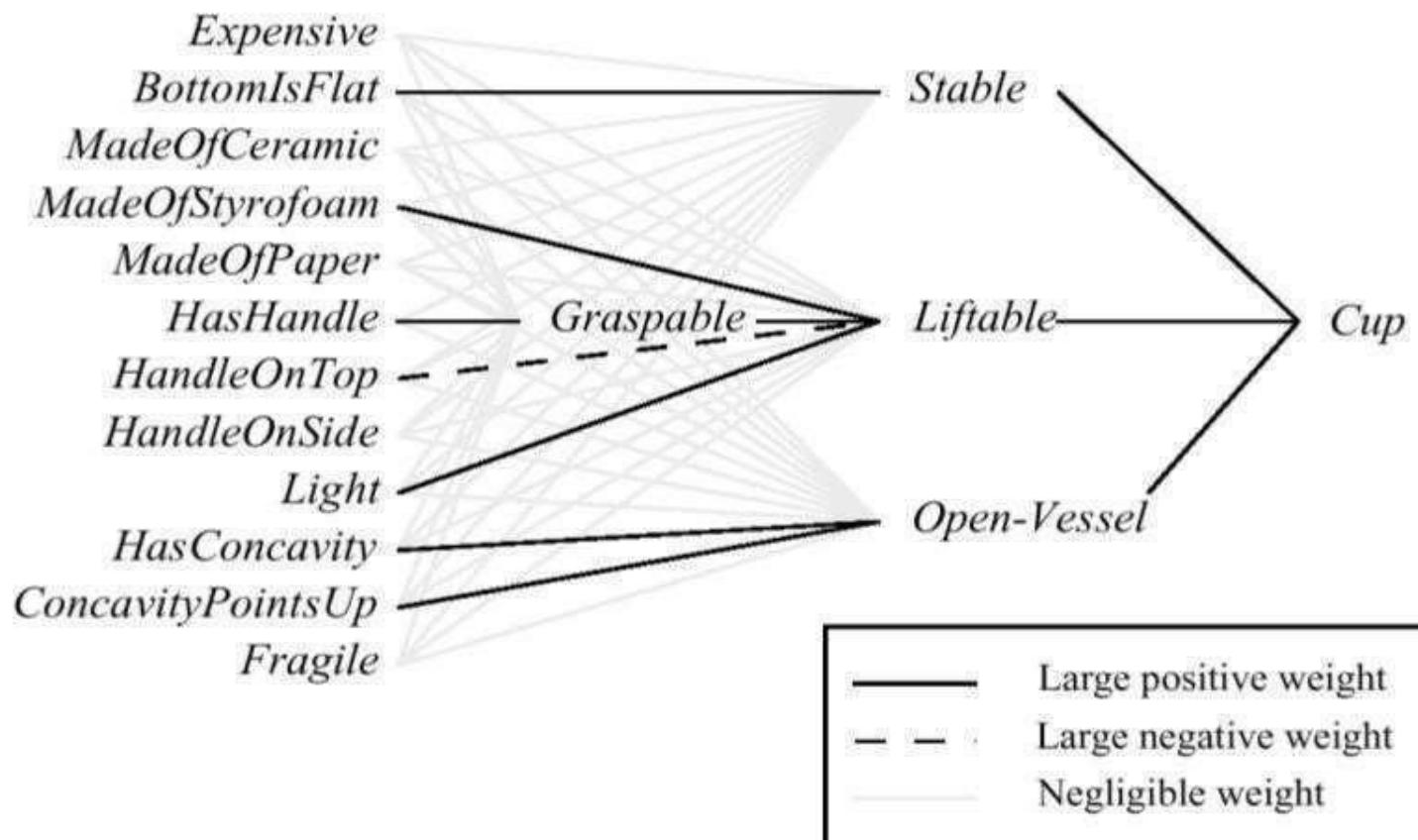
- Connect unit inputs to corresponding clause antecedents
- For each non-negated antecedent, corresponding input weight  $w \leftarrow W$ , where  $W$  is some constant
- For each negated antecedent, input weight  $w \leftarrow -W$
- Threshold weight  $w_0 \leftarrow -(n-.5)W$ , where  $n$  is number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

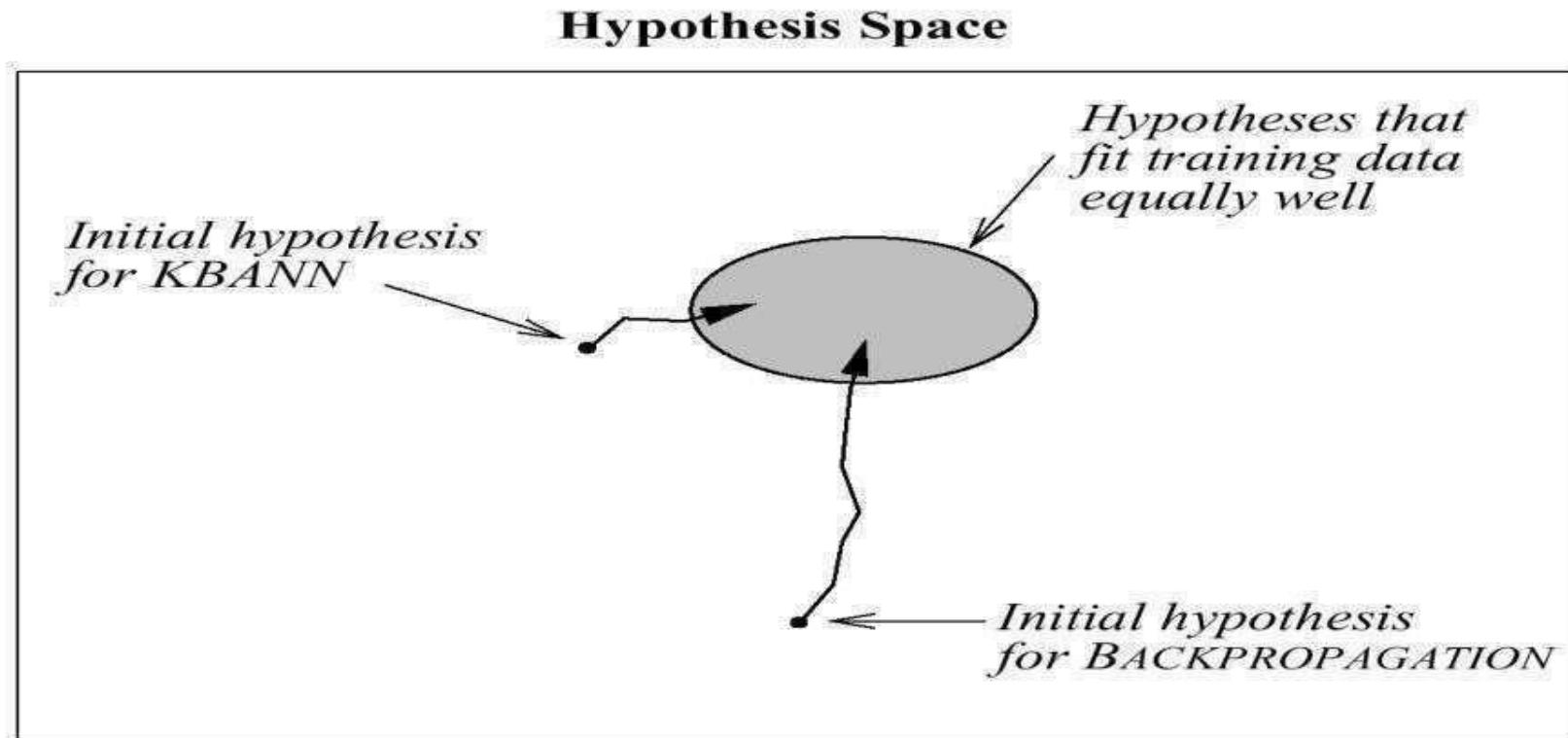
Classifying promoter regions in DNA leave one out testing:

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

Similar improvements on other classification, control tasks.



# Hypothesis space search in KBANN





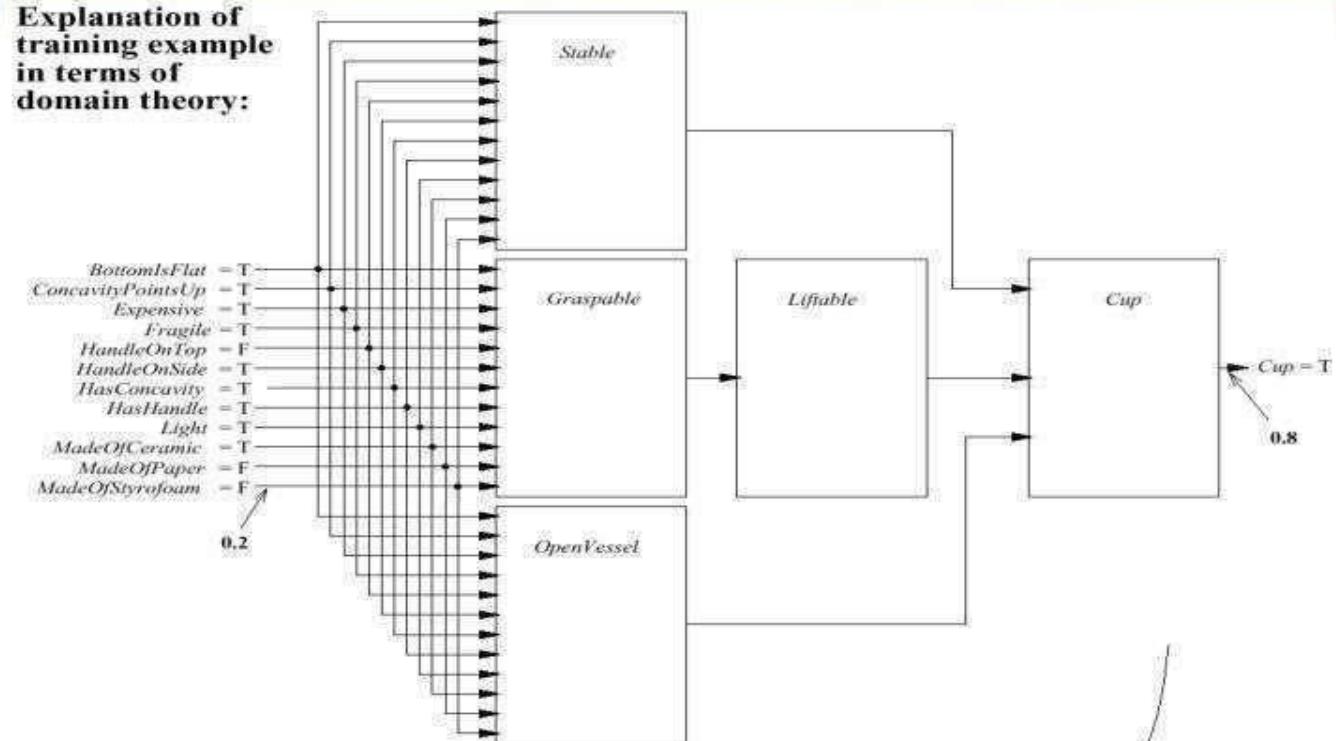
# EBNN

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:**





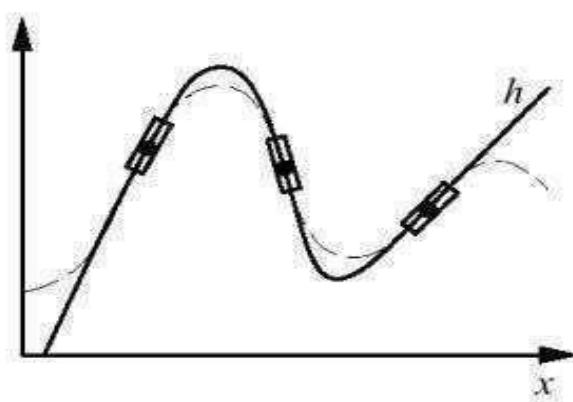
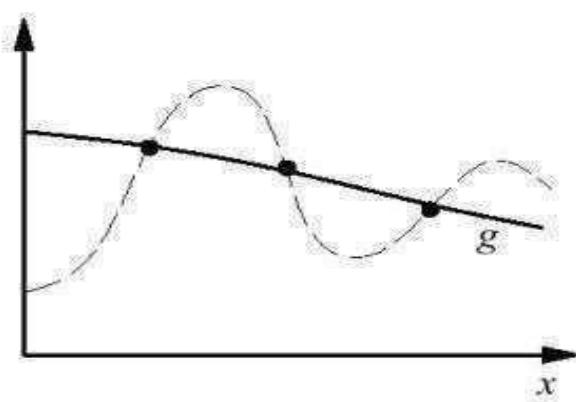
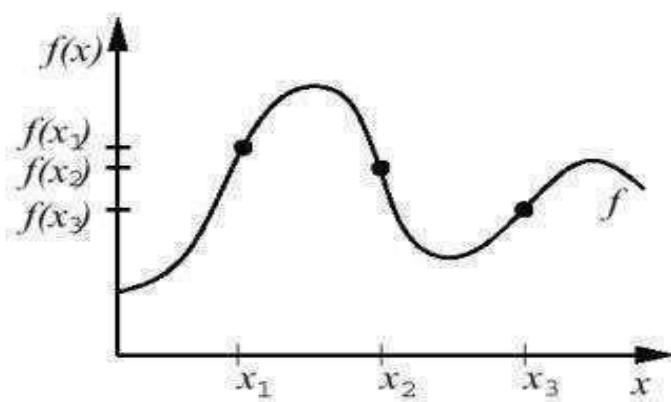
# Modified Objective for Gradient Descent

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

where

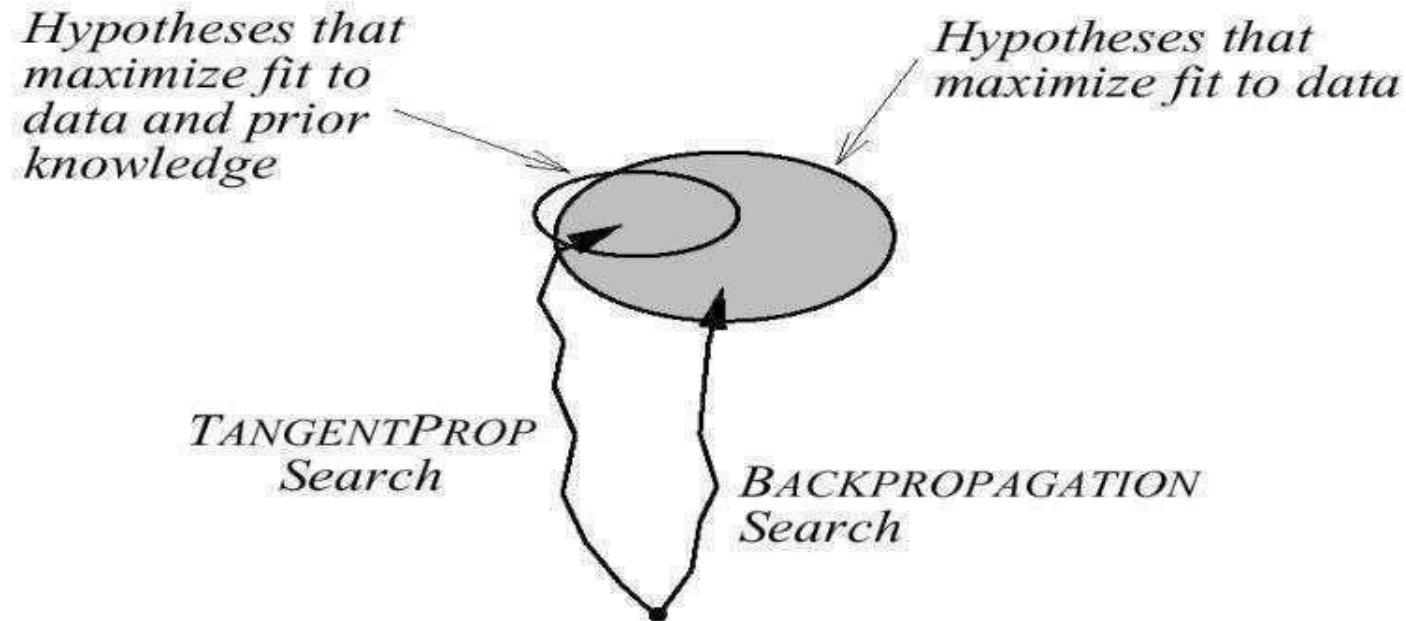
$$\mu_i \equiv 1 - \frac{|A(x_i) - f(x_i)|}{c}$$

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

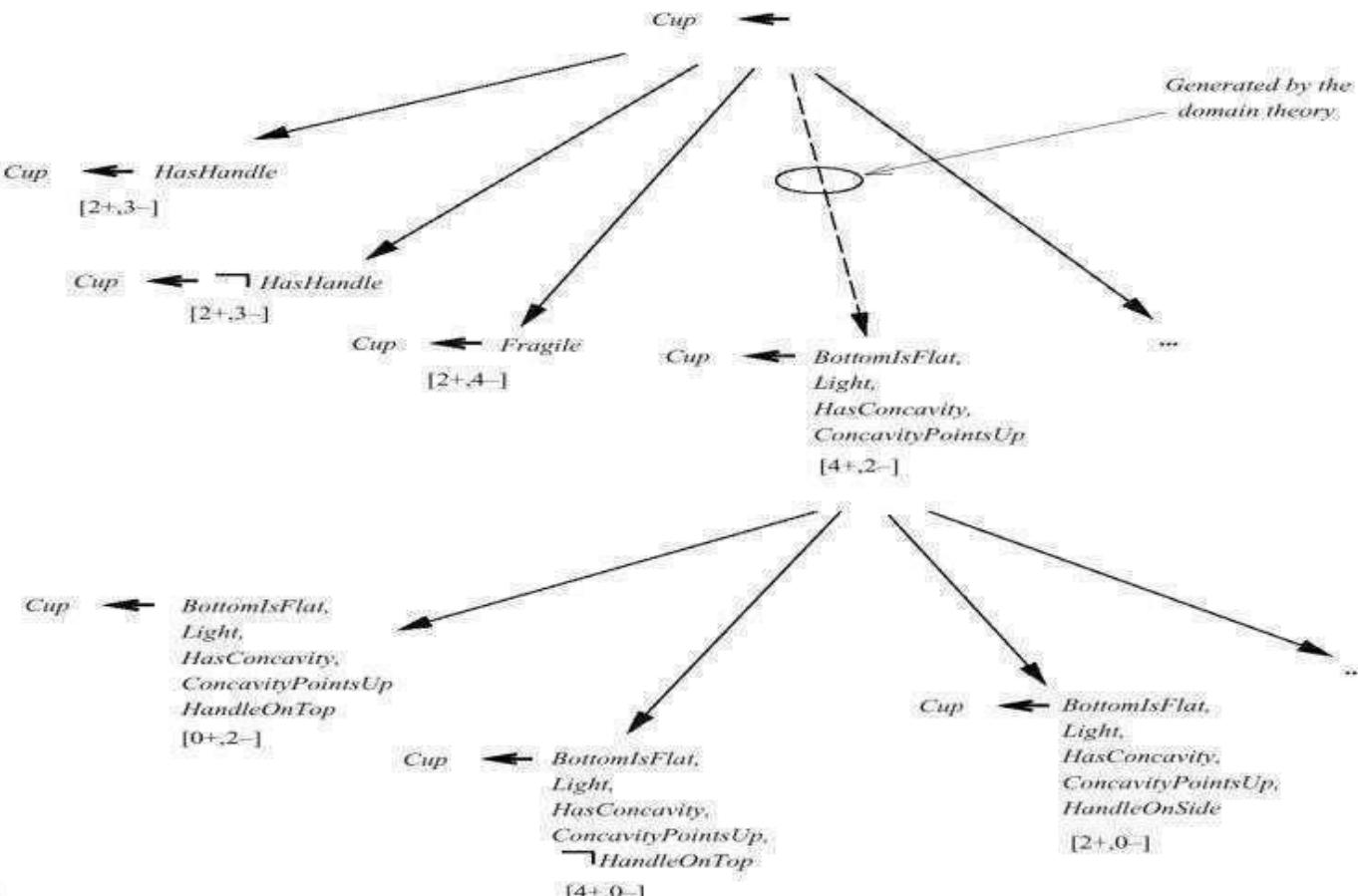


# Hypothesis Space Search in EBNN

## Hypothesis Space



# Search in FOCL



# THANKYOU