

# TDT4173 Machine Learning and Case-Based Reasoning

## Lecture 1 – Introduction

Norwegian University of Science and Technology

Helge Langseth and Anders Kofod-Petersen



- 1 Introduction to Machine learning
  - Machine learning overview
  - Examples
  - The Learning Problem
- 2 Practical information
  - About TDT4173
  - The other stuff
- 3 Checkers
- 4 Learning From Examples
  - EnjoySport - example
  - The Inductive Learning Hypothesis
  - Find-S
  - Version Spaces
  - Learning Bias
- 5 Summary

# The grand vision



An autonomous self-moving machine that **acts**, **reasons**, and **learns** like a human



We are still very far from achieving this. . .

# Why Machine Learning



- Recent progress in algorithms and theory
- Growing flood of online data
- Computational power is available
- Budding industry

## Three niches for machine learning:

- Data mining: using historical data to improve decisions
  - medical records → medical knowledge
- Software applications we can't program by hand
  - autonomous driving
  - speech recognition
- Self customizing programs
  - Recommendation systems

# Typical Datamining Task



## Data:

<i>Patient103</i> time=1	→	<i>Patient103</i> time=2	...	→	<i>Patient103</i> time=n
Age: 23		Age: 23			Age: 23
FirstPregnancy: no		FirstPregnancy: no			FirstPregnancy: no
Anemia: no		Anemia: no			Anemia: no
Diabetes: no		Diabetes: YES			Diabetes: no
PreviousPrematureBirth: no		PreviousPrematureBirth: no			PreviousPrematureBirth: no
Ultrasound: ?		Ultrasound: abnormal			Ultrasound: ?
Elective C-Section: ?		Elective C-Section: no			Elective C-Section: no
Emergency C-Section: ?		Emergency C-Section: ?			<b>Emergency C-Section: Yes</b>
...		...			...

## Given:

- 9714 patient records, each describing a pregnancy and birth
- Each patient record contains 215 features

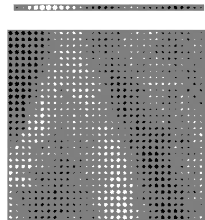
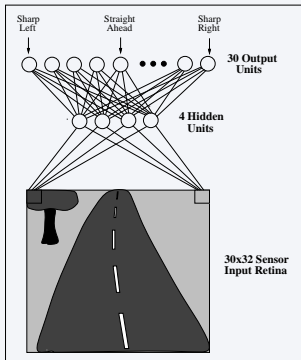
## Learn to predict:

- Classes of future patients at high risk for Emergency Cesarean Section

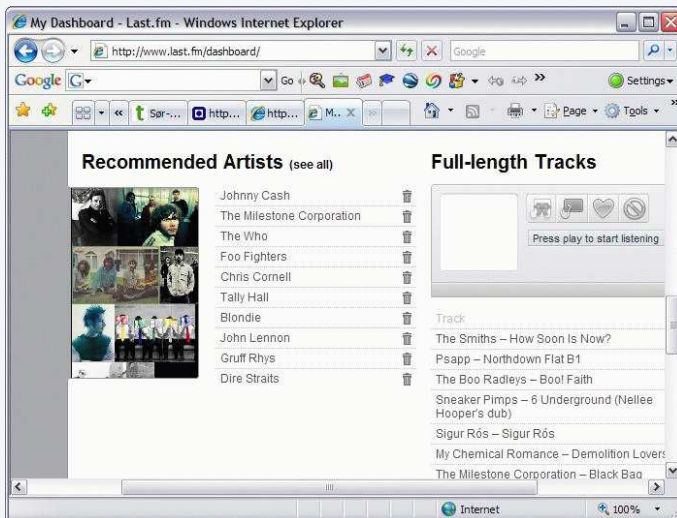
# Problems Too Difficult to Program by Hand



ALVINN [Pomerleau] drives 70 mph on highways



# Software that Customizes to User



<http://www.last.fm>

# Where Is this Headed?



## Today: tip of the iceberg:

- First-generation algorithms: neural nets, decision trees, regression ...
- Applied to well-formatted database
- Budding industry

## Opportunity for tomorrow: enormous impact:

- Learn across full mixed-media data
- Learn by active experimentation
- Cumulative, lifelong learning
- Programming languages with learning embedded?
- ... etc. (Only your imagination limits this list!)



# What is Machine Learning?

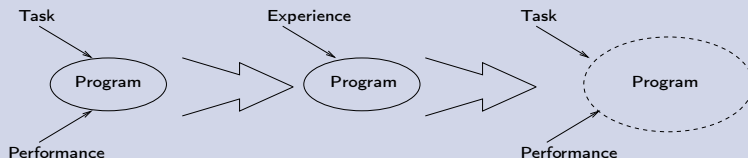


**Learning = Improving with experience at some task**

- Improve over task  $T$ ,
- with respect to performance measure  $P$ ,
- based on experience  $E$ .

## Definition

Machine Learning Methods and techniques that makes computer systems able to update its knowledge and problem-solving ability.



## Goals of the course:

*The course will give a basic insight into principles and methods for how computer systems can learn from its own experience.*

## Syllabus:

- The text-book “Machine Learning” by Tom Mitchell.
- A number of papers top be decided and made available

### How to get it...

Book available at Tapir. Papers will be made available for downloaded from our webpage

# Exercises



- Designed to give hands-on experience with the different machine learning methods we talk about
- Will contain both coding tasks as well as requirements towards discussions
- Typically given with a one-and-a-half week deadline. Time of delivery Thursday at 20:00.
- 1 “big” and 4 “small” assignments.

## NB! Counts towards final grade

All exercises count towards the final grade: If you fail one assignment, you will automatically take off 3.3% (6.7% for the big one) of the total available score.

# Paper presentation



- A number of “classic texts”
- Papers to be presented by students
- Will make up the last part of the semester course

**NB! Counts towards final grade**

Each student must participate in presenting at least one paper. Otherwise, 3.3% of the total available score will be taken off.

# Examination



- **December 6th**
- Written, 4 hours
- Counts for 80% of the final grade; the other 20% are determined by no. assignments handed in + participation in paper presentation.

# Getting information



Sources for information:

- Check the web-page  
<http://www.idi.ntnu.no/emner/tdt4173/>
- It's Learning
- Ask (Contact info on web-page.)

# Reference group



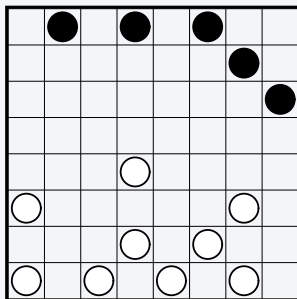
If you want, we can have a reference group.

Not much work (if all goes well):

- Evaluation meeting(s)
- Evaluation report
- Students' spokesman if there is something I should take into account

It is not required, but if two students to volunteer to be in the ref.grp., it may smooth of any problems we meet as we move along.

# Learning to Play Checkers



- *T*: Play checkers
- *P*: Percent of games won in world tournament
- *E*: opportunity to play against self



# Design Choices



- What experience can we learn from?
- What exactly should be learned?
- How shall it be represented?
- Target function:
  - collection of rules?
  - neural network ?
  - polynomial function of board features?
  - ...
- What specific algorithm can we use to learn it?

# Type of knowledge learned



We wish to **learn** a function that for any given board position  $\mathcal{B}$  chooses the best move  $\mathcal{M}$ , **ChooseMove**:  $\mathcal{B} \rightarrow \mathcal{M}$ .

# Type of knowledge learned

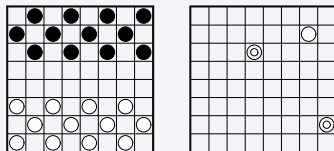


We wish to **learn** a function that for any given board position  $\mathcal{B}$  chooses the best move  $\mathcal{M}$ , **ChooseMove**:  $\mathcal{B} \rightarrow \mathcal{M}$ .

**Direct training:** Examples of individual checkers board states and the correct move for each.

**Indirect training:** Examples of sequences of moves and final outcomes of the various games played.

Indirect training makes **ChooseMove** impractical to learn:



If we end up winning, is the first move then optimal?

# Approximation – The start of the learning work



Instead of ChooseMove, we establish a *Value function*  $V$ :

$$V : \mathcal{B} \rightarrow \mathbb{R}$$

that maps legal board states  $\mathcal{B}$  into some real value.

**Playing rule:** For any board position, choose the move that maximizes the value of the resulting board position.

# Approximation – The start of the learning work



Instead of ChooseMove, we establish a *Value function*  $V$ :

$$V : \mathcal{B} \rightarrow \mathbb{R}$$

that maps legal board states  $\mathcal{B}$  into some real value.

**Playing rule:** For any board position, choose the move that maximizes the value of the resulting board position.

- ❶ if  $b$  is a final board state that is won, then  $V(b) = 100$
- ❷ if  $b$  is a final board state that is lost, then  $V(b) = -100$
- ❸ if  $b$  is a final board state that is drawn, then  $V(b) = 0$
- ❹ if  $b$  is not a final state in the game, then  $V(b) = ??$

# Approximation – The start of the learning work



Instead of ChooseMove, we establish a *Value function*  $V$ :

$$V : \mathcal{B} \rightarrow \mathbb{R}$$

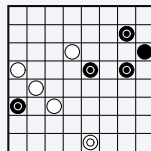
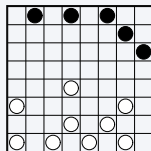
that maps legal board states  $\mathcal{B}$  into some real value.

**Playing rule:** For any board position, choose the move that maximizes the value of the resulting board position.

- ❶ if  $b$  is a final board state that is won, then  $V(b) = 100$
- ❷ if  $b$  is a final board state that is lost, then  $V(b) = -100$
- ❸ if  $b$  is a final board state that is drawn, then  $V(b) = 0$
- ❹ if  $b$  is not a final state in the game, then  $V(b) = V(b')$ ,  
 $b'$  is the best final board state that can be achieved starting from  $b$  and playing optimally until the end of the game.

It is still not trivial...

# What is of importance?



$x_1$ : # black pieces

$x_3$ : # black kings

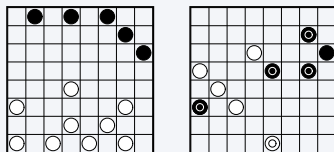
$x_5$ : # white pieces threatened

$x_2$ : # white pieces

$x_4$ : # white kings

$x_6$ : # black pieces threatened

# What is of importance?



$x_1$ : # black pieces

$x_3$ : # black kings

$x_5$ : # white pieces threatened

$x_2$ : # white pieces

$x_4$ : # white kings

$x_6$ : # black pieces threatened

## Approximation:

$\hat{V}(b) = w_0 + w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + w_5x_5 + w_6x_6$ ,  
 where  $w_i$  is the weight assigned to  $x_i$ .

## Learning task:

Determine the weights  $w_0$ ,  $w_1$ ,  $w_2$ ,  $w_3$ ,  $w_4$ ,  $w_5$ , and  $w_6$ .



# How to learn



In order to learn  $\hat{V}$ , we require a set of training examples, each describing a board state  $b$  and a training value  $V_{\text{train}}$  for  $b$ :

- **Current weights:**  $w_0, w_1, w_2, w_3, w_4, w_5$ , and  $w_6$  yield  $\hat{V}$ .
- **New game:**  $b_1, b_2, \dots, b_{\text{end}}$

What value should  $V_{\text{train}}(b_i)$  should we attach to position  $b_i$ ?

# How to learn



In order to learn  $\hat{V}$ , we require a set of training examples, each describing a board state  $b$  and a training value  $V_{\text{train}}$  for  $b$ :

- **Current weights:**  $w_0, w_1, w_2, w_3, w_4, w_5$ , and  $w_6$  yield  $\hat{V}$ .
- **New game:**  $b_1, b_2, \dots, b_{\text{end}}$

What value should  $V_{\text{train}}(b_i)$  should we attach to position  $b_i$ ?

## Idea:

When faced with a situation  $b_k$ , both players do the best they can – resulting in  $b_{\text{end}}$ .

# How to learn



In order to learn  $\hat{V}$ , we require a set of training examples, each describing a board state  $b$  and a training value  $V_{\text{train}}$  for  $b$ :

- **Current weights:**  $w_0, w_1, w_2, w_3, w_4, w_5$ , and  $w_6$  yield  $\hat{V}$ .
- **New game:**  $b_1, b_2, \dots, b_{\text{end}}$

What value should  $V_{\text{train}}(b_i)$  should we attach to position  $b_i$ ?

## Idea:

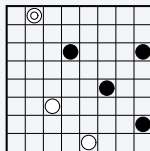
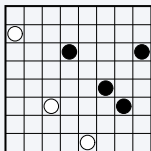
When faced with a situation  $b_k$ , both players do the best they can – resulting in  $b_{\text{end}}$ .

## In general:

$$V_{\text{train}}(b_i) \leftarrow \hat{V}(b_{i+1})$$

This makes sense if  $\hat{V}$  is more accurate for board states closer to the end of the game.

# Ehhh... And what does this mean?



Current state:  $b_i$     Next state:  $b_{i+1}$

- The (system believes that) situation  $b_i \rightsquigarrow b_{i+1}$
- Therefore  $V_{\text{train}}(b_i) = V(b_{i+1})$
- $V(b_{i+1})$  is unknown, but assuming the system is very good, we have  $\hat{V}(b_{i+1}) \approx V(b_{i+1})$ . Thus, we decide that

$$V_{\text{train}}(b_i) \leftarrow \hat{V}(b_{i+1}).$$

# And will it work?



Can we get reasonable training data?

- We know what  $V(b_{\text{end}})$  is for any state  $b_{\text{end}}$
- Using the previous setup, we should therefore be able to value situations that are *onestep* away from being finished!
- ...and using the same setup again, we should next be able to value situations that are *twosteps* away from being finished
- ...and so on

This should work, but we need to be able to *use* the training data...

# How to learn the weights



- Current weights:  $w_0, w_1, w_2, w_3, w_4, w_5$ , and  $w_6$  yield  $\hat{V}$ .
- New game:  $\langle b_1, V_{\text{train}}(b_1) \rangle, \dots, \langle b_{\text{end}}, V_{\text{train}}(b_{\text{end}}) \rangle$ .

# How to learn the weights



- **Current weights:**  $w_0, w_1, w_2, w_3, w_4, w_5$ , and  $w_6$  yield  $\hat{V}$ .
- **New game:**  $\langle b_1, V_{\text{train}}(b_1) \rangle, \dots, \langle b_{\text{end}}, V_{\text{train}}(b_{\text{end}}) \rangle$ .

**Idea:** Introduce *error function*  $E$ , and change weights such that the total error over all training examples is minimal.

$$E = \sum_{\langle b, V_{\text{train}}(b) \rangle \in \text{training examples}} \left( V_{\text{train}}(b) - \hat{V}(b) \right)^2$$

**Note:**

$E$  is a function of the weights,  $E = E(w_0, w_1, w_2, w_3, w_4, w_5, w_6)$ , and we will change the weights to make  $E$  obtain its minimal value.

# LMS Weight update rule



For each training set  $\langle b, V_{\text{train}}(b) \rangle$  do:

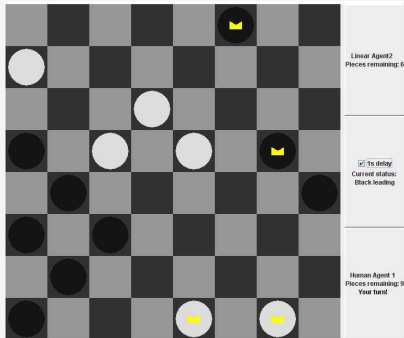
- Use the current weights to calculate  $\hat{V}(b)$ .
- For each weight  $w_i$  do:

$$w_i \leftarrow w_i + \mu \cdot x_i \cdot \left( V_{\text{train}}(b) - \hat{V}(b) \right)$$

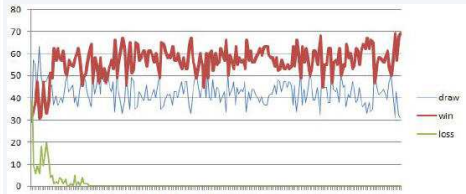
where  $\mu$  is the *learning rate*.



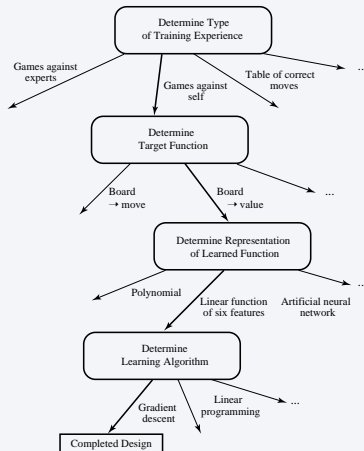
# Implementation: A 5th year Student Project



## Results



# Design Choices



# Training Examples for *EnjoySport*



Index	Sky	Temp	Humid	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

What is the general concept?

# Training Examples for *EnjoySport*



Index	Sky	Temp	Humid	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

What is the general concept?

- “Sky = Sunny”?
- “Sky = Sunny **AND** Temp = Warm”?
- “Forecast = Same **OR** Water = Cool”?
- “When Index is written in binary digits it requires *one*1”?

# Representing Hypotheses



Many possible representations

Here,  $h$  is conjunction of constraints on attributes

Each constraint can be

- a specific value (e.g., "Water = Warm")
- don't care (e.g., "Water = ?")
- no value allowed (e.g., "Water =  $\emptyset$ ")

For example,

Sky	AirTemp	Humid	Wind	Water	Forecast
⟨ Sunny	?	?	Strong	?	Same ⟩

# Prototypical Concept Learning Task



## Given:

- Instances  $X$ : Possible days, each described by the attributes Sky, AirTemp, Humidity, Wind, Water, Forecast
- Target function  $c$ : EnjoySport:  $X \rightarrow \{0, 1\}$
- Hypotheses  $H$ : Conjunctions of literals, e.g.,  $\langle ?, \text{Cold}, \text{High}, ?, ?, ? \rangle$ .
- Training examples  $D$ : Positive and negative examples of the target function  $\langle x_1, c(x_1) \rangle, \dots, \langle x_m, c(x_m) \rangle$

**Determine** a hypothesis  $h \in H$  such that  $\forall x \in D : h(x) = c(x)$ .

# Prototypical Concept Learning Task



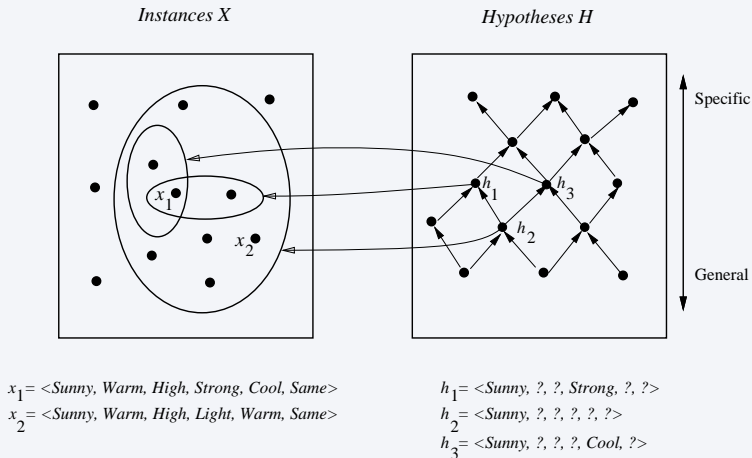
## Given:

- Instances  $X$ : Possible days, each described by the attributes Sky, AirTemp, Humidity, Wind, Water, Forecast
- Target function  $c$ : EnjoySport:  $X \rightarrow \{0, 1\}$
- Hypotheses  $H$ : Conjunctions of literals, e.g.,  $\langle ?, \text{Cold}, \text{High}, ?, ?, ? \rangle$ .
- Training examples  $D$ : Positive and negative examples of the target function  $\langle x_1, c(x_1) \rangle, \dots \langle x_m, c(x_m) \rangle$

**Determine** a hypothesis  $h \in H$  such that  $\forall x \in D : h(x) = c(x)$ .

**The inductive learning hypothesis:** 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.

# Instance, Hypotheses, and More-General-Than



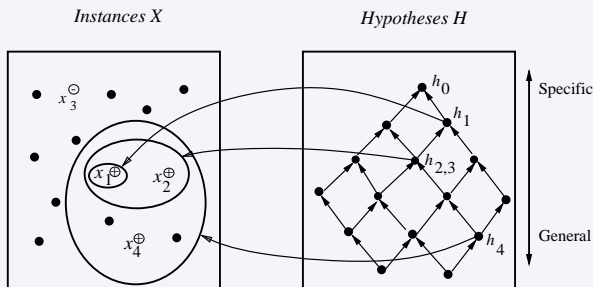


# Find-S Algorithm



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

## Hypothesis Space Search by Find-S



$x_1 = \langle \text{Sunny Warm Normal Strong Warm Same} \rangle, +$   
 $x_2 = \langle \text{Sunny Warm High Strong Warm Same} \rangle, +$   
 $x_3 = \langle \text{Rainy Cold High Strong Warm Change} \rangle, -$   
 $x_4 = \langle \text{Sunny Warm High Strong Cool Change} \rangle, +$

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

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

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

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

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

# Complaints about Find-S



- Can't tell whether it has learned concept
- Can't tell when training data inconsistent
- Picks a maximally specific  $h$  (why?)
- Depending on  $H$ , there might be several!

# Version Spaces



A hypothesis  $h$  is **consistent** with a set of training examples  $D$  of target concept  $c$  if and only if  $h(x) = c(x)$  for each training example  $\langle x, c(x) \rangle$  in  $D$ .

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

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

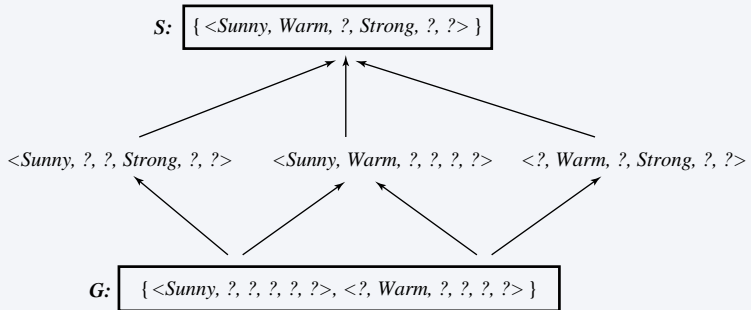
$$VS_{H,D} \equiv \{h \in H \mid \text{Consistent}(h, D)\}$$

# The List-Then-Eliminate Algorithm



- 1  $\text{VersionSpace} \leftarrow$  a list containing every hypothesis in  $H$
- 2 For each training example,  $\langle x, c(x) \rangle$   
remove from  $\text{VersionSpace}$  any hypothesis  $h$  for which  
 $h(x) \neq c(x)$
- 3 Output the list of hypotheses in  $\text{VersionSpace}$

# Example Version Space



# Representing Version Spaces



- The **General boundary**,  $G$ , of version space  $VS_{H,D}$  is the set of its maximally general members
- The **Specific boundary**,  $S$ , of version space  $VS_{H,D}$  is the set of its maximally specific members
- Every member of the version space lies between these boundaries

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

where  $x \geq y$  means  $x$  is more general or equal to  $y$

# Candidate Elimination Algorithm



$G \leftarrow$  maximally general hypotheses in  $H$

$S \leftarrow$  maximally specific hypotheses in  $H$

For each training example  $d$ , do

- If  $d$  is a positive example
  - Remove from  $G$  any hypothesis inconsistent with  $d$
  - For each hypothesis  $s$  in  $S$  that is not consistent with  $d$ 
    - Remove  $s$  from  $S$
    - Add to  $S$  all minimal *generalizations*  $h$  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  $S$

[CONT'd]



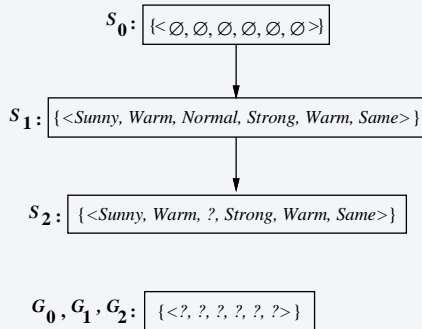
# Candidate Elimination Algorithm



## [...FROM PREVIOUS SLIDE]

- If  $d$  is a negative example
  - Remove from  $S$  any hypothesis inconsistent with  $d$
  - For each hypothesis  $g$  in  $G$  that is not consistent with  $d$ 
    - Remove  $g$  from  $G$
    - Add to  $G$  all minimal *specializations*  $h$  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  $G$

# Example Trace



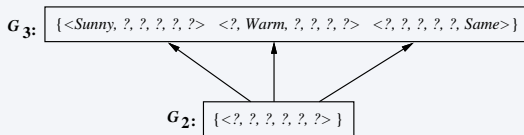
Training examples:

1.  $\langle \text{Sunny}, \text{Warm}, \text{Normal}, \text{Strong}, \text{Warm}, \text{Same} \rangle, \text{Enjoy Sport} = \text{Yes}$
2.  $\langle \text{Sunny}, \text{Warm}, \text{High}, \text{Strong}, \text{Warm}, \text{Same} \rangle, \text{Enjoy Sport} = \text{Yes}$

# Example Trace



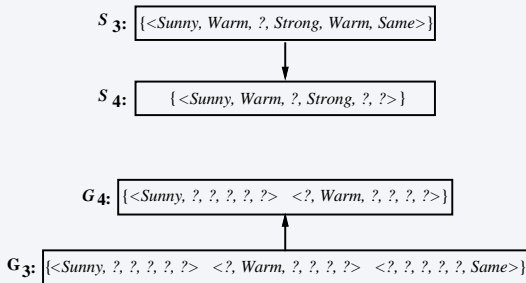
$S_2, S_3$ : { <Sunny, Warm, ?, Strong, Warm, Same> }



Training Example:

3. <Rainy, Cold, High, Strong, Warm, Change>, EnjoySport=No

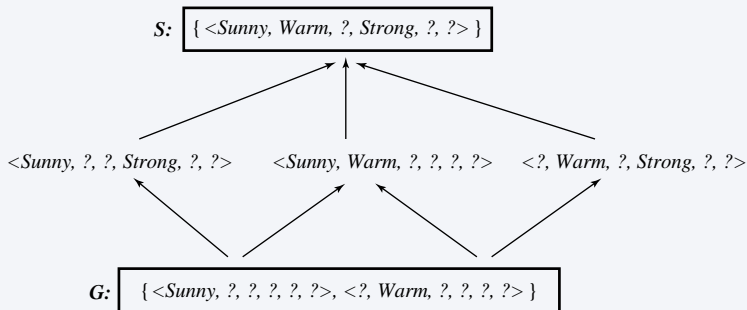
# Example Trace



Training Example:

4.<Sunny, Warm, High, Strong, Cool, Change>, EnjoySport = Yes

# How Should These Be Classified??



$\langle \text{Sunny Warm Normal Strong Cool Change} \rangle$

$\langle \text{Rainy Cool Normal Light Warm Same} \rangle$

$\langle \text{Sunny Warm Normal Light Warm Same} \rangle$

# What Justifies this Inductive Leap?



- +  $\langle \text{Sunny Warm Normal Strong Cool Change} \rangle$
  - +  $\langle \text{Sunny Warm Normal Light Warm Same} \rangle$
- 

$S$  :  $\langle \text{Sunny Warm Normal ? ? ?} \rangle$

Why believe we can classify the unseen

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

# Inductive Bias



## Consider

- concept learning algorithm  $L$
- instances  $X$ , target concept  $c$
- training examples  $D_c = \{\langle x, c(x) \rangle\}$
- let  $L(x_i, D_c)$  denote the classification assigned to the instance  $x_i$  by  $L$  after training on data  $D_c$ .

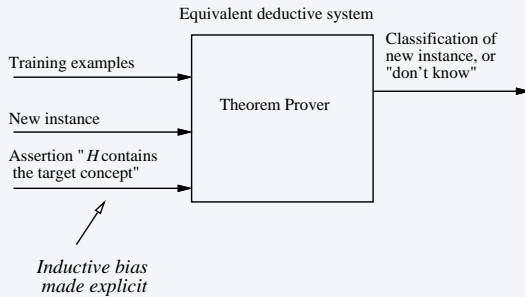
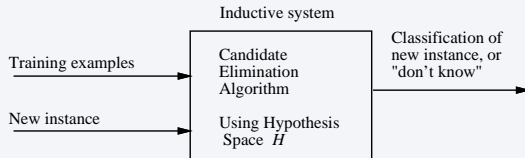
## Definition:

The **inductive bias** of  $L$  is any minimal set of assertions  $B$  such that for any target concept  $c$  and corresponding training examples  $D_c$

$$(\forall x_i \in X)[(B \wedge D_c \wedge x_i) \vdash L(x_i, D_c)]$$

where  $A \vdash B$  means  $A$  logically entails  $B$ .

# Inductive and Equivalent Deductive Systems





# Three Learners with Different Biases



- ① *Rote learner*: Store examples, Classify  $x$  iff it matches previously observed example.
- ② *Version space candidate elimination algorithm*
- ③ *Find-S*

# Summary Points



- 1 Concept learning as search through  $H$
- 2 General-to-specific ordering over  $H$
- 3 Version space candidate elimination algorithm
- 4  $S$  and  $G$  boundaries characterize learner's uncertainty
- 5 Learner can generate useful queries
- 6 Inductive leaps possible only if learner is biased
- 7 Inductive learners can be modelled by equivalent deductive systems