

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 *Air Temp*, *Humidity*, *Wind*, *Water Forecast* each have two possible values, the instance space X contains exactly

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

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

1 + (4.3.3.3.3.3) = 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 set 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 .

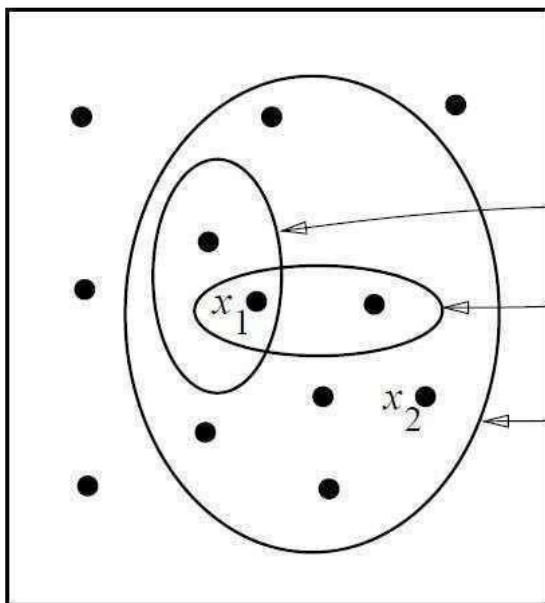
General-to-Specific Ordering of Hypotheses

- 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

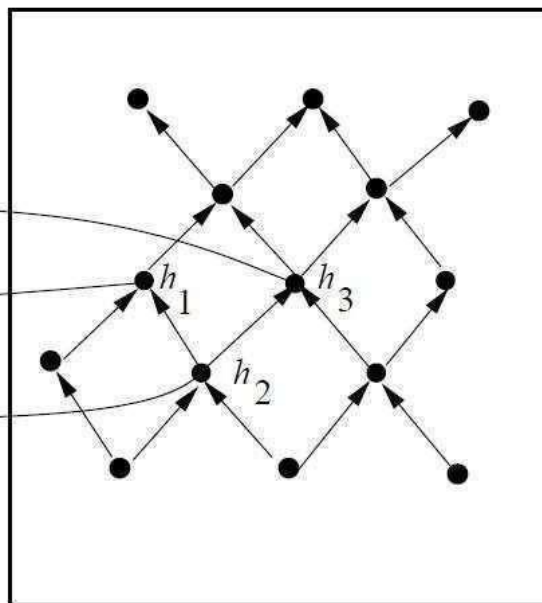
Definition: Let h_j and h_k be Boolean-valued functions defined over X . Then h_j is more general-than-or-equal-to h_k (written $h_j \geq h_k$) if and only if

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

Instances X



Hypotheses H



Specific
↑
↓
General

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

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

- 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 is to initialize h to the most specific hypothesis in H

$$h = (\emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset)$$

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

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

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

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

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

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

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

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

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

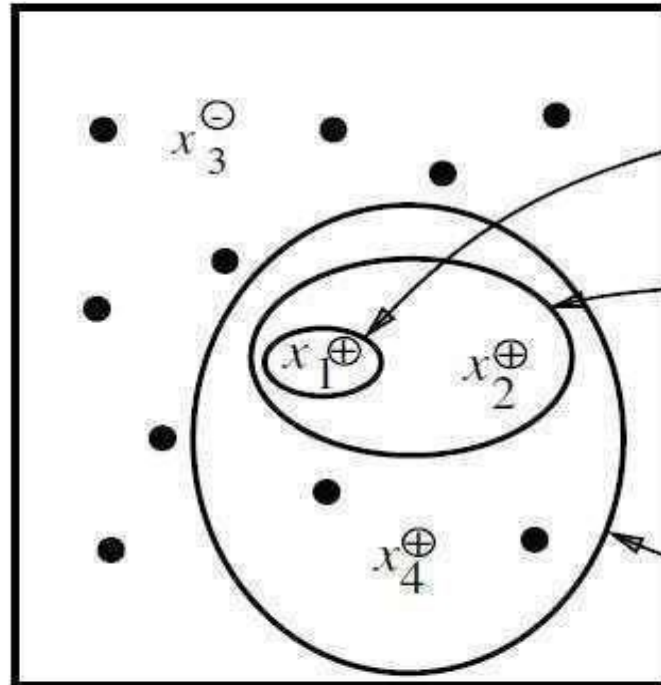
$h_3 = \langle \text{Sunny Warm? Strong Warm Same} \rangle$

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

The fourth example leads to a further generalization of h

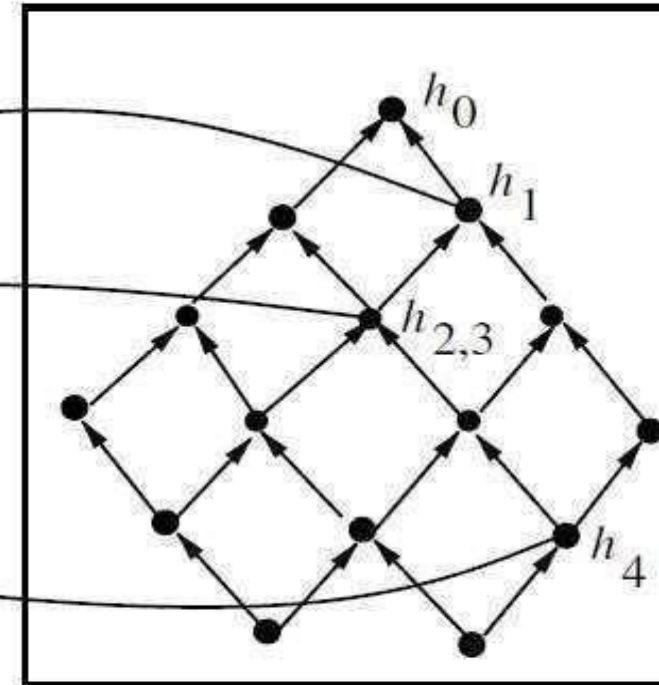
$h_4 = \langle \text{Sunny Warm? Strong??} \rangle$

Instances X



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

Hypotheses H



- $h_0 = \langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle$
- $h_1 = \langle \text{Sunny Warm Normal Strong Warm Same} \rangle$
- $h_2 = \langle \text{Sunny Warm ? Strong Warm Same} \rangle$
- $h_3 = \langle \text{Sunny Warm ? Strong Warm Same} \rangle$
- $h_4 = \langle \text{Sunny Warm ? Strong ? ?} \rangle$

The key property of the FIND-S algorithm is

- FIND-S is guaranteed to output the most specific hypothesis within H that is consistent 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 the training examples*

Representation

- **Definition:** A hypothesis 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 the difference between the definitions of *consistent* and *satisfies*

- An example x 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 x is said to *consistent* with hypothesis h iff $h(x) = c(x)$

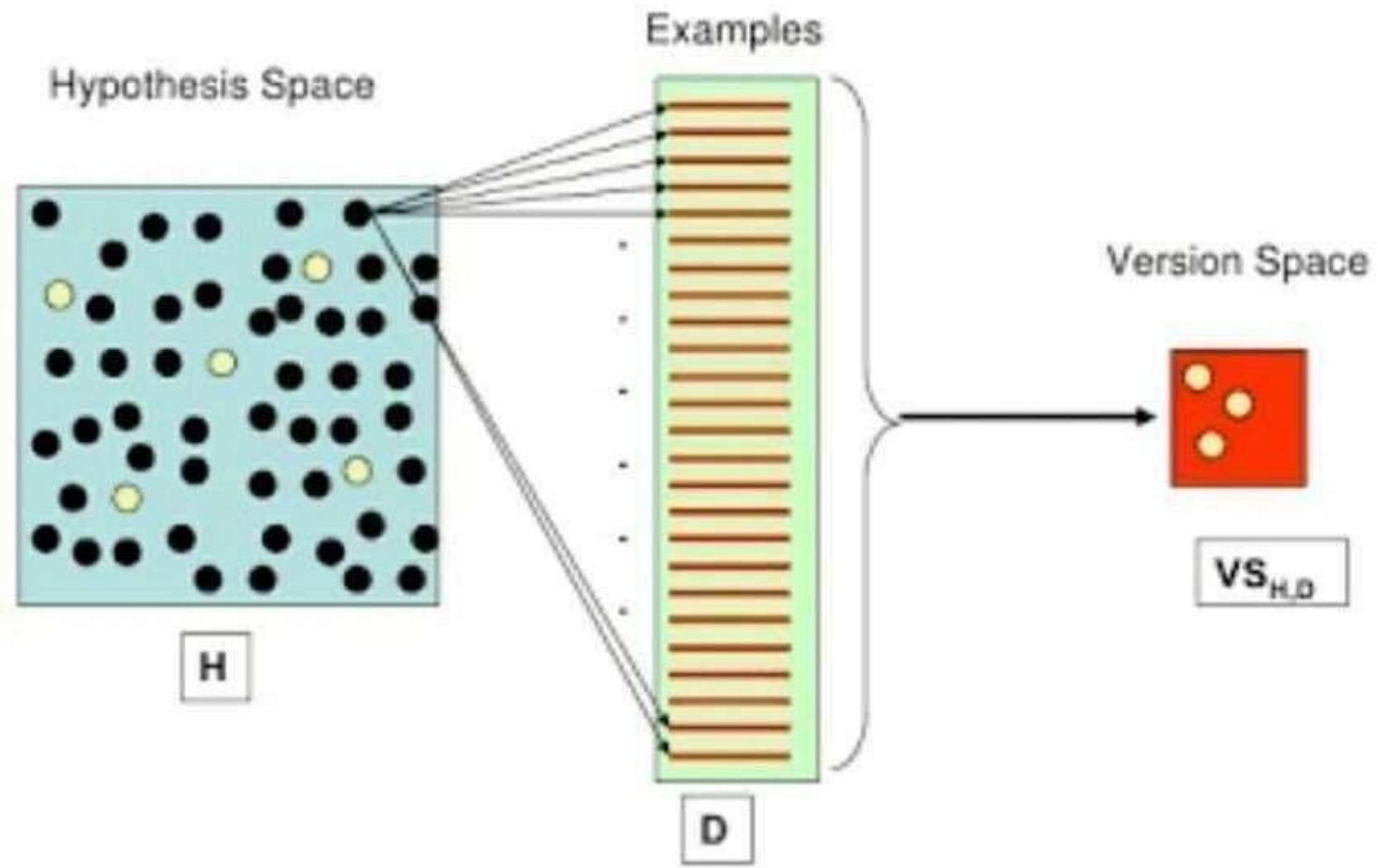
VersionSpace



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

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

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



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.

The LIST-THEN-ELIMINATE

Algorithm



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*

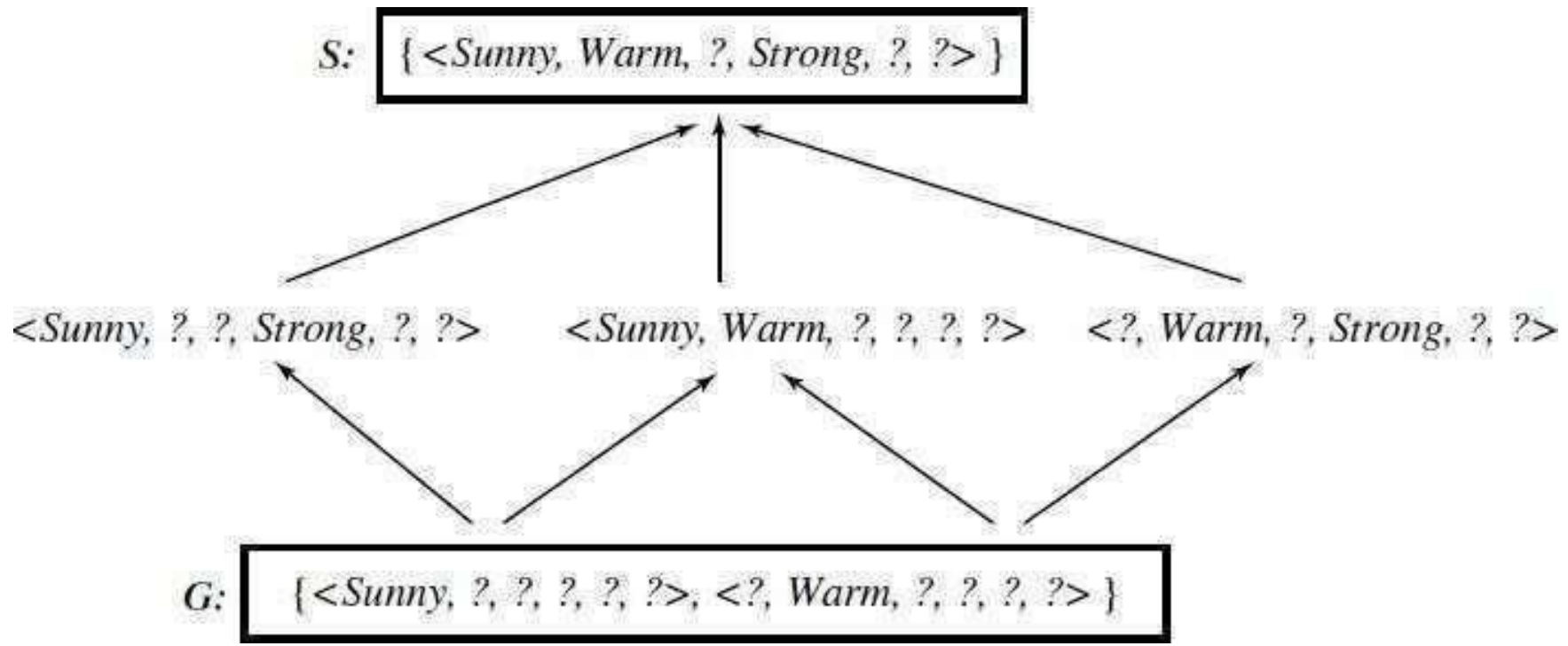
The LIST-THEN-ELIMINATE Algorithm

- **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.



- A version space with its general and specific boundary sets
- The version space includes all six hypotheses shown here, but can be represented more simply by S and G.
- Arrows indicate instances of the more-general-than relation. This is the version space for the EnjoySport concept learning problem and training examples described in the below table

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

Deepak D, Asst. Prof., Dept. of CSE, Canara Engg. College

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

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 .

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

VersionSpace 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,

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

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

To Prove:

1. Every h satisfying the right hand side of the above expression is in $VS_{H,D}$
2. Every member of $VS_{H,D}$ 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 $g \geq_g h$ 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 $VS_{H,D}$

2. It can be proven by assuming some h in $VS_{H,D}$, 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
 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

- 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

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 example,

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

S_0

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



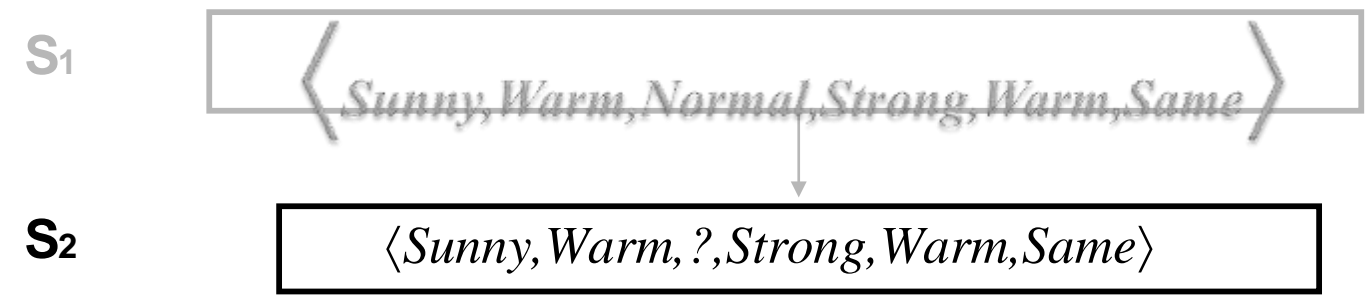
S_1

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

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

For training example,

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



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

For training example,

⟨Rainy,Cold,High,Strong,Warm,Change⟩-

S₂, S₃

⟨~~Sunny,Warm,?,Strong~~, Warm, Same⟩

G₃

⟨Sunny,?,?,?,?,?⟩⟨?, Warm,?,?,?,?,?⟩⟨?,?,?,?,?,Same⟩

G₂

⟨?,?,?,?,?,?⟩

For training example,

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

S₃ $\langle \text{Sunny, Warm, ?, Strong, Warm, Same} \rangle$



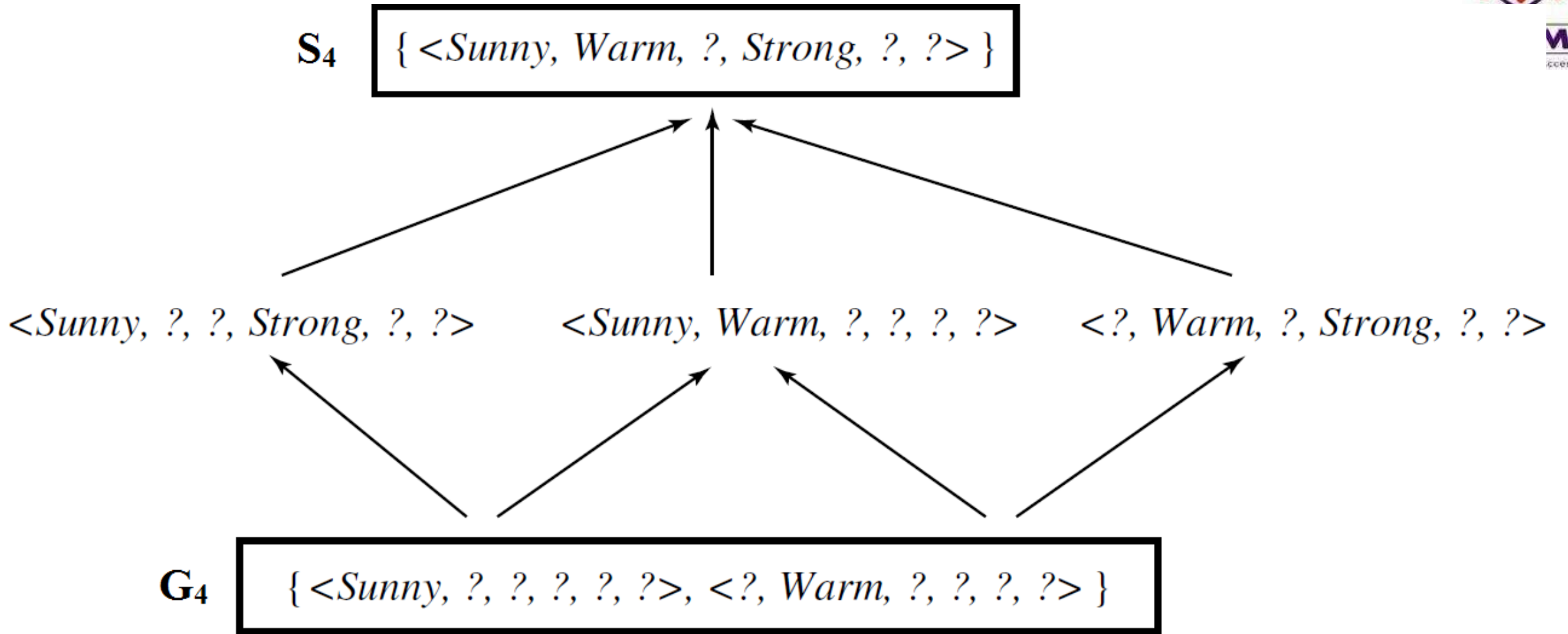
S₄ $\langle \text{Sunny, Warm, ?, Strong, ?, ?} \rangle$

G₄ $\langle \text{Sunny, ?, ?, ?, ?, ?} \rangle \langle \text{?, Warm, ?, ?, ?, ?} \rangle$



G₃ $\langle \text{Sunny, ?, ?, ?, ?, ?} \rangle \langle \text{?, Warm, ?, ?, ?, ?} \rangle \langle \text{?, ?, ?, ?, ?, Same} \rangle$





The final version space for the *EnjoySport* 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 instance 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 power set of X) is $2^{|X|}$
- Recall that there are 96 instances in *EnjoySport*; 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 *Enjoysport* 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 power set of X .
- One way to define such an H' is to allow arbitrary disjunctions, conjunctions, and negations of four earlier hypotheses.

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

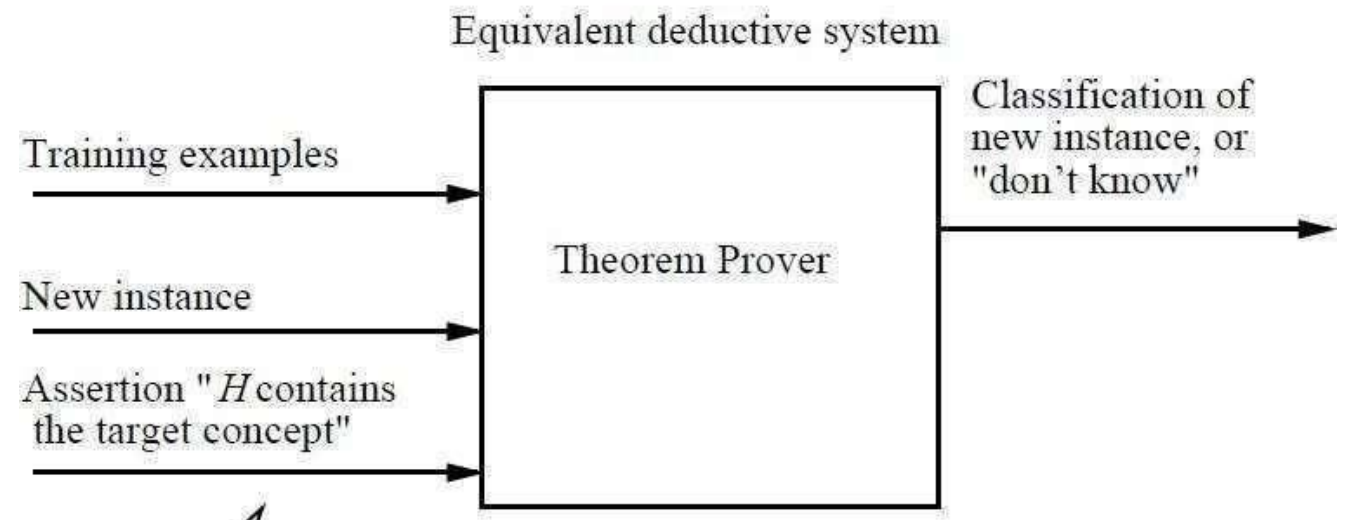
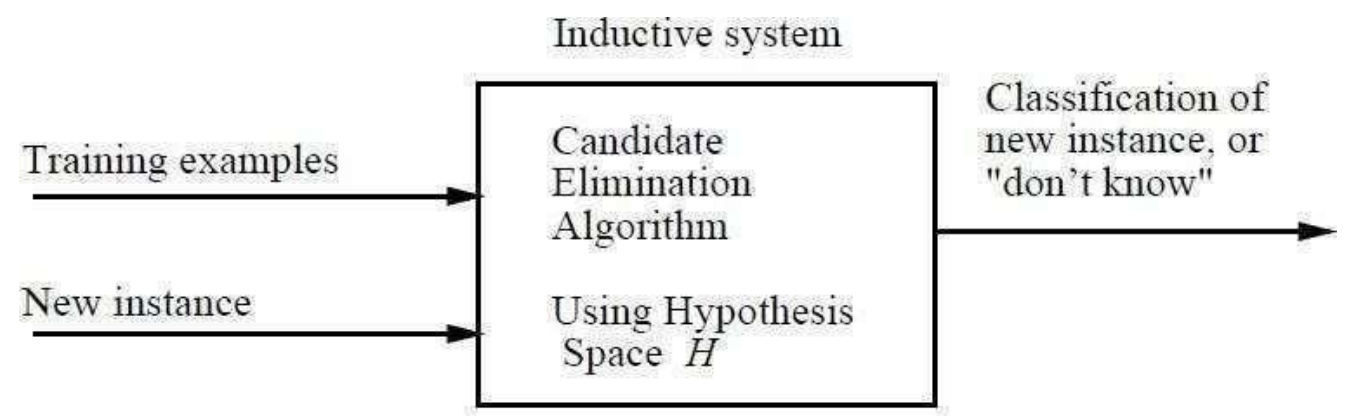
$$(\mathit{Sunny}, ?, ?, ?, ?, ?) \vee (\mathit{Cloudy}, ?, ?, ?, ?, ?)$$

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

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



*Inductive bias
made explicit*

Modelling inductive systems by equivalent deductive systems.

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.