

---

# Inductive Classification

Based on the ML lecture by Raymond J. Mooney  
University of Texas at Austin

# Classification (Categorization)

---

- Given:
  - A description of an instance,  $x \in X$ , where  $X$  is the *instance language* or *instance space*.
  - A fixed set of categories:  $C = \{c_1, c_2, \dots, c_n\}$
- Determine:
  - The category of  $x$ :  $c(x) \in C$ , where  $c(x)$  is a categorization function whose domain is  $X$  and whose range is  $C$ .
  - If  $c(x)$  is a binary function  $C = \{0, 1\}$  ( $\{\text{true}, \text{false}\}$ ,  $\{\text{positive}, \text{negative}\}$ ) then it is called a *concept*.

# Learning for Categorization

---

- A training example is an instance  $x \in X$ , paired with its correct category  $c(x)$ :  $\langle x, c(x) \rangle$  for an unknown categorization function,  $c$ .
- Given a set of training examples,  $D$ .
- Find a hypothesized categorization function,  $h(x)$ , such that:

$$\forall \langle x, c(x) \rangle \in D : h(x) = c(x)$$

*Consistency*

# Sample Category Learning Problem

---

- Instance language:  $\langle \text{size, color, shape} \rangle$ 
  - $\text{size} \in \{\text{small, medium, large}\}$
  - $\text{color} \in \{\text{red, blue, green}\}$
  - $\text{shape} \in \{\text{square, circle, triangle}\}$
- $C = \{\text{positive, negative}\}$

•  $D$ :

Example	Size	Color	Shape	Category
1	small	red	circle	positive
2	large	red	circle	positive
3	small	red	triangle	negative
4	large	blue	circle	negative

# Hypothesis Selection

---

- Many hypotheses are usually consistent with the training data.
  - red & circle
  - (small & circle) or (large & red)
  - (small & red & circle) or (large & red & circle)

# Generalization

---

- Hypotheses must generalize to correctly classify instances not in the training data.
- Simply memorizing training examples is a consistent hypothesis that does not generalize. But ...

# Hypothesis Space

---

- For learning concepts on instances described by  $n$  discrete-valued features, consider the space of conjunctive hypotheses represented by a vector of  $n$  constraints  $\langle c_1, c_2, \dots, c_n \rangle$  where each  $c_i$  is either:
  - $?$ , a wild card indicating no constraint on the  $i$ th feature
  - A specific value from the domain of the  $i$ th feature
  - $\emptyset$  indicating no value is acceptable
- Other notations
  - (Size = big) AND (Color = red)
  - size(Id, big), size(Id, red) . . . In Prolog
- Sample conjunctive hypotheses are
  - $\langle \text{big}, \text{red}, ? \rangle$
  - $\langle ?, ?, ? \rangle$  (most general hypothesis)
  - $\langle \emptyset, \emptyset, \emptyset \rangle$  (most specific hypothesis)

# Inductive Learning Hypothesis

---

- Any function that is found to approximate the target concept well on a sufficiently large set of training examples will also approximate the target function well on unobserved examples.
- Assumes that the training and test examples are drawn independently from the same underlying distribution.
- This is a fundamentally unprovable hypothesis unless additional assumptions are made about the target concept and the notion of “approximating the target function well on unobserved examples” is defined appropriately (cf. computational learning theory).



# Evaluation of Classification Learning

---

- Classification accuracy (% of instances classified correctly).
  - Measured on an independent test data.
- Training time (efficiency of training algorithm).
- Complexity of the hypothesis that has been learned
- Testing time (efficiency of subsequent classification).

# Category Learning as Search

---

- Category learning can be viewed as searching the hypothesis space for one (or more) hypotheses that are consistent with the training data.
- Consider an instance space consisting of  $n$  binary features which therefore has  $2^n$  instances.
- For conjunctive hypotheses, there are 4 choices for each feature:  $\emptyset$ , T, F, ?, so there are  $4^n$  syntactically distinct hypotheses.
- However, all hypotheses with 1 or more  $\emptyset$ s are equivalent, so there are  $3^{n+1}$  semantically distinct hypotheses.
- The target binary categorization function in principle could be any of the possible  $2^{2^n}$  functions on  $n$  input bits.
- Therefore, conjunctive hypotheses are a small subset of the space of possible functions, but both are intractably large.
- All reasonable hypothesis spaces are intractably large or even infinite.

# Learning by Enumeration

---

- For any finite or countably infinite hypothesis space, one can simply enumerate and test hypotheses one at a time until a consistent one is found.

For each  $h$  in  $H$  do:

    If  $h$  is consistent with the training data  $D$ ,  
    then terminate and return  $h$ .

- This algorithm is guaranteed to terminate with a consistent hypothesis if one exists; however, it is obviously computationally intractable for almost any practical problem.

# Efficient Learning

---

- Is there a way to learn conjunctive concepts without enumerating them?
- How do human subjects learn conjunctive concepts?
- Is there a way to efficiently find an unconstrained boolean function consistent with a set of discrete-valued training instances?
- If so, is it a useful/practical algorithm?

# Conjunctive Rule Learning

- Conjunctive descriptions are easily learned by finding all commonalities shared by all positive examples.

Example	Size	Color	Shape	Category
1	small	red	circle	positive
2	large	red	circle	positive
3	small	red	triangle	negative
4	large	blue	circle	negative

Learned rule: red & circle → positive

- Must check consistency with negative examples. If inconsistent, **no** conjunctive rule exists.

# Limitations of Conjunctive Rules

- If a concept does not have a single set of necessary and sufficient conditions, conjunctive learning fails.

Example	Size	Color	Shape	Category
1	small	red	circle	positive
2	large	red	circle	positive
3	small	red	triangle	negative
4	large	blue	circle	negative
5	medium	red	circle	negative

Learned rule: red & circle → positive

Inconsistent with negative example #5!

# Disjunctive Concepts

- Concept may be disjunctive.

Example	Size	Color	Shape	Category
1	small	red	circle	positive
2	large	red	circle	positive
3	small	red	triangle	negative
4	large	blue	circle	negative
5	medium	red	circle	negative

Learned rules: small & circle → positive  
large & red → positive

# Using the Generality Structure

---

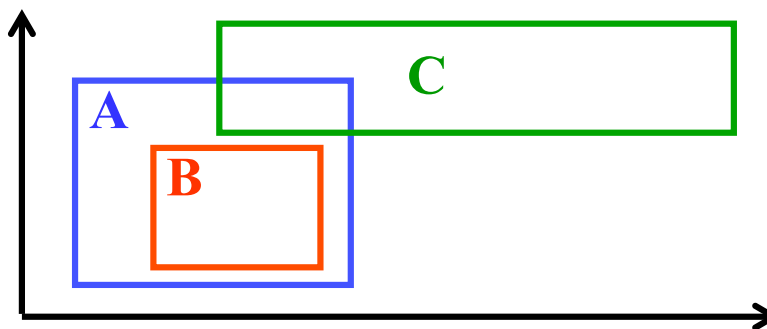
- By exploiting the structure imposed by the generality of hypotheses, an hypothesis space can be searched for consistent hypotheses without enumerating or explicitly exploring all hypotheses.
- An instance,  $x \in X$ , is said to *satisfy* an hypothesis,  $h$ , iff  $h(x)=1$  (positive)
- Given two hypotheses  $h_1$  and  $h_2$ ,  $h_1$  is *more general than or equal to*  $h_2$  ( $h_1 \geq h_2$ ) iff every instance that satisfies  $h_2$  also satisfies  $h_1$ .
- Given two hypotheses  $h_1$  and  $h_2$ ,  $h_1$  is *(strictly) more general than*  $h_2$  ( $h_1 > h_2$ ) iff  $h_1 \geq h_2$  and it is not the case that  $h_2 \geq h_1$ .
- Generality defines a partial order on hypotheses.



# Examples of Generality

---

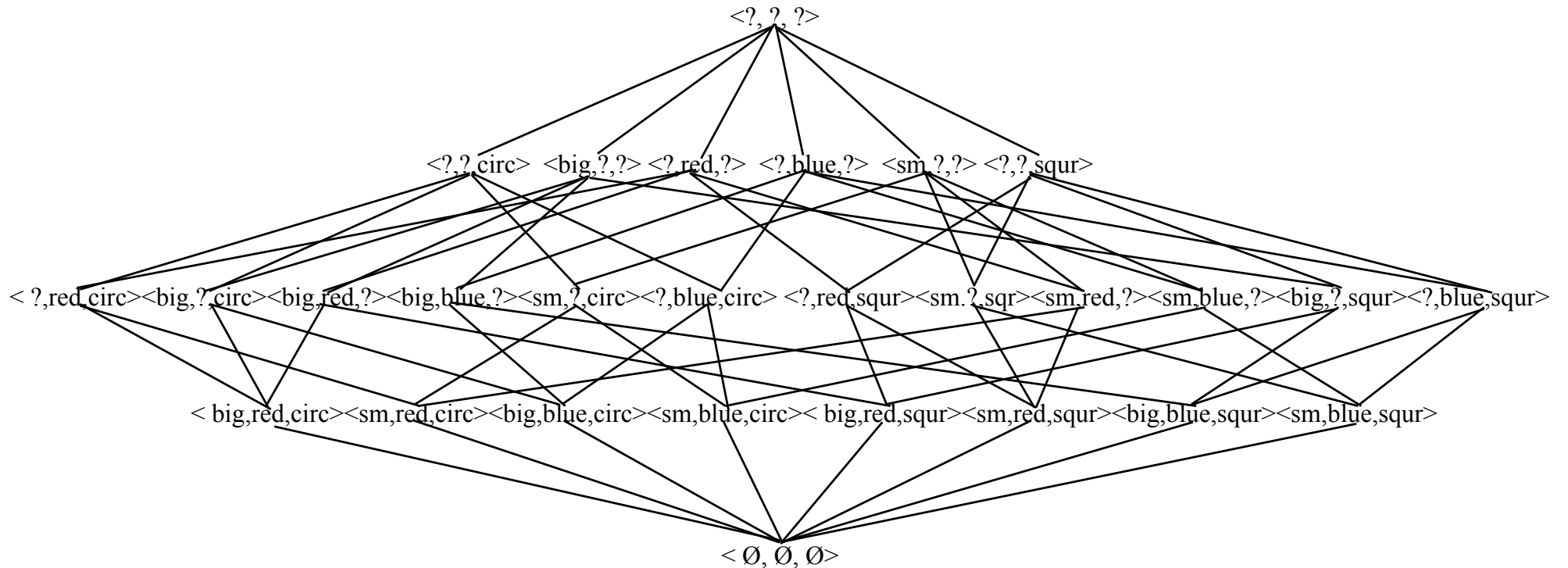
- Conjunctive feature vectors
  - $\langle ?, \text{red}, ? \rangle$  is more general than  $\langle ?, \text{red}, \text{circle} \rangle$
  - Neither of  $\langle ?, \text{red}, ? \rangle$  and  $\langle ?, ?, \text{circle} \rangle$  is more general than the other.
- Axis-parallel rectangles in 2-d space



- A is more general than B
- Neither of A and C are more general than the other.

# Sample Generalization Lattice

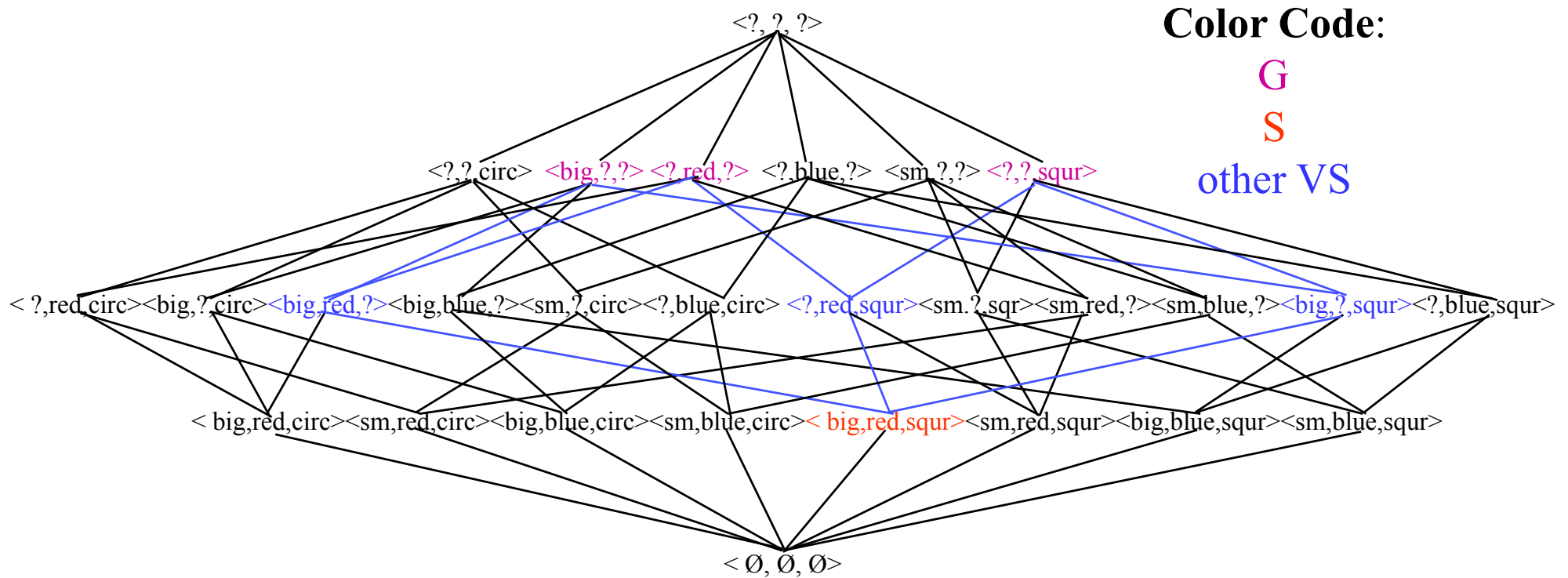
Size: {sm, big}    Color: {red, blue}    Shape: {circ, squar}



$$\text{Number of hypotheses} = 3^3 + 1 = 28$$

# Version Space Lattice

Size: {sm, big}    Color: {red, blue}    Shape: {circ, squar}



<<big, red, squar> positive>

<<sm, blue, circ> negative>

# No Panacea

---

- No Free Lunch (NFL) Theorem (Wolpert, 1995)  
Law of Conservation of Generalization Performance (Schaffer, 1994)
  - One can prove that improving generalization performance on unseen data for some tasks will always decrease performance on other tasks (which require different labels on the unseen instances).
  - Averaged across all possible target functions, no learner generalizes to unseen data any better than any other learner.
- There does not exist a learning method that is uniformly better than another for all problems.
- Given any two learning methods  $A$  and  $B$  and a training set,  $D$ , there always exists a target function for which  $A$  generalizes better (or at least as well) as  $B$ .

# Logical View of Induction

---

- Deduction is inferring sound specific conclusions from general rules (axioms) and specific facts.
- Induction is inferring general rules and theories from specific empirical data.
- Induction can be viewed as inverse deduction.
  - Find a hypothesis  $h$  from data  $D$  such that
    - $h \cup B \vdash D$   
where  $B$  is optional background knowledge
- **Abduction** is similar to induction, except it involves finding a specific hypothesis,  $h$ , that best *explains* a set of evidence,  $D$ , or inferring cause from effect. Typically, in this case  $B$  is quite large compared to induction and  $h$  is smaller and more specific to a particular event.

# Induction and the Philosophy of Science

---

- Bacon (1561-1626), Newton (1643-1727) and the sound deductive derivation of knowledge from data.
- Hume (1711-1776) and the *problem of induction*.
  - Inductive inferences can never be proven and are always subject to disconfirmation.
- Popper (1902-1994) and *falsifiability*.
  - Inductive hypotheses can only be falsified not proven, so pick hypotheses that are most subject to being falsified.
- Kuhn (1922-1996) and *paradigm shifts*.
  - Falsification is insufficient, an alternative paradigm must be available that is clearly elegant and more explanatory must be available.
    - Ptolmaic epicycles and the Copernican revolution
    - Orbit of Mercury and general relativity
    - Solar neutrino problem and neutrinos with mass
- Postmodernism: Objective truth does not exist; relativism; science is a social system of beliefs that is no more valid than others (e.g. religion).

# Ockham (Occam)'s Razor

---

- William of Ockham (1295-1349) was a Franciscan friar who applied the criteria to theology:
  - “Entities should not be multiplied beyond necessity” (Classical version but not an actual quote)
  - “The supreme goal of all theory is to make the irreducible basic elements as simple and as few as possible without having to surrender the adequate representation of a single datum of experience.” (Einstein)
- Requires a precise definition of simplicity.
- Acts as a bias which assumes that nature itself is simple.
- Role of Occam's razor in machine learning remains controversial.