



A THEORY AND METHODOLOGY OF INDUCTIVE LEARNING

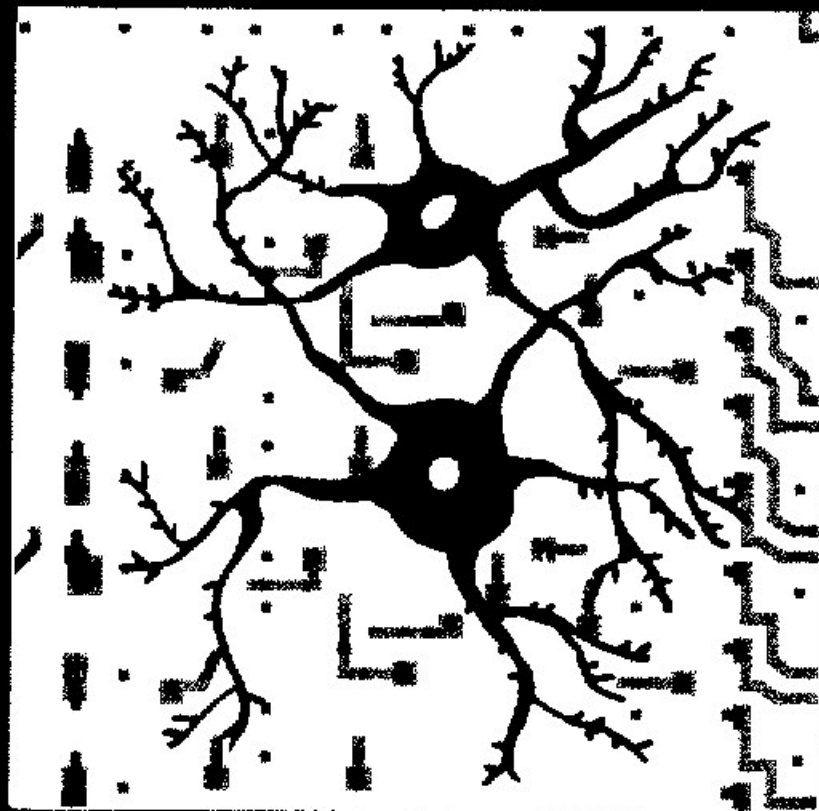
by

R. S. Michalski

Chapter In The Book, Machine Learning: An Artificial Intelligence Approach, R. S. Michalski, J. Carbonell And T. Mitchell (Eds.), Tioga Publishing Co., Pp 83-134. Palo Alto, 1983.

MACHINE LEARNING

An Artificial Intelligence Approach



Ryszard S. Michalski
Jaime G. Carbonell
Tom M. Mitchell

MACHINE LEARNING

An Artificial Intelligence Approach

MACHINE LEARNING

An Artificial Intelligence Approach

Contributing authors:

John Anderson
Ranan Banerji
Gary Bradshaw
Jaime Carbonell
Thomas Dietterich
Norman Haas
Frederick Hayes-Roth
Gary Hendrix
Patrick Langley
Douglas Lenat

Ryszard Michalski
Tom Mitchell
Jack Mostow
Bernard Nudel
Michael Rychener
Ross Quinlan
Herbert Simon
Derek Sleeman
Robert Stepp
Paul Utgoff

Editors:

Ryszard S. Michalski
*University of Illinois
at Urbana-Champaign, IL*

Jaime G. Carbonell
*Carnegie-Mellon University
Pittsburgh, PA*

Tom M. Mitchell
*Rutgers University
New Brunswick, NJ*



Library of Congress Cataloging in Publication Data

Main entry under title:

Machine learning.

Bibliography: p.

Includes index.

1. Machine learning. 2. Artificial intelligence.

I. Michalski, Ryszard Spencer, date

II. Carbonell, Jaime Guillermo, date

III. Mitchell, Tom Michael, date

Q325.M32 1983 001.53'5 82-10654

ISBN 0-935382-05-4

© 1983 by Tioga Publishing Company, P. O. Box 98, Palo Alto, CA 94302

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording, or otherwise, without the prior written permission of the publisher. Printed in the United States of America, Library of Congress Catalog Card Number 82-10654.

This book was set in Times Roman by *Fast* on a Mergenthaler Omnitech™/2000 phototypesetter driven by the Scribe™ document production system.

ISBN 0-935382-05-4
ABCDEFG-DO-876543

4

A THEORY AND METHODOLOGY OF INDUCTIVE LEARNING

Ryszard S. Michalski
*University of Illinois
at Urbana-Champaign*

ABSTRACT

The presented theory views inductive learning as a heuristic search through a space of symbolic descriptions, generated by an application of various inference rules to the initial observational statements. The inference rules include generalization rules, which perform generalizing transformations on descriptions, and conventional truth-preserving deductive rules. The application of the inference rules to descriptions is constrained by problem background knowledge, and guided by criteria evaluating the “quality” of generated inductive assertions.

Based on this theory, a general methodology for learning structural descriptions from examples, called *Star*, is described and illustrated by a problem from the area of conceptual data analysis.

4.1 INTRODUCTION

“...Scientific knowledge through demonstration¹ is impossible unless a man knows the primary immediate premises ... We must get to know the primary premises by induction; for the method by which even sense-perception implants the universal is inductive...”—Aristotle, Posterior Analytics, Book II, Chapter 19 (circa 330 B.C.)

The ability of people to make accurate generalizations from a few scattered

¹That is, what we now call “deduction”.

facts or to discover patterns in seemingly chaotic collections of observations is a fascinating research topic of long-standing interest. The understanding of this ability is now also of growing practical importance, as it holds the key to an improvement of methods by which computers can acquire knowledge. A need for such an improvement is evidenced by the fact that knowledge acquisition is presently the most limiting "bottleneck" in the development of modern knowledge-intensive artificial intelligence systems.

The above ability is achieved by a process called *inductive learning*, that is, inductive inference from facts provided by a teacher or the environment. The study and modeling of this form of learning is one of the central topics of machine learning. This chapter outlines a theory of inductive learning and then presents a methodology for acquiring general concepts from examples.

Before going further into this topic, let us first discuss the potential for applications of inductive learning systems. One such application is an automated construction of knowledge bases for expert systems. The present approach to constructing knowledge bases involves a tedious process of formalizing experts' knowledge and encoding it in some knowledge representation system, such as production rules [Shortliffe, 1976; Davis & Lenat, 1981] or a semantic network [Brachman, 1979; Gaschnig, 1980]. Inductive learning programs could provide both an improvement of the current techniques and a basis for developing alternative knowledge acquisition methods.

In appropriately selected small domains, inductive programs are already able to determine decision rules by induction from examples of expert decisions. This process greatly simplifies the transfer of knowledge from an expert into a machine. The feasibility of such inductive knowledge acquisition has been demonstrated in the expert system PLANT/DS, for the diagnosis of soybean diseases. In this system, the diagnostic rules were developed in two ways: by formalizing experts' diagnostic processes and by induction from examples. In an experiment where both types of diagnostic rules were tested on a few hundred disease cases, the inductively-derived rules outperformed the expert-derived ones [Michalski & Chilausky, 1980]. Another example is an inductive acquisition of decision rules for a chess end-game [Michalski & Negri, 1977; Quinlan, 1979; Shapiro & Niblett, 1982; O'Rorke, 1982]. (See also Chapter 15 of this book.)

A less direct, but potentially promising, use of inductive learning is for the refinement of knowledge-bases initially developed by human experts. Here, inductive learning programs could be used to detect and rectify inconsistencies, to remove redundancies, to cover gaps, and to simplify expert-derived decision rules. By applying an inductive inference program to the data, consisting of original rules and examples of correct and incorrect results of these rules' performance in new situations, the rules could be incrementally improved with little or no human assistance.

Another important application of inductive programs is in various experimental sciences, such as biology, chemistry, psychology, medicine, and genetics. Here they could assist a user in detecting interesting conceptual pat-

terns or in revealing structure in collections of observations. The widely used traditional mathematical and statistical data analysis techniques, such as regression analysis, numerical taxonomy, or factor analysis, are not sufficiently powerful for this task. Methods for *conceptual* data analysis are needed, that generate not merely mathematical formulas but logic-style descriptions, characterizing data in terms of high-level, human-oriented concepts and relationships. An early example of such an application is the META-DENDRAL program [Buchanan & Feigenbaum, 1978], which infers cleavage rules for mass-spectrometer simulation. (See its analysis in Chapter 3 of this book.)

There are two basic modes in which inductive programs can be utilized: as interactive tools for acquisition of knowledge from specific facts or examples, or as parts of machine-learning systems. In the first mode, a user supplies learning examples and exercises strong control over the way the program is used (for example, [Michalski, 1975a; Quinlan, 1979; Michalski & Chilausky, 1980] and Chapter 15 of this book).

In the second mode, an inductive program is a component of an integrated learning system whose other components generate the needed learning examples [Buchanan *et al.*, 1979]. Such examples—positive and negative—constitute the feedback from the system's attempts to perform a desired task. An example of the second mode is the learning system LEX for symbolic integration (see Chapter 6 of this book), where a "generalizer" module performs inductive inference on instances provided by a "critic" module. Another example is discussed in Chapter 5 of this book, in the context of analogy-based learning.

From the viewpoint of applications, such as aiding the construction of expert systems or conceptual analysis of experimental data, the most relevant is *conceptual inductive learning*. We use this term to designate a type of inductive learning whose final products are symbolic descriptions expressed in high-level, human-oriented terms and forms (more details are given in Section 4.3.1). The descriptions typically apply to real world objects or phenomena, rather than abstract mathematical concepts or computations. This paper is concerned specifically with conceptual inductive learning.

The most frequently studied type of such learning is *concept learning from examples* (called also *concept acquisition*), whose task is to induce general descriptions of concepts from specific instances of these concepts. The early studies of this subject go back to the fifties, for example, [Hovland, 1952; Bruner *et al.*, 1956; Newell *et al.*, 1960; Amarel, 1960; Feigenbaum, 1963; Kochen, 1960; Banerji, 1962; Simon & Kotovsky, 1963; Hunt *et al.*, 1966; Hájek *et al.*, 1966; Bongard, 1970]. Among more recent contributions there are, for instance, [Winston, 1970; Waterman, 1970; Michalski, 1972; Hayes-Roth, 1973; Simon & Lea, 1974; Stoffel, 1974; Vere, 1975; Larson & Michalski, 1977; Mitchell, 1978; Quinlan, 1979; Moraga, 1981]. An important variant of concept learning from examples is the *incremental concept refinement*, where the input information includes, in addition to the training examples, previously-learned hypotheses, or human-provided initial hypotheses that may be partially

incorrect or incomplete [Michalski & Larson, 1978]. Chapter 3 of this book discusses various evaluation criteria and several methods for concept learning from examples.

Another type of conceptual inductive learning is *concept learning from observation* (or *descriptive generalization*), concerned with establishing new concepts or theories characterizing given facts. This area includes such topics as automated theory formation (for example, [Lenat, 1976] and Chapter 9 of this book), discovery of relationships in data (for example, [Hájek & Havránek, 1978; Pokorny, 1980; Zagoruiċo, 1981] and Chapter 10 of this book), or an automatic construction of taxonomies (for example, Chapter 11 of this book). Differences between concept learning from examples and concept learning from observation are discussed in more detail in the next section.

Conceptual inductive learning has a strong cognitive science flavor. Its emphasis on inducing human-oriented, rather than machine-oriented descriptions, and its primary interest in nonmathematical domains distinguishes it from other types of inductive learning, such as grammatical inference and program synthesis. In grammatical inference, the task is to determine a formal grammar that can generate a given set of symbol strings (for example, [Solomonoff, 1964; Biermann & Feldman, 1972; Yau & Fu, 1978; Gaines, 1979]). In program synthesis the objective is to construct a computer program from I/O pairs or computational traces, or to transform a program from one form to another by applying correctness-preserving transformation rules (for example, [Shaw *et al.*, 1975; Burstall & Darlington, 1977; Case & Smith, 1981; Biermann, 1978; Jouannaud & Kodratoff, 1980; Smith, 1980; Pettorossi, 1980]). The final result of such learning is a computer program, in an assumed programming language, destined for machine rather than human "consumption". For example, the method of "model inference" by Shapiro [1981] constructs a PROLOG program characterizing a given set of mathematical facts.

Recent years have witnessed the development of a number of task-oriented inductive learning systems that have demonstrated an impressive performance in their specific domain of application. Major weaknesses, however, persist in much of the research in this area. Most systems lack generality and extensibility. The theoretical principles upon which they are built are rarely well explained. Lack of common terminology and an adequate formal theory makes it difficult to compare different learning methods.

In the following sections we formulate logical foundations of inductive learning, define various types of such learning, present inference rules for generalizing concept descriptions, and finally describe a general methodology, called Star, for learning structural descriptions from examples. To improve the readability of this chapter, Table 4-1 provides a list of basic symbols used, with a short explanation. The Appendix gives the details of the description language used (the annotated predicate calculus).

Table 4-1: A Table of Basic Symbols

| | |
|-------------------|--|
| \sim | negation |
| $\&$ | conjunction (logical product) |
| \vee | disjunction (logical sum) |
| \Rightarrow | implication |
| \Leftrightarrow | logical equivalence |
| \leftrightarrow | term rewriting |
| \setminus | exception (symmetric difference) |
| F | a set of facts (formally, a predicate that is true for all the facts) |
| H | a hypothesis (an inductive assertion) |
| $ >$ | specialization |
| $ <$ | generalization |
| $:=$ | reformulation |
| $\exists v_i$ | existential quantifier over v_i |
| $\exists(I) v_i$ | numerical quantifier over v_i (I is a set of integers) |
| $\forall v_i$ | universal quantifier over v_i |
| D_i | a concept description |
| K_i | a predicate asserting the name of a concept (a class) |
| $::>$ | the implication linking a concept description with a concept name |
| e_i | an event (a description of an object or a situation) |
| E_i | a predicate that is true only for the training events of concept K_i |
| x_i | an attribute (zero- or one-argument descriptor) |
| LEF | a lexicographic evaluation functional |
| DOM(p) | the domain of descriptor p . |

4.2 TYPES OF INDUCTIVE LEARNING

4.2.1 Inductive Paradigm

As mentioned before, inductive learning is a process of acquiring knowledge by drawing inductive inferences from teacher- or environment- provided facts. Such a process involves operations of generalizing, specializing, transforming, correcting and refining knowledge representations. Although it is one of the most common forms of learning, it has one fundamental weakness: except for special cases, the acquired knowledge cannot, in principle, be completely validated. This predicament, observed by the Scottish philosopher David Hume in the 18th century, is due to the fact that inductively-acquired assertions are hypotheses with a potentially infinite number of consequences, while only a finite number of confirming tests can be performed.

Traditional inquiries into inductive inference have therefore dealt with questions of what are the best criteria for guiding the selection of inductive assertions, and how these assertions can be confirmed. These are difficult problems, permeating all scientific activities. The search for answers has turned inductive inference into a battlefield of philosophers and logicians. There was even doubt whether it would ever be possible to formalize inductive inference and perform it on a machine. For example, philosopher Karl Popper [1968] believed that inductive inference requires an irrational element. Bertrand Russell [1946] stated:

“So far no method has been found which would make it possible to invent hypotheses by rule.” George Polya [1954] in his pioneering and now classic treatise on plausible inference (of which inductive inference is a special case) observed: “A person has a background, a machine has not; indeed, you can build a machine to draw demonstrative conclusions for you, but I think you can never build a machine that will draw plausible inferences.”

The above pessimistic prospects are now being revised. With the development of modern computers and subsequent advances in artificial intelligence research, it is now possible to provide a machine with a significant amount of background information. Also, the problem of automating inductive inference can be simplified by concentrating on the subject of hypothesis generation, while ascribing to humans the question of how to adequately validate them. Some successful inductive inference systems have already been built and a body of knowledge is emerging about the nature of this inference. The rest of this section will analyze the logical basis for inductive inference, and then Section 4.5 will present various generalization rules, which can be viewed as inductive inference rules.

In contrast to deduction, the starting premises of induction are specific facts rather than general axioms. The goal of inference is to formulate plausible general assertions that explain the given facts and are able to predict new facts. In other words, inductive inference attempts to derive a complete and correct description of a given phenomenon from specific observations of that phenomenon or of parts of it. As mentioned earlier, of the two aspects of inductive inference—the generation of plausible hypotheses and their validation (the establishment of their truth status)—only the first is of primary interest to inductive learning research. The problem of hypothesis validation, a subject of various philosophical inquiries (for example, [Carnap, 1962]) is considered to be of lesser importance, because it is assumed that the generated hypotheses are judged by human experts, and tested by known methods of deductive inference and statistics.

As described in Chapter 1 of this book, there are several different methods by which a human (or a machine) can acquire knowledge, such as rote learning (or learning by being programmed), learning from instruction (or learning by being told), learning from teacher-provided examples (concept acquisition), and learning by observing the environment and making discoveries (learning from observation and discovery).

Although all of these ways except the first involve some amount of inductive inference, in the last two, that is, in learning from examples and in learning from observation, this inference is the central operation. These two forms are therefore considered to be the major forms of inductive learning. In order to explain them, let us formulate a general paradigm for inductive inference:

Given:

- *Observational statements* (facts), F , that represent specific knowledge about some objects, situations, processes, and so on,

- A *tentative inductive assertion* (which may be null),
- *Background knowledge* that defines the assumptions and constraints imposed on the observational statements and generated candidate inductive assertions, and any relevant problem domain knowledge. The last includes the *preference criterion* characterizing the desirable properties of the sought inductive assertion.

Find:

- An *inductive assertion* (hypothesis), H , that tautologically or weakly implies the observational statements, and satisfies the background knowledge.

A hypothesis H tautologically implies facts F if F is a logical consequence of H , that is, if the expression $H \Rightarrow F$ is true under all interpretations (" \Rightarrow " denotes logical implication). This is expressed as follows:

$$H \mid > F \text{ (read: } H \text{ specializes to } F) \quad (1)$$

or

$$F \mid < H \text{ (read: } F \text{ generalizes to } H) \quad (2)$$

Symbols $\mid >$ and $\mid <$ are called the *specialization* and *generalization* symbols, respectively. If $H \Rightarrow F$ is valid, and H is true, then by the law of detachment (*modus ponens*) F must be true. Deriving F from H (deductive inference), is, therefore, truth-preserving. In contrast, deriving H from F (inductive inference) is not truth-preserving, but falsity-preserving; that is, if some facts falsify F , then they also must falsify H . (More explanation on this topic is given in Section 4.5.)

The condition that H *weakly implies* F means that facts F are not certain but only plausible or partial consequences of H . By allowing weak implication, this paradigm includes methods for generating "soft" hypotheses, which hold only probabilistically, and partial hypotheses, which account for some but not all of the facts (for example, hypotheses representing "dominant patterns" or characterizing inconsistent data). In the following we will limit our attention to hypotheses that tautologically imply facts.

For any given set of facts, a potentially infinite number of hypotheses can be generated that imply these facts. Background knowledge is therefore necessary to provide the constraints and a preference criterion for reducing the infinite choice to one hypothesis or a few most preferable ones.

A typical way of defining such a criterion is to specify the preferable properties of the hypothesis—for example, to require that the hypothesis is the shortest or the most economical description consistent with all the facts (as, for example, in [Michalski, 1973]). Such a "biased-choice" criterion is necessary when the description language is complete, that is, able to express any possible hypothesis. An alternative is to use a "biased-language" criterion [Mitchell, 1978], restricting the description language in which hypotheses are expressed

(that is, to use an incomplete description language). Although in many methods the background knowledge is not explicitly stated, the authors make implicit assumptions serving the same purpose. More details on the criteria for selecting hypotheses are given in Section 4.4.7.

4.2.2 Concept Acquisition versus Descriptive Generalization

As mentioned in the Introduction, one can distinguish between two major types of inductive learning: *learning from examples* (concept acquisition) and *learning from observation* (descriptive generalization). In concept acquisition, the observational statements are characterizations of some objects (situations, processes, and so on) preclassified by a teacher into one or more classes (concepts). The induced hypothesis can be viewed as a concept recognition rule, such that if an object satisfies this rule, then it represents the given concept. For example, a recognition rule for the concept “philosopher” might be:

“A person who pursues wisdom and gains the knowledge of underlying reality by intellectual means and moral self-discipline is a philosopher.”

In descriptive generalization the goal is to determine a general description (a law, a theory) characterizing a collection of observations. For example, observing that the philosophers Aristotle, Plato, and Socrates were Greek, but that Spencer was British, one might conclude:

“Most philosophers were Greek.”

Thus, in contrast to concept acquisition that produces descriptions for classifying objects into classes on the basis of the objects’ properties, descriptive generalization produces descriptions specifying properties of objects belonging to a certain class. Here are some example problems belonging to the above two categories:

1. Concept Acquisition:

- Learning a *characteristic description* of a class of objects, that specifies all common properties of known objects in the class, and by that defines the class in the context of an unlimited number of other object classes (for example, [Bongard, 1967; Winston, 1970; Stoffel, 1974; Vere, 1975; Cohen, 1977; Hayes-Roth & McDermott, 1978; Mitchell, 1978; Stepp, 1978; Michalski, 1980a] and Chapter 3 of this book).
- Learning a *discriminant description* of a class of objects that distinguishes the given class from a limited number of other classes (for example, [Michalski, 1973; Quinlan, 1979; Michalski, 1980a] and Chapter 15 of this book).
- Inferring *sequence extrapolation rules* (for example, [Simon & Kotovsky, 1963; Dietterich, 1979]) able to predict the next element (a symbol, a number, an object, and so on) in a given sequence.

2. Descriptive Generalization:

- Formulating a theory characterizing a collection of entities (for example, a number theory, as in [Lenat, 1976] and Chapter 9 of this book).
- Discovering patterns in observational data (for example, [Soloway & Riseman, 1977; Hájek & Havránek, 1978; Pokorný, 1980; Zago-ruiko, 1981] and Chapter 10 of this book).
- Determining a taxonomic description (classification) of a collection of objects (for example, [Michalski, 1980c; Michalski *et al.*, 1981] and Chapter 11 of this book).

This paper is concerned primarily with problems of concept acquisition. In this case, the set of observational statements F can be viewed as a collection of implications:

$$F : \{e_{ik} ::> K_i\}, i \in I \quad (3)$$

where e_{ik} (a *training event*) denotes a description of the k^{th} example of *concept* (class) asserted by predicate K_i (for short, class K_i) and I is a set indexing classes K_i . It is assumed here that any given event represents only one concept. Symbol $::>$ is used here, and will be used henceforth, to denote the implication linking a concept description with a predicate asserting the concept name (in order to distinguish this implication from the implication between arbitrary descriptions). The inductive assertion H can be characterized as a set of concept recognition rules:

$$H : \{D_i ::> K_i\}, i \in I \quad (4)$$

where D_i is a concept description of class K_i , that is, an expression of conditions, such that when they are satisfied by an object, the object is considered an instance of class K_i .

According to the definition of inductive assertion, we must have:

$$H \mid > F \quad (5)$$

By substituting (3) and (4) for F and H , respectively, in (5), and making appropriate transformations, one can derive the following conditions to be satisfied in order that (5) holds:

$$\forall i \in I (E_i \Rightarrow D_i) \quad (6)$$

and

$$\forall i, j \in I (D_i \Rightarrow \sim E_j), \text{ if } j \neq i \quad (7)$$

where E_i , $i \in I$, is a description satisfied by all training events of class K_i , and only by such events (the logical disjunction of training events).

Expression (6) is called the *completeness condition*, and (7) the *consistency condition*. These two conditions are the requirements that must be satisfied for an inductive assertion to be acceptable as a concept recognition rule. The com-

pleteness condition states that every training event of some class must satisfy the description D_i of the same class (since the opposite does not have to hold, D_i is equivalent to, or more general than, E_i). The consistency condition states that if an event satisfies a description of some class, then it cannot be a member of a training set of any other class. In learning a concept from examples and counter-examples, the latter constitute the "other" class.

The completeness and consistency conditions provide the logical foundation of algorithms for concept learning from examples. We will see in Section 4.5 that to derive D_i satisfying the completeness condition one can adopt some inference rules of formal logic.

4.2.3 Characteristic versus Discriminant Descriptions

The completeness and consistency conditions allow us to clearly explain the distinction between the previously mentioned characteristic and discriminant descriptions. A characteristic description of a class of objects (also known as *conjunctive generalization*) is an expression that satisfies the completeness condition or is the logical product of such expressions. It is typically a conjunction of some simple properties common to all objects in the class. From the applications viewpoint, the most interesting are *maximal characteristic descriptions* (maximal conjunctive generalizations) that are the most specific (that is, the longest) logical products characterizing all objects in the given class, using terms of the given language. Such descriptions are intended to discriminate the given class from all other possible classes (for illustration see Section 4.7.2).

A discriminant description is an expression that satisfies the completeness and consistency condition, or is the logical disjunction of such expressions. It specifies one or more ways to distinguish the given class from a fixed number of other classes. The most interesting are *minimal discriminant descriptions* that are the shortest (that is, have the minimum number of descriptors) expressions distinguishing all objects in the given class from objects of the other classes. Such descriptions are intended to specify the minimum information sufficient to identify the given class among a fixed number of other classes (for illustration see Section 4.7.1).

4.2.4 Single- versus Multiple-concept Learning

It is instructive to distinguish between learning a single concept, and learning a collection of concepts. In *single concept learning*, one can distinguish two cases: (i) when observational statements are just examples of the concept to be learned (learning from "positive" instances only); and (ii) when they are examples and counter-examples of the concept (learning from "positive" and "negative" instances).

In the first case, because of the lack of counter-examples, the consistency condition (7) is not applicable, and there is no natural limit to which description D_i (here, $i = 1$) can be generalized. One way to impose such a limit is to specify

restrictions on the form and properties of the sought description. For example, one may require that it be the maximal characteristic description, that is, the longest conjunctive statement satisfying the completeness condition (for example, [Vere, 1975; Hayes-Roth & McDermott, 1978]). Another way is to require that the description not exceed a given degree of generality, measured, for example, by the ratio of the number of all distinct events which could potentially satisfy the description to the number of training instances [Stepp, 1978].

In the second case, when the teacher also provides counter-examples of the given concept, the learning process is considerably simplified. These counter-examples can be viewed as representing a "different class", and the consistency condition (7) provides an obvious limit on the extent to which a hypothesis can be generalized. The most useful counter-examples are the so-called "near misses" that only slightly differ from positive examples [Winston, 1970, 1977]. Such examples place stronger constraints on the generalization process than randomly-generated examples.

In *multiple-concept learning* one can also distinguish two cases: (i) when descriptions D_i of different classes are required to be mutually disjoint, that is, no event can satisfy more than one description; and (ii) when they are overlapping. In an overlapping generalization an event may satisfy more than one description. In some situations this is desirable. For example, if a patient has two diseases, his symptoms should satisfy the descriptions of both diseases, and in this case the consistency condition is not applicable.

An overlapping generalization can be interpreted in such a way that it always indicates only one decision class. For example, the concept recognition rules, $D_i ::> K_i$, can be applied in a linear order, and the first rule satisfied generates the decision. In this case, if a concept description D_i for class K_i contains a conjunctively-linked condition A , and precedes the rule for class K_j that contains condition $\sim A$, then the condition $\sim A$ is superfluous and can be removed. As a result, the linearly-ordered recognition rules can be significantly simplified. For example, the set of linearly-ordered rules:

$$\begin{aligned} D_1 &::> K_1 \\ D_2 &::> K_2 \\ D_3 &::> K_3 \end{aligned}$$

is logically equivalent to the set of (unordered) rules:

$$\begin{aligned} D_1 &::> K_1 \\ \sim D_1 \& D_2 &::> K_2 \\ \sim D_1 \& \sim D_2 \& D_3 &::> K_3 \end{aligned}$$

There are also other ways to derive a single decision from overlapping rules, such as those given in [Davis & Lenat, 1981].

The above forms of multiple-concept learning have been implemented in inductive programs AQVAL/1 [Michalski, 1973] and AQ11 [Michalski & Larson, 1978].

4.3 DESCRIPTION LANGUAGE

4.3.1 Bias Toward Comprehensibility

In concept acquisition, the main interest is in derivation of symbolic descriptions that are human-oriented, that is, that are easy to understand and easy to use for creating mental models of the information they convey. A tentative criterion for judging inductive assertions from such a viewpoint is provided by the following *comprehensibility postulate*:

The results of computer induction should be symbolic descriptions of given entities, semantically and structurally similar to those a human expert might produce observing the same entities. Components of these descriptions should be comprehensible as single "chunks" of information, directly interpretable in natural language, and should relate quantitative and qualitative concepts in an integrated fashion.

As a practical guide, one can assume that the components of descriptions (single sentences, rules, labels on nodes in a hierarchy, and so on) should be expressions that contain only a few (say, less than five) conditions in a conjunction, few single conditions in a disjunction, at most one level of bracketing, at most one implication, no more than two quantifiers, and no recursion (the exact numbers may be disputed,² but the principle is clear). Sentences are kept within such limits by substituting names for appropriate subcomponents. Any operators used in descriptions should have a simple intuitive interpretation. Conceptually related sentences are organized into a simple data structure, preferably a shallow hierarchy or a linear list, such as a frame [Minsky, 1975]. (See also Chapter 9 of this book.)

The rationale behind this postulate is to ensure that descriptions generated by inductive inference bear similarity to human knowledge representations [Hintzman, 1978], and therefore, are easy to comprehend. This requirement is very important for many applications. For example, in developing knowledge bases for expert systems, it is important that human experts can easily and reliably verify the inductive assertions and relate them to their own domain knowledge. Satisfying the comprehensibility postulate will also facilitate debugging or improving the inductive programs themselves. When the complexity of problems undertaken by computer induction becomes very great, the comprehensibility of the generated descriptions will likely be a crucial criterion. This research orientation fits well within the role of artificial intelligence envisaged by Michie [1977] to study and develop methods for man-machine conceptual interface and knowledge refinement.

²The numbers mentioned seem to apply to the majority of human descriptive sentences.

4.3.2 Language of Assertions

One of the difficulties with inductive inference is its open-endedness. This means that when one makes an inductive assertion about some aspect of reality there is no natural limit to the level of detail in which this reality may be described, or to the richness of forms in which this assertion can be expressed. Consequently, when conducting research in this area, it is necessary to circumscribe very carefully the goals and the problem to be solved. This includes defining the language and the scope of allowed forms in which assertions will be expressed, as well as the modes of inference which will be used. The description language should be chosen so that crucial features can be easily encoded while peripheral or irrelevant information ignored.

An instructive criterion for classifying inductive learning methods is therefore the type of language used to express inductive assertions. Many authors use a restricted form of predicate calculus or closely related notation (for example, [Plotkin, 1971; Fikes *et al.*, 1972; Morgan, 1975; Vere, 1975; Banerji, 1980; Michalski, 1980a; Sammut, 1981; Zagoruiko, 1981]). Some other formalisms include decision trees [Hunt *et al.*, 1966; Quinlan, 1979] (see also Chapter 15 of this book), production rules (for example, [Waterman, 1970; Hedrick, 1974] (see also Chapter 16 of this book), semantic nets (Chapter 13), and frames (Chapter 9). In his earlier work (for example, [Michalski, 1972, 1973, 1975a, 1975b]) this author used a multiple-valued logic propositional calculus with typed variables, called VL_1 (the variable-valued logic system one). Later on an extension of the predicate calculus, called VL_2 , was developed, that was especially oriented to facilitate inductive inference [Michalski, 1980a].

Here we will use a somewhat modified and extended version of the latter language, to be called the *annotated predicate calculus* (APC). The APC adds to predicate calculus additional forms and new concepts that increase its expressive power and facilitate inductive inference. The major differences between the annotated predicate calculus and the conventional predicate calculus can be summarized as follows:

1. Each predicate, variable and function (referred to collectively as a *descriptor*) is assigned an *annotation* that contains relevant problem-oriented information. The annotation may contain the definition of the concept represented by a descriptor, a characterization of its relationship to other concepts, a specification of the set over which the descriptor ranges (when it is a variable or a function), a characterization of the structure of this set, and so on (see Section 4.4).
2. In addition to predicates, the APC also includes *compound predicates*. Arguments of such predicates can be *compound terms*, composed of two or more ordinary terms.
3. Predicates that express relations $=$, \neq , \geq , $>$, \leq and $<$ between terms or between compound terms are expressed explicitly as *relational statements*, also called *selectors*.

4. In addition to the universal and existential quantifiers, there is also a *numerical quantifier* that expresses quantitative information about the objects satisfying an expression.

The concept of annotation is explained in more detail in the next section. Other aspects of the language are described in the Appendix. (The reader interested in a thorough understanding of this work is encouraged to read the Appendix at this point.)

4.4 PROBLEM BACKGROUND KNOWLEDGE

4.4.1 Basic Components

As we mentioned earlier, given a set of observational statements, one may construct a potentially infinite number of inductive assertions that imply these statements. It is therefore necessary to use some additional information, *problem background knowledge*, to constrain the space of possible inductive assertions and locate the most desirable one(s). In this section, we shall look at various components of the problem background knowledge employed in the inductive learning methodology called Star, described in Section 4.6. These components include:

- Information about descriptors (i.e., predicates, variables, or functions) used in observational statements. This information is provided by an annotation assigned to each descriptor (Section 4.4.3).
- Assumptions about the form of observational and inductive assertions.
- A preference criterion that specifies the desirable properties of inductive assertions sought.
- A variety of inference rules, heuristics, and specialized procedures, general and problem-dependent, that allow a learning system to generate logical consequences of given assertions and new descriptors.

Before we examine these components in greater detail, let us first consider the problem of how the choice of descriptors in the observational statements affects the generated inductive assertions.

4.4.2 Relevance of the Initial Descriptors

A fundamental problem underlying any machine inductive learning task is that of what information is provided to the machine and what information the machine is expected to produce or learn. As specified in the inductive paradigm, the major component of the input to a learning system is a set of observational statements. The descriptors used in those statements are observable characteristics and available measurements of objects under consideration. These descriptors are selected as relevant to the learning task by a teacher specifying the problem.

Determining these descriptors is a major part of any inductive learning problem. If they capture the essential properties of the objects, the role of the learning process is simply to arrange these descriptors into an expression constituting an appropriate inductive assertion. If the selected descriptors are completely irrelevant to the learning task (as the color, weight, or shape of men in chess is irrelevant to deciding the right move), no learning system will be able to construct a meaningful inductive assertion.

There is a range of intermediate possibilities between the above two extremes. Consequently, learning methods can be characterized on the basis of the degree to which the initial descriptors are relevant to the learning problem.

Three cases can be distinguished:

1. **Complete relevance**—In this case all descriptors in the observational statements are assumed to be directly relevant to the learning task. The task of the learning system is to formulate an inductive assertion that is a mathematical or logical expression of some assumed general form that properly relates these descriptors (for example, a regression polynomial).
2. **Partial relevance**—Observational statements may contain a large number of irrelevant or redundant descriptors. Some of the descriptors, however, are relevant. The task of the learning system is to select the most relevant ones and construct from them an appropriate inductive assertion.
3. **Indirect relevance**—Observational statements may contain no directly-relevant descriptors. However, among the initial descriptors there are some that can be used to construct derived descriptors that are directly relevant. The task of the learning system is to construct those derived descriptors and formulate an appropriate inductive assertion. A simple form of this case occurs, for example, when a relevant descriptor is the volume of an object, but the observational statements contain only the information about the object's dimensions (and various irrelevant facts).

The above three cases represent problem statements that put progressively less demand on the relevance of the initial descriptors (that is, that require less work from the person defining the problem) and more demand on the learning system. Early work on adaptive control systems and concept formation represents case 1. More recent research has dealt with case 2, which is addressed in *selective inductive learning*. A method of such learning must possess efficient mechanisms for determining combinations of descriptors that are relevant and sufficient for the learning task. Formal logic provides such mechanisms, and therefore it has become the major underlying formalism for selective methods.

An example of a selective learning method is the one implemented in program AQ11 [Michalski & Larson, 1978] that inductively determined soybean disease diagnostic rules for the system PLANT/DS, mentioned in the Introduction. A different type of selective method was implemented in program ID3 (Chapter 15) that determines a decision tree for classifying a large number of events. A comparison between these two programs is described by O'Rorke [1982].

Case 3 represents the task of *constructive inductive learning*. Here, a method must be capable of formulating new descriptors (that is, new concepts, new variables, and the like), of evaluating their relevance to the learning task and of using them to construct inductive assertions. There has been relatively little done in this area. The “automated mathematician” program AM (Chapter 9) can be classified as a domain-specific system of this category. Some constructive learning capabilities have been incorporated in system BACON that automatically formulates mathematical expressions encapsulating chemical and other laws [Langley *et al.*, 1980] (see also Chapter 10). The general-purpose INDUCE program for learning structural descriptions from examples incorporates several constructive generalization techniques [Larson, 1977; Michalski, 1980a]. Sections 4.5 and 4.6 give more details on this subject.

4.4.3 Annotation of Descriptors

An *annotation* of a descriptor (that is, of a predicate, variable, or function) is a store of background information about this descriptor tailored to the learning problem under consideration. It may include:

- A specification of the *domain* and the *type* of the descriptor (see below).
- A specification of operators applicable to it.
- A specification of the constraints and the relationships between the descriptor and other descriptors.
- For numerical descriptors, the mean, the variance, or the complete probability distribution of values for the problem under consideration.
- A characterization of objects to which the descriptor is applicable (such as a characterization of its possible arguments).
- A specification of a descriptor class containing the given descriptor, that is, the parent node in a generalization hierarchy of descriptors (for example, for descriptors “length”, “width”, and “height”, the parent node would be the “dimensions”).
- Synonyms that can be used to denote the descriptor.
- A definition of a descriptor (when it is derived from some other descriptors).
- If a descriptor denotes a class of objects, typical examples of this class can be specified.

Let us consider some of the above components of the annotation in greater detail.

4.4.4 The Domain and Type of a Descriptor

Given a specific problem, it is usually possible to specify the set of values each descriptor could potentially adopt in characterizing any object in the population under consideration. Such a set is called the *domain* (or the *value set*) of

the descriptor. The domain is used to constrain the extent to which a descriptor can be generalized. For example, the information that the temperature of a living human being may vary, say, only between 34°C and 44°C prevents the system from considering inductive assertions in which the descriptor "body temperature" would assume values beyond these limits.

Other important information for conducting the generalization process is concerned with the structure of the domain, that is, with the relationship existing among the elements of the domain. For numerical descriptors, such relationships are specified by the measurement scale. Depending on the structure of the descriptor domain, we distinguish among three basic types of descriptors:

1. **Nominal (categorical) descriptors**—The value set of such descriptors consists of independent symbols or names, that is, no structure is assumed to relate the values in the domain. For example, "blood-type(person)" and "name(person)" are unary nominal descriptors. Predicates, that is, descriptors with the value set {True, False}, and n-ary functions whose ranges are unordered sets, are also nominal descriptors. An example of a two-argument nominal descriptor is "license-plate-number(car, owner)", which denotes a function assigning to a specific car of the given owner a license plate number.
2. **Linear descriptors**—The value set of linear descriptors is a totally ordered set. For example, a person's military rank or the temperature, weight, or number of items in a set is such a descriptor. Variables measured on ordinal, interval, ratio, and absolute scales are special cases of a linear descriptor. Functions that map a set into a totally-ordered set are also linear descriptors, for example, "distance(P_1, P_2)".
3. **Structured descriptors**—The value set of such descriptors has a tree-oriented graph structure that reflects the generalization relation between the values, that is, is a *generalization hierarchy*. A parent node in such a structure represents a more general concept than the concepts represented by its children nodes. For example, in the value set of descriptor "place", "U.S.A." would be a parent node of the nodes "Indiana", "Illinois", "Iowa", and so on. The domain of structured descriptors is defined by a set of inference rules specified in the problem background knowledge (see, for example, descriptor "shape(B_i)" in Section 4.7).

Structured descriptors can be further subdivided into ordered and unordered structured descriptors (see Chapter 11).

Sometimes, descriptors themselves can also be organized into a generalization hierarchy. For example, as already mentioned, the descriptors "length", "width", and "depth" belong to a class of "dimensions". Information about the type of a descriptor is useful as it determines the operations applicable to a descriptor.

4.4.5 Constraints on the Description Space

For a given induction problem there may exist a variety of constraints on the space of the acceptable concept descriptions, due to the specific properties and relationships among descriptors. Here are a few examples of such relationships:

- **Interdependence among values**—In many practical problems some variables specify a state of an object, and some other variables characterize the state. Depending on the values of the state-specifying variables, the variables characterizing a state may or may not be needed. For example, if a descriptor “state(plant’s leaf)” takes on value “diseased”, then a descriptor “leaf discoloration” will be used to characterize the change of the leaf’s color. When the descriptor “state(plant’s leaf)” takes on value “normal”, then obviously the “leaf discoloration” descriptor is irrelevant. Such information can be represented by an implication:

$$[\text{state(plant's leaf)} = \text{normal}] \Rightarrow [\text{discoloration(plant's leaf)} = \text{NA}]$$

where NA is a special value meaning “not applicable”.

- **Properties of descriptors**—Descriptors that are relations between objects may have certain general properties—they can be reflexive, symmetric, transitive, and so on. All such properties are defined as assertions in the annotated predicate calculus (see the Appendix). For example, the transitivity of relation “above(P_1, P_2)” can be defined as:

$$\forall P_1, P_2, P_3, (\text{above}(P_1, P_2) \ \& \ \text{above}(P_2, P_3)) \Rightarrow \text{above}(P_1, P_3)$$

- **Interrelationships among descriptors**—In some problems there may exist relationships between descriptors that constrain their values. For example, the length of an object is assumed always to be greater than or equal to its width:

$$\forall P, \text{length}(P) \geq \text{width}(P)$$

Also, descriptors may be related by known equations. For example, the area of a rectangle is the arithmetic product of its length and width:

$$\forall P, ([\text{shape}(P) = \text{rectangle}] \Rightarrow [\text{area}(P) = \text{length}(P) \times \text{width}(P)])$$

The infix operator “ \times ” is used to simplify notation of the term $\text{multiply}(\text{length}(P), \text{width}(P))$.

4.4.6 The Form of Observational and Inductive Assertions

The basic form of assertions in the Star methodology is a *c-expression*, defined as a conjunctive statement:

$$\langle \text{quantifier form} \rangle \langle \text{conjunction of relational statements} \rangle \quad (8)$$

where $\langle \text{quantifier form} \rangle$ stands for zero or more quantifiers, and $\langle \text{relational}$

statements> are predicates in a special form, as defined in the Appendix. The following is an example of a c-expression:

$\exists P_0, P_1, P_2, P_3 ([\text{contains}(P_0, P_1, P_2, P_3)][\text{ontop}(P_1 \& P_2, P_3)][\text{length}(P_1) = 3..5]$
 $[\text{weight}(P_1) > \text{weight}(P_2)][\text{color}(P_1) = \text{red} \vee \text{blue}][\text{shape}(P_1 \& P_2 \& P_3) = \text{box}]$

that can be paraphrased in English:

An object P_0 contains parts P_1 , P_2 and P_3 and only these parts. Parts P_1 & P_2 are on top of part P_3 , length of P_1 is between 3 and 5, the weight of P_1 is greater than that of P_2 , the color of P_1 is red or blue, and the shape of all three parts is box.

An important special case of a c-expression is an *a-expression* (an atomic expression), in which there is no "internal disjunction" (see the Appendix).

Note that due to the use of internal disjunction a c-expression represents a more general concept than a universally quantified conjunction of predicates, used in typical production rules.

Progressively more complex forms of expressions are described below:

- A *case expression* is a logical product of implications:

$$[L = a_i] \Rightarrow \text{Exp}_i, i = 1, 2, \dots$$

where a_i are single elements or disjoint subsets of elements from the domain of descriptor L , and Exp_i are c-expressions.

A case expression describes a class of objects by splitting it into separate cases, each represented by a different value(s) of a certain descriptor.

- An *implicative expression* (i-expression):

$$C \& (C_1 \Rightarrow C_2) \tag{9}$$

where C , C_1 and C_2 are c-expressions.

This form of description is very useful when the occurrence of some properties (defined in C_2) depends on the occurrence of some other properties (defined in C_1). Typical production rules used in expert systems are a special case of (9), where C is omitted and no internal logical operators are used. When $(C_1 \Rightarrow C_2)$ is omitted, then the conditional expression becomes a c-expression.

- A *disjunctive expression* (d-expression), defined as a disjunction of implicative expressions.
- An *exception-based expression* (e-expression). In some situations it is simpler to formulate a somewhat overgeneralized statement and indicate exceptions than to formulate a precise statement. The following form is used for such purposes:

$$D_1 \setminus D_2$$

where D_1 and D_2 are d-expressions. This expression is equivalent to $(\sim D_2 \Rightarrow D_1) \& (D_2 \Rightarrow \sim D_1)$.

Observational assertions are formulated as a set of rules:

$$\{\text{a-expression} \Rightarrow K_i\} \quad (10)$$

Inductive assertions are expressed as a set of rules:

$$\{\text{EXP} \Rightarrow \text{c-expression}\} \quad (11)$$

where EXP is a c-expression or any of the more complex expressions described above. It is also assumed that the left side and the right side of (11) satisfy the principle of comprehensibility described in Section 4.2.

4.4.7 The Preference Criterion

In spite of the constraints imposed by the above components of the background knowledge, the number of inductive assertions consistent with observational statements may still be unlimited. The problem then arises of choosing the most desirable inductive assertion(s). In making such a choice, one must take into consideration the aspects of the particular inductive learning problem; therefore the definition of a "preference criterion" for selecting a hypothesis is a part of the problem background knowledge. Typically, the inductive assertions are chosen on the basis of some simplicity criterion (such as given in [Kemeni, 1953; Post, 1960]).

In the context of scientific discovery, philosopher Karl Popper [1968] has advocated constructing hypotheses that are both simple and easy to refute. By generating such hypotheses and conducting experiments aimed at refuting them, he argues, one has the best chance of ultimately formulating the true hypothesis. In order to use this criterion for automated inductive inference, it is necessary to define it formally. This, however, is not easy because there does not seem to exist any universal measure of hypothesis simplicity and refutability.

Among more specific measures for evaluating the "quality" of inductive assertions one may list:

- An overall simplicity for human comprehension, measured, for example, by the number of descriptors and number of operators used in an inductive assertion.
- The degree of "fit" between the inductive and observational assertions (measured, for example, by the degree of generalization, defined as the amount of uncertainty that any given description satisfying the inductive assertion corresponds to some observational statement [Michalski, 980c]).
- The cost of measuring values of descriptors used in the inductive assertion.
- The computational cost of evaluating the inductive assertion.
- The memory required for storing the inductive assertion.
- The amount of information needed for encoding the assertion using predefined operators [Coulon & Kayser, 1978].

The importance given to each such measure depends on the ultimate pur-

pose of constructing the inductive assertions. For that reason, the Star methodology allows a user to build a global preference criterion as a function of such measures, tailored to a specific inductive problem. Since some of the above measures are computationally costly, simpler measures are used, called *elementary criteria*. Among such criteria are: the number of c-expressions in the assertion, the total number of relational statements, the ratio of possible but unseen events implied by an assertion to the total number of training events (a simple measure of generalization), and the total number of different descriptors. The global preference criterion is formulated by selecting from the above list those elementary criteria that are most relevant to the problem, and then arranging them into a *lexicographic evaluation functional* (LEF). A LEF is defined as a sequence of criterion-tolerance pairs:

$$\text{LEF: } (c_1, \tau_1), (c_2, \tau_2)\dots \quad (12)$$

where c_i is an elementary criterion selected from the available "menu", and τ_i is a *tolerance threshold* for criterion c_i ($\tau_i \in [0..100\%]$).

Given a set of inductive assertions, the LEF determines the most preferable one(s) in the following way:

In the first step, all assertions are evaluated from the viewpoint of criterion c_1 , and those which score best, or within the range defined by the threshold τ_1 from the best, are retained. Next the retained assertions are evaluated from the viewpoint of criterion c_2 and reduced similarly as above, using tolerance τ_2 . This process continues until either the subset of retained assertions contains only one assertion (the "best" one) or the sequence of criterion-tolerance pairs is exhausted. In the latter case, the retained set contains assertions that are equivalent from the viewpoint of the LEF.

An important and somewhat surprising property of such an approach is that the same learning system can generate either characteristic or discriminant descriptions of object classes by properly defining the preference criterion (see Section 4.7).

4.5 GENERALIZATION RULES

4.5.1 Definitions and an Overview

Constructing an inductive assertion from observational statements can be conceptually characterized as a heuristic state-space search [Nilsson, 1980], where:

- *states* are symbolic descriptions; the initial state is the set of observational statements.
- *operators* are inference rules, specifically, *generalization*, *specialization* and *reformulation* rules, as defined below.
- the *goal* state is an inductive assertion that implies the observational state-

ments, satisfies the problem background knowledge and maximizes the given preference criterion.

A generalization rule is a transformation of a description into a more general description, one that tautologically implies the initial description. A specialization rule makes an opposite transformation: given a description, it generates a logical consequence of it. A reformulation rule transforms a description into another, logically-equivalent description. A reformulation rule can be viewed as a special case of a generalization and a specialization rule.

Specialization and reformulation rules are the conventional truth-preserving inference rules used in deductive logic. In contrast to them, the generalization rules are not truth-preserving but falsity preserving. This means that if an event falsifies some description, then it also falsifies a more general description. This is immediately seen by observing that $H \Rightarrow F$ is equivalent to $\sim F \Rightarrow \sim H$ (the law of contraposition). To illustrate this point, suppose that a statement "some water birds in this lake are swans" has been generalized to "all water birds in this lake are swans." If there are no water birds in the lake that are swans, then this fact falsifies not only the first statement but also the second. Falsifying the second statement, however, does not imply the falsification of the first.

In concept acquisition, as explained in Section 4.2.2, transforming a rule $E \Rightarrow K$ into a more general rule $D \Rightarrow K$ means that description E must imply description D :

$$E \Rightarrow D \quad (13)$$

(recall expression (6)). Thus, to obtain a generalization rule for concept acquisition, one may use a tautological implication of formal logic. The premise and consequence of such an implication must, however, be interpretable as a description of a class of objects. For example, the known law of simplification:

$$P \& Q \Rightarrow P \quad (14)$$

can be turned into a generalization rule:

$$P \& Q \Rightarrow K \quad | \quad P \Rightarrow K \quad (15)$$

If P stands for "round objects", Q for "brown objects" and K for "balls", then rule (15) states that the expression "round and brown objects are balls" can be generalized to "round objects are balls." Thus, in concept acquisition, the generalization operation has a simple set-theoretical interpretation: a description is more general if it is satisfied by a larger number of objects. (Such an interpretation does not apply, however, to descriptive generalization, as shown below.)

In order to obtain a rule for descriptive generalization, implication (14) is reversed, and P and Q are interpreted as properties of objects of some class K :

$$P(K) \quad | \quad P(K) \& Q(K) \quad (16)$$

If $P(K)$ stands for "balls are round" and $Q(K)$ for "balls are brown," then according to rule (16), the statement "balls are round and brown" is a generaliza-

tion of the statement "balls are round" (because from the former one can deduce the latter). We can see that the notion "the number of objects satisfying a description" is not applicable here. Generalizing means here adding (hypothesizing) properties that are ascribed to a class of objects.

After this informal introduction we shall now present various types of generalization rules, concentrating primarily on the rules for concept acquisition. These rules will be expressed using the notation of the annotated predicate calculus (see the Appendix). The reverse of these rules are specialization rules and, as special cases, reformulation rules. With regard to other specialization and reformulation rules we shall refer the reader to a standard book on predicate calculus (such as [Suppes, 1957]). Some reformulation rules of the annotated predicate calculus that do not occur in ordinary predicate calculus are given in the Appendix.

We will restrict our attention to generalization rules that transform one or more statements into a single more general statement:

$$\{D_i \Rightarrow K\}_{i \in I} \mid \leftarrow D \Rightarrow K \quad (17)$$

Such a rule states that if an event (a symbolic description of an object or situation) satisfies any description D_i , $i \in I$, then it also satisfies description D (the reverse may not be true). A basic property of the generalization transformation is that the resulting description has "unknown" truth-status, that is, is a hypothesis that must be tested on new data. A generalization rule does not guarantee that the obtained description is useful or plausible.

We distinguish between two types of generalization rules, *selective* and *constructive*. If every descriptor used in the generated concept description D is among descriptors occurring in the initial concept descriptions D_i , $i = 1, 2, \dots$, then the rule is selective, otherwise it is constructive.

4.5.2 Selective Generalization Rules

In the rules presented below, CTX , CTX_1 and CTX_2 stand for some arbitrary expressions (context descriptions) that are augmented by additional components to formulate a concept description.

- The *dropping condition* rule—This rule is a generalized version of the previously described rule (15):

$$CTX \& S \Rightarrow K \mid \leftarrow CTX \Rightarrow K \quad (18)$$

where S is an arbitrary predicate or logical expression.

This rule states that a concept description can be generalized by simply removing a conjunctively-linked expression. This is one of the most commonly-used rules for generalizing information.

- The *adding alternative* rule:

$$CTX_1 \Rightarrow K \mid \leftarrow CTX_1 \vee CTX_2 \Rightarrow K \quad (19)$$

A concept description can be generalized by adding, through the use of

logical disjunction, an alternative to it. An especially useful form of this rule is when the alternative is added by extending the scope of permissible values of one specific descriptor. Such an operation can be expressed very simply by using the internal disjunction operator of the annotated predicate calculus. For example, suppose that a concept description is generalized by allowing objects to be not only red but also blue. This can be expressed as follows:

$$\text{CTX} \& [\text{color}=\text{red}] \Rightarrow K \quad | \quad \text{CTX} \& [\text{color}=\text{red} \vee \text{blue}] \Rightarrow K \quad (20)$$

(Forms in brackets are selectors; the expressions on the right of '=' are called references—see the Appendix)

Because of the importance of this special case, it will be presented as a separate general rule.

- The *extending reference* rule:

$$\text{CTX} \& [L=R_1] \Rightarrow K \quad | \quad \text{CTX} \& [L=R_2] \Rightarrow K \quad (21)$$

where $R_1 \subseteq R_2 \subseteq \text{DOM}(L)$ and $\text{DOM}(L)$ denotes the domain of L .

In this rule, L is a term, and R_1 and R_2 (references) are internal disjunctions of values of L . References R_1 and R_2 can be interpreted as sets of values that descriptor L can take in order to satisfy the concept description.

The rule states that a concept description can be generalized by enlarging the reference of a descriptor ($R_2 \supseteq R_1$). The elements added to R_2 must, however, be from the domain of L .

If R_2 is extended to be the whole domain, that is, $R_2 = \text{DOM}(L)$, then the selector $[L = \text{DOM}(L)]$ is always true, and therefore can be removed. In this case, the extending reference rule becomes the dropping condition rule. There are two other special cases of the extending reference rule. They take into consideration the type of the descriptor L [defined by the structure of $\text{DOM}(L)$]. They are presented as separate rules below.

- The *closing interval* rule:

$$\begin{array}{l} \text{CTX} \& [L=a] \Rightarrow K \\ + \\ \text{CTX} \& [L=b] \Rightarrow K \end{array} \quad \left| \quad \text{CTX} \& [L=a..b] \Rightarrow K \quad (22)$$

where L is a linear descriptor, and a and b are some specific values of descriptor L . The two premises are assumed to be connected by the logical conjunction (this convention holds for the remaining rules as well).

The rule states that if two descriptions of the same class (the premises of the rule) differ in the values of only one linear descriptor, then the descriptions can be replaced by a single description in which the reference of the descriptor is the interval linking these two values.

To illustrate this rule, consider as objects two states of a machine, and K as a class of *normal* states. The rule says that if a machine is in the normal state for two different temperatures, say a and b , then a hypothesis is made that all states in which the temperature falls into the interval $[a,b]$ are also normal.

Thus, this rule is not only a logically-valid generalization rule, but expresses also some aspect of plausibility.

- The *climbing generalization tree* rule

$$\begin{array}{l}
 \text{(one or} \\
 \text{more} \\
 \text{statements)}
 \end{array}
 \left.
 \begin{array}{l}
 \text{CTX \& [L=a] } ::> \text{ K} \\
 \text{CTX \& [L=b] } ::> \text{ K} \\
 \cdot \\
 \cdot \\
 \cdot \\
 \text{CTX \& [L=i] } ::> \text{ K}
 \end{array}
 \right|
 < \text{CTX \& [L=s] } ::> \text{ K} \quad (23)$$

where L is a structured descriptor, and s represents the lowest parent node whose descendants include nodes a, b, ... and i, in the generalization tree domain of L.

The rule is applicable only to descriptions involving structured descriptors, and is used in various forms in, for example [Winston, 1977; Hedrick, 1974; Lenat, 1976] (see also Chapters 11 and 6 of this book). The following example illustrates the rule:

$$\begin{array}{l}
 \exists P, \text{CTX \& [shape(P)=triangle] } ::> \text{ K} \\
 \exists P, \text{CTX \& [shape(P)=rectangle] } ::> \text{ K}
 \end{array}
 \left|
 < \exists P, \text{CTX \& [shape(P)=polygon] } ::> \text{ K}$$

Paraphrasing this rule in English: if an object of class K is triangular and another object of this class is rectangular, then the rule generates a statement that objects of class k are polygonal.

- The *turning constraints into variables* rule—This rule is best known for the case of descriptive generalization:

$$\begin{array}{l}
 \text{(one or} \\
 \text{more} \\
 \text{statements)}
 \end{array}
 \left.
 \begin{array}{l}
 \text{F[a]} \\
 \text{F[b]} \\
 \cdot \\
 \cdot \\
 \cdot \\
 \text{F[i]}
 \end{array}
 \right|
 < \forall v, \text{F[v]} \quad (24)$$

where F[v] stands for some description (formula) dependent on variable v, and a, b, ... are constants.

If some description F[v] holds for v being a constant a or constant b, and so on, then the rule generalizes these observations into a statement that F[v] holds for every value of v. This is the rule used most often in methods of inductive inference employing predicate calculus.

A corresponding rule for concept acquisition is:

$$\text{F[a] \& F[b] \& ... } ::> \text{ K} \quad | < \exists v, \text{F[v] } ::> \text{ K} \quad (25)$$

To illustrate this version, assume that a, b, and so on, are parts of an object of class K that have a property F. Rule (25) generalizes these facts into an

assertion that if any part of an object has property F then the object belongs to class K.

- The *turning conjunction into disjunction* rule:

$$F_1 \& F_2 \Rightarrow K \quad | < \quad F_1 \vee F_2 \Rightarrow K \quad (26)$$

where F_1 and F_2 are arbitrary descriptions.

A concept description can be generalized by replacing the conjunction operator by the disjunction operator.

- The *extending the quantification domain* rule—In the simplest case, the rule changes the universal quantifier into the existential quantifier:

$$\forall v, F[x] \Rightarrow K \quad | < \quad \exists v, F[v] \Rightarrow K \quad (27)$$

This rule can be viewed as a generalization of the previous rule (26). Using the concept of numerical quantifier (see the Appendix) this rule can be expressed in an even more general way:

$$\exists(I_1)v, F[v] \Rightarrow K \quad | < \quad \exists(I_2)v, F[v] \Rightarrow K \quad (28)$$

where I_1, I_2 are the quantification domains (sets of integers) satisfying relation $I_1 \subseteq I_2$.

For example, the statement “if an object has two parts ($I_1 = \{2\}$) with property F, then it belongs to class K” can be generalized by rule (28) to a statement “if an object has two or more parts ($I_2 = \{2, 3, \dots\}$) with property F then it belongs to class K.”

- The *inductive resolution* rule

- As applied to concept acquisition

The deductive inference rule, called the resolution principle, widely used in automatic theorem proving, can be adopted as a rule of generalization for concept acquisition. In propositional form, the resolution principle can be expressed as:

$$(P \Rightarrow F_1) \& (\sim P \Rightarrow F_2) \quad | > \quad F_1 \vee F_2 \quad (29)$$

where P is a predicate and F_1 and F_2 are arbitrary formulas. By interpreting both sides of (29) as concept descriptions, and making appropriate transformations we obtain:

$$\begin{array}{l} P \& F_1 \Rightarrow K \\ \sim P \& F_2 \Rightarrow K \end{array} \quad \left| \quad < \quad F_1 \vee F_2 \Rightarrow K \quad (30)$$

To illustrate this rule, assume that K is the set of situations when John goes to a movie. Suppose that it has been observed that he goes to a movie when he has company (P) and the movie has high rating (F_1), or when he does not have company ($\sim P$), but has plenty of time (F_2). Rule (30) generalizes these two observations to a statement “John goes to a movie when either the movie has high rating or he has plenty of time.”

(ii) As applied to descriptive generalization

By applying logical equivalence $(Q \supset P) \Leftrightarrow (\sim P \supset \sim Q)$ (the law of contraposition) to expression (29), then reversing the obtained rule and substituting the negative literals by the positive, we obtain:

$$P \& F_1 \vee \sim P \& F_2 \supset F_1 \& F_2 \quad (31)$$

This version has been formulated by Morgan (1975).

Both versions, (i) and (ii), can be generalized by applying the full-fledged resolution principle that uses predicates with arguments, and the unification algorithm to unify these arguments (for example, [Chang & Lee, 1973]).

• The *extension against* rule:

$$\begin{array}{l} \text{CTX}_1 \& [L=R_1] \quad \Rightarrow \quad K \\ \text{CTX}_2 \& [L=R_2] \quad \Rightarrow \quad \sim K \end{array} \quad \left| \quad \begin{array}{l} < \quad [L \neq R_2] \quad \Rightarrow \quad K \end{array} \quad (32)$$

where sets R_1 and R_2 are assumed to be disjoint.

Given a description of an object belonging to class K (a positive example), and a description of an object not belonging to this class (a negative example), the rule produces the most general statement consistent with these two descriptions. It is an assertion that classifies an object as belonging to class K if descriptor L does not take any value from the set R_2 , thus ignoring context descriptions CTX_1 and CTX_2 . This rule is the basic rule for learning discriminant descriptions from examples used in the previously-mentioned inductive program AQ11 [Michalski & Larson, 1978]. Various modifications of this rule can be obtained by replacing reference R_2 in the output assertion by some superset of it that does not intersect with R_1 .

4.5.3 Constructive Generalization Rules

Constructive generalization rules generate inductive assertions that use descriptors not present in the original observational statements. This means that the rules perform a transformation of the original representation space. The following is a general constructive rule that makes such a transformation by applying the knowledge of a relationship between different concepts. It is assumed that this relationship is known to the learning system as background knowledge, as a previously-learned concept, or that it is computed according to user-defined procedures.

$$\begin{array}{l} \text{CTX} \& F_1 \quad \Rightarrow \quad K \\ F_1 \Rightarrow F_2 \end{array} \quad \left| \quad \begin{array}{l} < \quad \text{CTX} \& F_2 \quad \Rightarrow \quad K \end{array} \quad (33)$$

The rule states that if a concept description contains a part F_1 (a concept, a subdescription, and so on) that is known to imply some other concept F_2 , then a more general description is obtained by replacing F_1 by F_2 . For example, sup-

pose a learning system is told that if an object is black, wide and long, then it belongs to class K (for example, it is a blackboard). This can be expressed in the annotated predicate calculus:

$$\exists P, [\text{color}(P) = \text{black}][\text{width}(P) \ \& \ \text{length}(P) = \text{large}] \ ::> \ K$$

Suppose the learner already knows that:

$$\forall P, ([\text{width}(P) \ \& \ \text{length}(P) = \text{large}] \Rightarrow [\text{area}(P) = \text{large}])$$

Then rule (33) produces a generalization:

$$\exists P, [\text{color}(P) = \text{black}][\text{area}(P) = \text{large}] \ ::> \ K$$

As another example, suppose the system is given a description of an object classified as an arch. This description states that a horizontal bar is on top of two equal objects placed apart, B_1 and B_2 , having certain color, weight, shape, and so on. Suppose now that characterizations of B_1 and B_2 in this description satisfy a previously-learned concept of a block. Then rule (33) generates an assertion that an arch is a bar on top of two blocks placed apart. This rule is the basis for an interactive concept learning system developed by Sammut [1981].

Specific constructive generalization rules can be obtained from (33) by evoking procedures computing new descriptors in expression F_2 as functions of initial or previously-derived descriptors (contained in F_1). Here are some examples of rules for generating new descriptors.

- *Counting arguments rules*

- (i) The CQ rule (count quantified variables)—If a concept description is in the form:

$$\exists v_1, v_2, \dots, v_k, F[v_1, v_2, \dots, v_k]$$

then the rule generates descriptors “#v-COND” representing the number of v_i 's that satisfy some condition COND. This condition expresses selected properties of v_i 's specified in the concept description. Since many such COND's can usually be formulated, the rule allows the system to generate a large number of such descriptors.

For example, if the COND is “[attribute₁(v_i) = R]”, then the generated descriptor will be “# v_i -attribute₁-R” counting the number of v_i 's that satisfy this condition. If the attribute₁ is, for instance, length, and R is [2..4], then the derived descriptor is “# v_i -length-2..4” (that is, it measures the number of v_i 's whose length is between 2 and 4, inclusively).

- (ii) The CA-rule (count arguments of a predicate)—If a descriptor in a description is a relation with several arguments, REL(v_1, v_2, \dots), the rule generates descriptors “#v-COND”, measuring the number of arguments in REL that satisfy some condition COND. As above, many such descriptors can be generated, each with different COND.

The annotation of a descriptor provides information about its properties. Such a property may be that a descriptor is, for instance, a

transitive relation, such as relation "above", "inside", "left-of", and "before". For example, if the relation is "contains(A,B₁,B₂,...)", stating that object A contains objects B₁,B₂,..., and COND is "large and red", then the derived descriptor "#B-large-red-A-contains" measures the number of B_i-s contained in A that are large and red.

- The *generating chain properties* rule—If the arguments of different occurrences of a transitive relation in a concept description form a chain (that is, form a sequence of consecutive objects ordered by this relation), the rule generates descriptors characterizing some specific objects in the chain. Such objects may be:

| | |
|-------------------------|---|
| LST-object | the "least object", or the object at the beginning of the chain (for example, the bottom object in the case of the relation "above"). |
| MST-object | the object at the end of the chain (for example, the top object). |
| MID-object | the objects in the middle of the chain. |
| N th -object | the object in the N th position in the chain (starting from LST-object). |

After identifying these objects, the rule investigates all known properties of them (as specified in the observational statements) in order to determine potentially relevant new descriptors. The rule also generates a descriptor characterizing the chain itself, namely:

REL-chain-length: the length of the chain defined by relation REL.

For example, if the REL is ON-TOP, then descriptor ON-TOP-chain-length would specify the height of a stack of objects. When a new description is generated and adopted, an annotation for it is also generated and filled out, as in Lenat [1976]. This rule can be extended to a partial order relation. In such a case it becomes the "find extrema of a partial order" rule.

- The *detecting descriptor interdependence* rule—Suppose that given is a set of objects exemplifying some concept, and that attribute descriptions are used to characterize these objects. Such descriptions specify only attribute values of the objects; they do not characterize the objects' structure. Suppose that the values a linear descriptor x takes on in all descriptions (events) are ordered in increasing order. If the corresponding values of another linear descriptor y exhibit an increasing or decreasing order, then a two-place descriptor:

M(x,y)

is created, signifying that x and y have a monotonic relationship. This descriptor has value ↑ when y values are increasing and value ↓ when they are decreasing.

The idea of the above M-descriptor can be extended in two directions. The first is to create M-descriptors dependent on some condition COND that must be satisfied by the events under consideration:

M(x,y)-COND

For example, descriptor:

M(length,weight)-red

states that length and weight have a monotonic relationship for red objects.

The second direction of extension is to relax the requirement for the monotonic relationship; that is, not to require that the order of y values is strictly increasing (or decreasing), but only approximately increasing (or decreasing). For example, the coefficient of statistical correlation between x and y can be measured, and when its absolute value is above a certain threshold, a descriptor $R(x,y)$ is created. The domain of this R -descriptor can also be $\{\uparrow, \downarrow\}$, indicating the positive or negative correlation, respectively, or it can have values representing several subranges of the correlation coefficient. Similarly, as in the case of M -descriptors, R -descriptors can be extended to R -COND descriptors.

The M - or R -descriptors can be used to generate new descriptors. For example, if $[M(x,y) = \uparrow]$, then a new descriptor $z = x/y$ can be generated. If z assumes a constant or nearly-constant value, then an important relationship has been discovered. Similarly, if $[M(x,y) = \downarrow]$ then a new descriptor $z = x \times y$ can be generated. These two techniques for generating new descriptors have been successfully used in the BACON system for discovering mathematical expressions representing physical or chemical laws, as described in Chapter 10 of this book.

The above ideas can be extended to structural descriptions. Such descriptions involve not only global properties of objects, but also properties of objects' parts and the relationships among the parts. Suppose that in a structural description of an object, existentially-quantified variables P_1, P_2, \dots, P_m denote its parts. If $x(P_i)$ and $y(P_i)$ are linear descriptors of P_i (for example, numerical attributes characterizing parts P_i , $i = 1, 2, \dots$), the above-described techniques for generating M - and R -descriptors can be applied.

4.6 THE STAR METHODOLOGY

4.6.1 The Concept of a Star

The methodology presented here for learning structural descriptions from examples receives its name from the major concept employed in it, that of a *star*. In the most general sense, a *star of an event e under constraints E* is a set of all possible alternative non-redundant descriptions of event e that do not violate constraints E . A somewhat more restrictive definition of a star will be used here. Let e be an example of a concept to be learned and E be a set of some counterexamples of this concept. A star of the event e *against* the event set E , denoted $G(e|E)$, is defined as the set of all maximally general c -expressions that cover (that is, are satisfied by) event e and that do not cover any of the negative events in E .

The *c*-expressions in a star may contain *derived* descriptors, that is, descriptors not present in the observational statements. In such a case, testing whether event *e* satisfies a given description requires that appropriate transformations be applied to the event. Such a process can be viewed as proving that the event implies the description, and therefore methods of automatic theorem proving could be used.

In practical problems, a star of an event may contain a very large number of descriptions. Consequently, such a theoretical star is replaced by a *bounded star* $G(e|E,m)$ that contains no more than a fixed number, *m*, of descriptions. These *m* descriptions are selected as the *m* most preferable descriptions, among the remaining ones, according to the preference criterion defined in the problem background knowledge. Variable *m* is a parameter of the learning program, defined either by the user or by the program itself, as a function of the available computational resources.

Chapter 11 of this book gives an illustration and an algorithm for generating a bounded star with *c*-expressions restricted to attribute expressions (that is, expressions involving only object attributes). Section 4.6.3 presents an algorithm for generating a bounded star consisting of regular *c*-expressions. The concept of a star is useful because it reduces the problem of finding a complete description of a concept to subproblems of finding consistent descriptions of single positive examples of the concept.

Since any single example of a concept can always be characterized by a conjunctive expression (a logical product of some predicates), elements of a star can always be represented by conjunctive descriptions. One should also notice that if the concept to be learned is describable by a *c*-expression, then this description clearly will be among the elements of a (non-bounded) star of any single positive example of the concept. Consequently, if there exists a positive example not covered by any description of such a star, then the complete concept description must be disjunctive, that is, must include more than one *c*-expression.

4.6.2 Outline of the General Algorithm

It is assumed that every observational statement is in the form:

$$\text{a-expression} ::= > K \quad (34)$$

where a-expression is an atomic expression describing an object (recall Section 4.4.6) and *K* is the concept exemplified by this object.

It is also assumed that inductive assertions are in the form of a single *c*-expression or the disjunction of *c*-expressions. For simplicity we will restrict our attention to only single-concept learning. In the case of multiple-concept learning, the algorithm is repeated for each concept with modifications depending on the assumed interdependence among the concept descriptions (Section 4.2.4).

Let POS and NEG denote sets of events representing positive and negative examples of a concept, respectively. A general and simplified version of the Star methodology can be described as follows:

1. Randomly select an event e from POS.
2. Generate a bounded star, $G(e|NEG,m)$, of the event e against the set of negative examples NEG, with no more than m elements. In the process of star generation apply generalization rules (both selective and constructive), task-specific rules, heuristics for generating new descriptors supplied by problem background knowledge, and definitions of previously-learned concepts.
3. In the obtained star, find a description D with the highest preference according to the assumed preference criterion LEF.
4. If description D covers set POS completely, then go to step 6.
5. Otherwise, reduce the set POS to contain only events not covered by D , and repeat the whole process from step 1.
6. The disjunction of all generated descriptions D is a complete and consistent concept description. As a final step, apply various reformulation rules (defined in the problem background knowledge) and "contracting" rules [equations (8) and (9) in the Appendix] in order to obtain a possibly simpler expression.

This algorithm is a simplified version of the general covering algorithm A9 [Michalski, 1975b]. The main difference is that algorithm A9 selects the initial events (if possible) from events not covered by any of the descriptions of generated stars, rather than not covered by only the selected descriptions D . This way the algorithm is able to determine a bound on the maximum number of separate descriptions in a disjunction needed to define the concept. Such a process may, however, be computationally costly.

The above algorithm describes only single-step learning. If, after generating a concept description, a newly-presented training event contradicts it, specialization or generalization rules are applied to generate a new consistent concept description. A method for such incremental learning is described in [Michalski & Larson, 1978]. (See also Chapter 8 of this book.)

The central step in the above methodology is the generation of a bounded star. This can be done using a variety of methods. Thus, the above Star methodology can be viewed as a general schema for implementing various learning methods and strategies. The next section describes one specific method of star generation.

4.6.3 Star Generation: The INDUCE Method

This method generates a bounded star $G(e|NEG,m)$ by starting with a set of expressions that are single selectors, either extracted from the event for which the star is generated or inferred from the event by applying constructive generalization rules or inference rules provided by background knowledge. These expressions are then specialized by adding other selectors until consistency is achieved (that is, until each expression does not intersect with set NEG).

Next, the obtained consistent expressions are generalized so that each achieves the maximum coverage of the remaining positive training examples. The best consistent m so obtained and the generalized c -expressions (if some are also complete, then they are alternative solutions) constitute the bounded star sought, $G(e|NEG, m)$. Specifically, the steps of the procedure are:

1. In the first step individual selectors of event e are put on the list called PS. This list is called a *partial star*, because its elements may cover some events in NEG. These initial elements of PS (single selectors from e) can be viewed as generalizations of event e obtained by applying in all possible ways the dropping condition generalization rule (each application drops all selectors except one). Elements of the partial star PS are then ordered from the most to the least preferred according to a preference criterion:

$$LEF_1 = \langle (-negcov, \tau_1), (poscov, \tau_2) \rangle \quad (35)$$

where $negcov$ and $poscov$ are numbers of negative and positive examples, respectively, covered by an expression in the star, and τ_1 and τ_2 are tolerances (recall Section 4.4.7).

The LEF_1 minimizes the $negcov$ (by maximizing the $-negcov$) and maximizes $poscov$.

2. The list PS is then expanded by adding new selectors obtained by applying the following inference rules to the event e :
 - a. the constructive generalization rules (Section 4.5.3)
 - b. the problem-specific heuristics defined in the background knowledge
 - c. the definitions of the previously-learned concepts (to determine whether parts of e satisfy some already known concepts)
3. Each new selector is inserted in the appropriate place in list PS, according to preference criterion LEF_1 . The size of PS is kept within the limit defined by parameter m by removing from PS all but the m most preferred selectors.
4. Descriptions in PS are tested for consistency and completeness. A description is consistent if $negcov = 0$ (that is, if it covers no events in NEG) and is complete if $poscov$ is equal to the total number of positive examples. Consistent and complete descriptions are removed from PS and put on the list called SOLUTIONS. If the size of the list SOLUTIONS is greater than a parameter #SOL, then the algorithm stops. Parameter #SOL determines the number of desired alternative concept descriptions. Incomplete but consistent descriptions are removed from the list PS and put on the list called CONSISTENT. If the size of the CONSISTENT list is greater than a parameter #CONS, then control is transferred to step 6.
5. Each expression in PS is specialized in various ways by appending to it a single selector from the original list PS. Appended selectors must be of lower preference than the last selector in the conjunctive expression

(initially, the expression has only one selector). Parameter %BRANCH specifies the percentage of the selectors ranked lower (by the preference criterion) than the last selector in the current conjunction. If %BRANCH = 100%, all lower preference selectors are singly appended—that is, the number of new expressions generated from this conjunction will be equal to the total number of selectors having lower preference than the last selector in the conjunction. All newly-obtained expressions are ranked by LEF_1 and only the m best are retained. This “expression growing” process is illustrated in Figure 4-1.

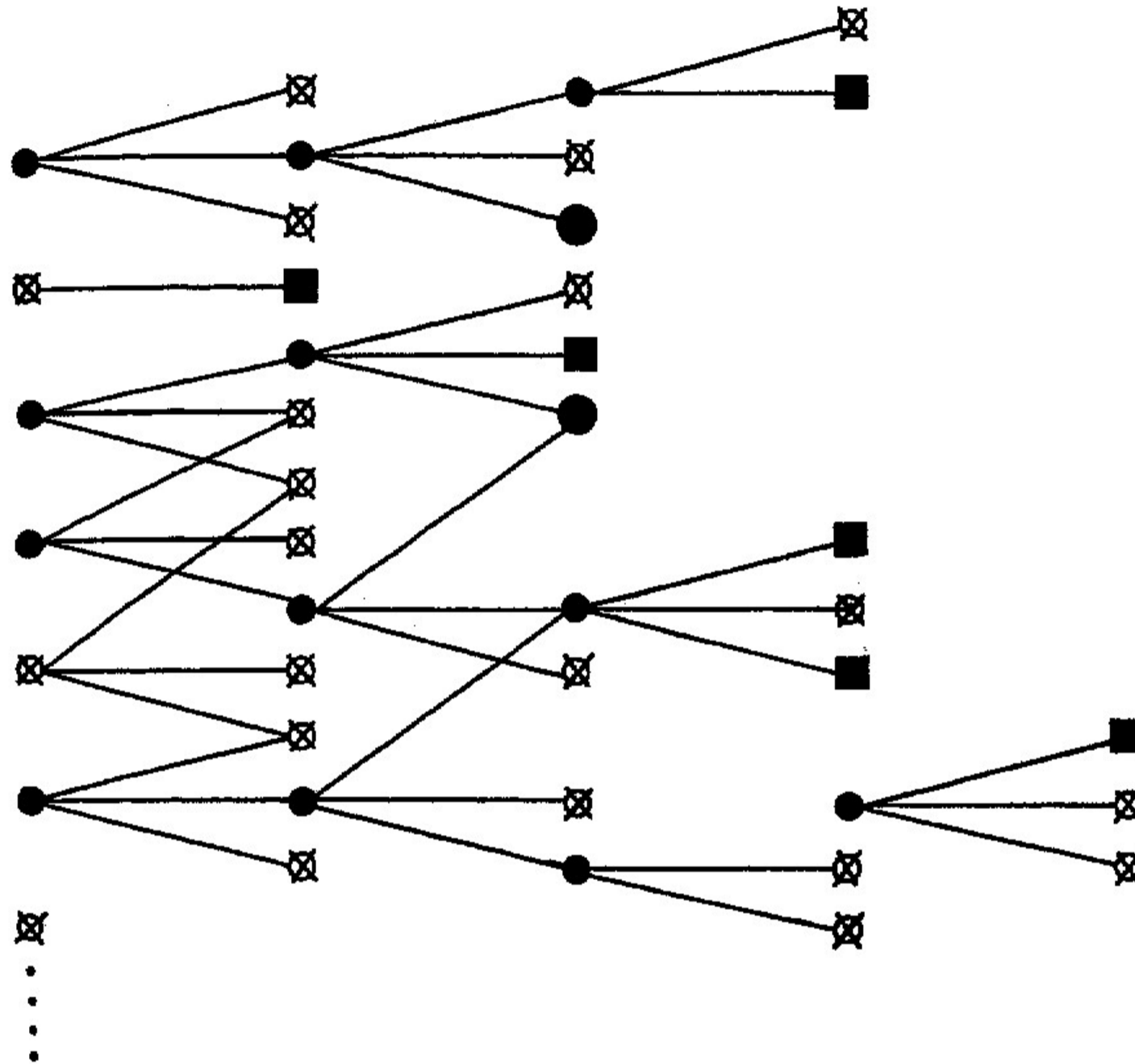
Steps 4 and 5 are repeated until the CONSISTENT list contains the number of expressions specified by parameter #CONS, or until the time allocated for this process is exhausted.

6. Each expression on the CONSISTENT list is generalized by applying the *extension against, closing the interval, and climbing generalization tree* generalization rules. An efficient way to implement such a process is to transform the original structural-description space into an attribute-description space. Attributes (that is, descriptions with zero arguments) defining this space are created from the descriptors in the given expression on the CONSISTENT list in a manner such as that described in Section 3.2.3.2 of Chapter 3 in this book. The generalization of the obtained attribute descriptions is accomplished by the star generation procedure, analogous to the one described in Chapter 11 of this book. Details of this process of transforming structural descriptions into attribute descriptions are described by Larson [1977]. The reason for such a transformation is that structural descriptions are represented as labeled graphs while attribute descriptions are represented as binary strings. It is computationally much more economical to handle binary strings than labeled graphs.
7. The obtained generalizations are ranked according to the global preference criterion LEF defined in the background knowledge. To obtain a discriminant description, a typical LEF is to maximize the number of events covered in POS set and to minimize the complexity of the expression (measured, for example, by the number of selectors it contains). The m best expressions so determined constitute the bounded star $G(e|NEG, m)$.

The Star algorithm and a somewhat restricted version of the above-described star generation algorithm has been implemented in various incarnations of the INDUCE learning program [Larson, 1977; Dietterich, 1978; Michalski, 1980a; Hoff *et al.*, 1982].

4.7 AN EXAMPLE

To illustrate the inductive learning methodology just presented, let us consider a simple problem in the area of conceptual data analysis. Suppose we are



- ☒ - a disregarded rule
- - an active rule
- - a terminal node denoting a consistent c-expression
- - a terminal node denoting a consistent and complete c-expression (a solution)

The nodes in the first column are selectors extracted from the event e or derived from e by applying inference rules. Each arc represents an operation of adding a new selector to the current c-expression.

Figure 4-1: Illustration of the process of generating a reduced star $RG(e|NEG,m)$.

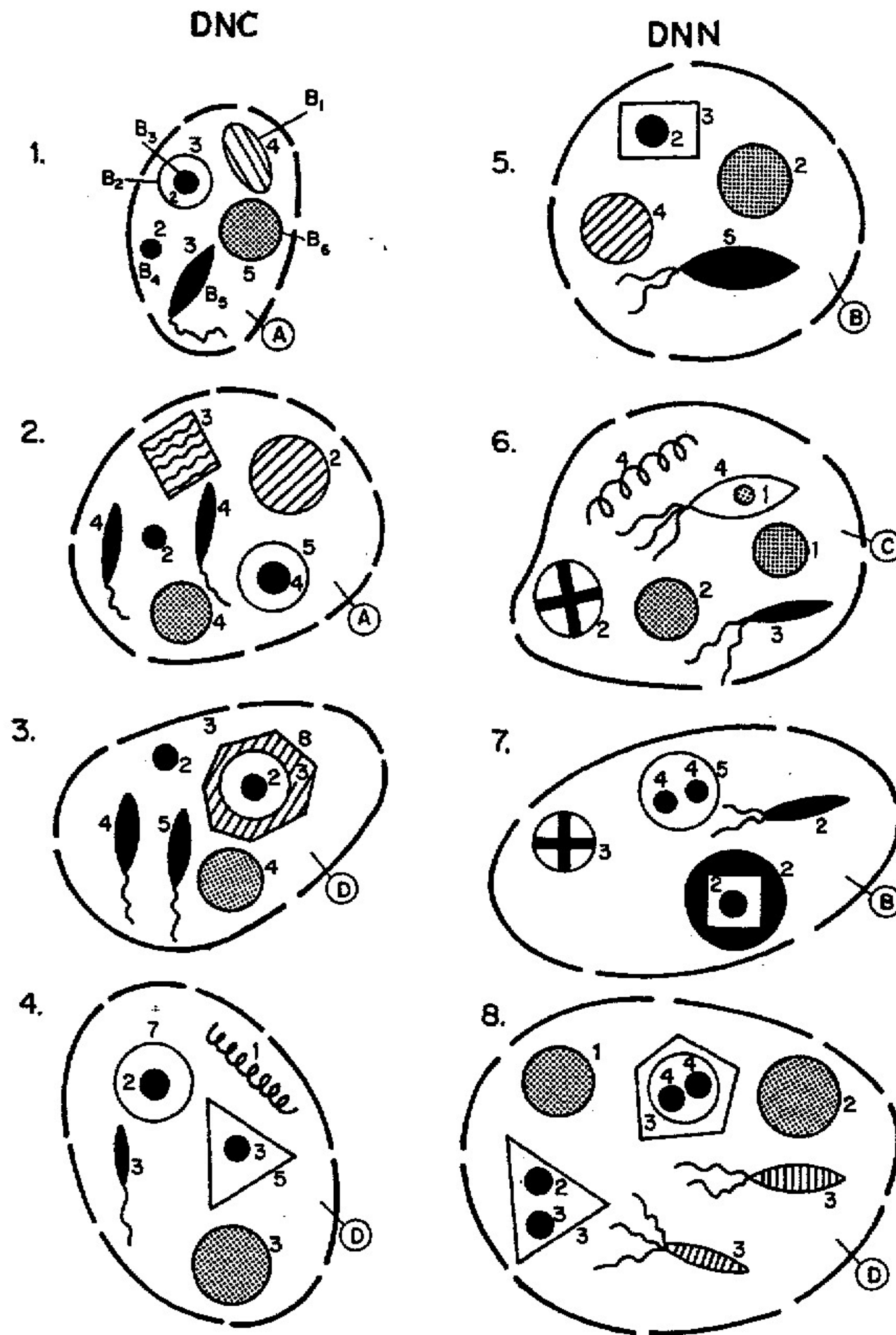


Figure 4-2: "Cancerous" and "Normal" cells.

given examples of "cancerous" and "normal" cells, denoted DNC and DNN, respectively, in Figure 4-2, and the task of the analysis is:

- to determine properties differentiating the two classes of cells (that is, to find discriminant descriptions of each class)
- to determine important common properties of the cancerous and the normal cells (that is, to find characteristic descriptions of each class).

An assumption is made that the properties to be discovered may involve both quantitative information about the cells and their components, and qualitative information, expressed by nominal variables and relationships existing among the components.

The solution to the problem posed (or similar problems) can be obtained by a successive repetition of the "focus attention→hypothesize→test" cycle described below.

The "focus attention" phase is concerned with defining the scope of the problem under consideration. This includes selecting descriptors appearing to be relevant, specifying underlying assumptions, and formulating the relevant problem knowledge. This first phase is performed by a researcher; it involves his/her technical knowledge and informal intuitions. The third, the "test" phase, examines the hypotheses and tests them on new data. This phase may require collecting new samples, performing laboratory experiments, and/or critically analyzing the hypotheses. This phase is likely to involve knowledge and abilities that go beyond currently-feasible computer systems.

It is the second, the "hypothesize" phase, in which an inductive learning system may play a useful role: the role of an assistant for conducting a search for the most plausible and/or most interesting hypotheses. This search may be a formidable combinatorial task for a researcher, if the data sample is large and if each item of the data (in this case, a cell) is described by many variables and/or relations.

Individual steps are as follows:

1. The user determines the set of initial descriptors and provides an annotation for each descriptor. We will assume that the annotation specifies the type, the domain, and any special properties of each descriptor (for example, the transitivity of a relation). In the case of structured descriptors, the annotation also specifies the structure of the domain. The specification of the annotation constitutes the first part of the problem background knowledge.

Suppose that for our simple example problem, the following descriptors are selected:

a. Global descriptors (those characterizing a whole cell)

- circ—the number of segments in the circumference of the cell

Type: linear

Domain: {1..10}

- pplasm—the type of protoplasm in the cell (marked by encircled capital letters in Figure 4-2)

Type: structured

Domain: a tree structure with a set of leaves {triangle, circle, ellipse, heptagon, square, boat, spring}

Non-leaf nodes are defined by rules:

[shape = circle \vee ellipse] \Rightarrow [shape = oval]

[shape = triangle \vee square \vee heptagon] \Rightarrow [shape = polygon]

[shape = oval \vee polygon] \Rightarrow [shape = regular]

[shape = spring \vee boat] \Rightarrow [shape = irregular]

- texture(B_i)—the texture of body B_i

Type: nominal

Domain: {blank, shaded, solid-black, solid-grey, stripes, crossed, wavy}

- weight (B_i)—the weight of body B_i

Type: linear

Domain: {1,2,...,5}

- orient (B_i)—the orientation of B_i

Type: linear-cyclic (the last element is followed by the first)

Domain: {N, NE, E, SE, S, SW, W, NW}

Condition of applicability: if [shape (B_i) = boat]

- contains (C, B_1, B_2, \dots)— C contains B_1, B_2, \dots

Type: nominal

Domain: {True,False}

Properties: transitive relation

- hastails (B, L_1, L_2, \dots)—a body B has tails L_1, L_2, \dots

Type: nominal

Domain: {True,False}

Condition of applicability: if [shape (B) = boat]

Note that the descriptors “contains” and “hastails” are predicates with a variable number of arguments. Descriptor “contains” is characterized as a transitive relation. Descriptors “hastails” and “orient” are applicable only under certain conditions.

2. The user formulates observational statements which describe cells in terms of selected descriptors and specify the class to which each cell belongs. For example, the following is an observational statement for the DNC cell 1:

$$\begin{aligned} \exists \text{CELL}_1, B_1, B_2, \dots, B_6 \text{ [contains(CELL}_1, B_1, \dots, B_6) \text{ [circ(CELL}_1) = 8] \& \\ \text{[ppiasm(CELL}_1) = A] \text{[shape(B}_1) = \text{ellipse] [texture(B}_1) = \text{stripes] \& \\ [weight(B}_1) = 4] \text{[orient(B}_1) = \text{NW] [shape(B}_2) = \text{circle] \& \\ \text{[contains(B}_2, B_3) \text{] [texture(B}_2) = \text{blank] [weight(B}_2) = 3] \dots \& \\ \text{[shape(B}_6) = \text{circle] [texture(B}_6) = \text{shaded] [weight(B}_6) = 5] \\ \text{::> [class = DNC]} \end{aligned}$$

3. To specify the second part of the problem background knowledge the user indicates which general rules of constructive induction (Section 4.5.3) are applicable, and also formulates any problem-specific rules.

The constructive rules will generate various derived descriptors. For example, the counting rule CQ will generate, among others, a descriptor:

- #B-black-boat—the number of bodies whose shape is “boat” and texture is “solid-black” (that is, assuming COND:
[texture(B) = solid-black] & [shape(B) = boat])

(For simplicity of notation, the name of this descriptor, as well as other descriptors below, has been abbreviated, so it does not follow strictly the naming convention described in Section 4.5.3.) The counting rule CA will generate such descriptors as:

- total-B—the total number of bodies in a cell (no COND is used)
- indep-B—the number of independent bodies in a cell (assuming the COND “bodies not contained in another body”)
- #contained-in-B—the number of smaller bodies contained in the body B
- #tails-boat-B—the number of tails in a body B, whose shape is “boat”

As advice to the system, the user may formulate arbitrary arithmetic expressions for generating possibly relevant descriptors. For example, the user may suggest a descriptor:

$$\text{weight(CELL)} = \sum_i \text{weight(B}_i)$$

where B_i , $i = 1, 2, \dots$ denote bodies in a cell.

The background knowledge may also contain special concepts, such as even or odd numbers, the definitions of the area and perimeter of a circle or rectangle, and so on.

4. Finally, as the last part of the background knowledge, the user specifies the type of description sought and the hypothesis preference criterion. Let us assume that both characteristic descriptions and discriminant descriptions are sought. We therefore choose as the preference criterion for constructing characteristic descriptions “maximize the length of the complete c-expressions,” and for constructing discriminant descriptions, “minimize the length of consistent and complete c-expressions.”

As illustration, we shall present here samples of discriminant descriptions and characteristic descriptions of the DNC “cells”, obtained by the INDUCE program.

4.7.1 Discriminant Descriptions of DNC Cells

Each of these descriptions is sufficient to discriminate all DNC cells from DNN cells. A concept description for class DNC can thus be any one of these descriptions or the disjunction of two or more of these descriptions.

- $\exists(1)B$ [texture(B) = shaded][weight(B) > 3]
Paraphrasing in English: "Every DNC cell, as opposed to DNN, has exactly one body with 'shaded' texture and weight at least 3."
- [circ = even]
"The number of segments in the circumference of every DNC cell is even." (The concept of "even" was determined by "climbing the generalization tree" rule.)
- $\exists(> 1)B$ [shape(B) = boat][orient(B) = N \vee NE]
"Every DNC cell has at least one 'boat' shape body with orientation N or NE."
- $\exists(> 1)B$ [#tails-boat-B = 1]
"Every DNC cell has at least one body with number of tails equal to 1."
- $\exists(1)B$ [shape(B) = circle][#contains-B = 1]
"Every DNC cell has a circle containing a single object."

Underscored descriptors are derived descriptors obtained through constructive generalization rules.

4.7.2 Characteristic Descriptions of DNC Cells

Every description below is a characterization of some pattern common to all DNC cells. Some of these patterns taken separately may cover one or more DNN cells. The length of each description has been maximized, rather than minimized, as in the case of discriminant descriptions.

- $\exists(1)B$ [weight(B) = 5]
Paraphrasing in English: "In every DNC cell there is one and only one body with weight 5."
- $\exists.B_1, B_2$ [contains(B₁, B₂)] [shape(B₁) & shape(B₂) = circle] & [texture(B₁) = blank] [weight(B₁) = odd] [texture(B₂) = solid-black] & [weight(B₂) = even] [#contained-in-B₁ = 1]
"In every cell there are two bodies of circle shape, one contained in another, of which the outside circle is blank and has 'odd' weight, the inside circle is solid-black and has 'even' weight. The number of bodies in the outside circle is only one." (This is also a non-minimal discriminant description.)
- $\exists(1)B$ [shape(B) = circle][texture(B) = shaded][weight(B) > 3]
"Every cell contains a circle with 'shaded' texture, whose weight is at least 3." (This is also a non-minimal discriminant description.)

- $\exists (> 1)B$ [shape(B)=boat][orient(B) = N \vee NE][#tails-boat(B)=1]
 “Every cell has at least one body of ‘boat’ shape with N or NE orientation, which has one tail.” (This is also a non-minimal discriminant description.)
- $\exists (2)B$ [shape(B)=circle][texture(B)=solid-black], or, alternatively, [#B-circle-solid-black=2]
 “Each cell has exactly two bodies that are solid black circles.” (This is also a non-minimal discriminant description.)
- [pplasm = A \vee D]
 “The protoplasm of every cell is of type A or D.”

The above example is too simple for really unexpected patterns to be discovered. But it illustrates well the potential of the learning program as a tool for searching for patterns in complex data, especially when the relevant properties involve both numerical and structural information about the objects under consideration. An application of this program to a more complex problem [Michalski, 1980a] did generate unexpected patterns.

4.8 CONCLUSION

A theory of inductive learning has been presented that views such learning as a heuristic search through a space of symbolic descriptions, generated by an application of certain inference rules to the initial observational statements (teacher-generated examples of some concepts or environment-provided facts). The process of generating the goal description—the most preferred inductive assertion—relies on the universally intertwined and complementary operations of specializing or generalizing the currently-held assertion in order to accommodate new facts. The domain background knowledge has been shown to be a necessary component of inductive learning, which provides constraints, guidance, and a criterion for selecting the most preferred assertion.

Such a characterization of inductive learning is conceptually simple, and constitutes a theoretical framework for describing and comparing learning methods, as well as developing new methods. The Star methodology for learning structural descriptions from examples, described in the second part of this chapter, represents a general approach to concept acquisition which can be implemented in a variety of ways and applied to different problem domains.

There are many important topics of inductive learning that have not been covered here. Among them is learning from incomplete or uncertain information, learning from descriptions containing errors, learning with a multitude of forms of observational statements, as well as multimodel-based inductive assertions, and learning general rules with exceptions. The problem of discovering new concepts, descriptors and, generally, various many-level transformations of the initial description space (that is, the problem of constructive inductive learning) has been covered only very superficially.

These and related topics have been given little attention so far in the field of machine learning. There is no doubt, however, that as the understanding of the fundamental problems in the field matures, these challenging topics will be given increasing attention.

ACKNOWLEDGMENTS

In the development of the ideas presented here the author benefited from discussions with Tom Dietterich and Robert Stepp. Proofreading and comments of Jaime Carbonell, Bill Hoff and Tom Mitchell were helpful in shaping up the final version of the chapter. Comments and suggestions of the reviewers of *Artificial Intelligence Journal*, where the original version of this chapter was submitted and accepted for publication, helped to improve its clarity and organization.

The author gratefully acknowledges the partial support of the research by the National Science Foundation under grant MCS 82-05166, and the Office of Naval Research under grant N00014-82-K-0186.

REFERENCES

- Amarel, S., "An approach to automatic theory formation," *Illinois Symposium on Principles of Self-Organization*, H. von Foerster (Ed.), 1960.
- Banerji, R. B., "The description list of concepts," *J.A.C.M.*, 1962.
- Banerji, R. B., *Artificial Intelligence: A Theoretical Perspective*, Elsevier North Holland, New York, 1980.
- Biermann, A. W., "The inference of regular LISP programs from examples," *IEEE Transactions on Systems, Man, and Cybernetics*, Vol. SMC-8, No. 8, pp. 585-600, August 1978.
- Biermann, A. and Feldman, J., *A survey of results in grammatical inference*, Academic Press, New York, 1972.
- Bongard, N., *Pattern Recognition*, Spartan Books, New York, 1970, (Translation from Russian original, published in 1967).
- Brachman, R. J., "On the epistemological status of semantic networks," *Associative Networks*, N. V. Findler (Ed.), New York: Academic Press, 1979.
- Bruner, J. S., Goodnow, J. J. and Austin, G. A., *A Study of Thinking*, Wiley, New York, 1956.
- Buchanan, B. G. and Feigenbaum, E. A., "DENDRAL and Meta-DENDRAL: their applications dimension," *Artificial Intelligence*, Vol. 11, pp. 5-24, 1978.
- Buchanan, B. G., Mitchell, T. M., Smith, R. G. and Johnson, C. R. Jr., "Models of Learning Systems", Technical Report STAN-CS-79-692, Stanford University, Computer Science Dept., January 1979.

- Burstall, R. M. and Darlington, J., "A transformation system for developing recursive programs," *Journal of the ACM*, Vol. 24, No. 1, pp. 44-67, 1977.
- Carnap, R., "The aim of inductive logic," *Logic, Methodology and Philosophy of Science*, Nagel, E., Suppes, P. and Tarski, A. (Eds.), Stanford University Press, Stanford, pp. 303-318, 1962.
- Case, J. and Smith, C., "Comparison of identification criteria for mechanized inductive inference", Technical Report TR-154, Dept. Computer Science., State U. of New York at Buffalo, 1981.
- Chang, C., and Lee, R. C., *Symbolic Logic and Mechanical Theorem Proving*, Academic Press, New York, 1973.
- Cohen, B. L., "A powerful and efficient structural pattern recognition system," *Artificial Intelligence*, Vol. 9, No. 3, December 1977.
- Coulon, D. and Kayser, D., "Learning criterion and inductive behavior," *Pattern Recognition*, Vol. 10, No. 1, pp. 19-25, 1978.
- Davis, R. and Lenat, D. B., *Knowledge Based Systems in Artificial Intelligence*, McGraw Hill, New York, 1981.
- Dietterich, T., "Description of inductive program INDUCE 1.1", Technical Report (Internal), Department of Computer Science, University of Illinois, Urbana-Champaign, October 1978.
- Dietterich, T. G., "The methodology of knowledge layers for inducing descriptions of sequentially ordered events," Master's thesis, University of Illinois, Urbana, October 1979.
- Feigenbaum E. A., "The simulation of verbal learning behavior," *Computers and Thought*, Feigenbaum, E. A. and Feldman, J. (Eds.), McGraw-Hill Book Company, New York, NY, 1963.
- Fikes, R. E., Hart, P. E. and Nilsson, N. J., "Learning and executing generalized robot plans," *Artificial Intelligence*, Vol. 3, pp. 251-288, 1972.
- Gaines, B. R., "Maryanski's grammatical inferencer," *IEEE Trans on Computers*, Vol. C-28, pp. 62-64, 1979.
- Gaschnig, J., "Development of Uranium Exploration Models for Prospector Consultant System", Internal, SRI International, March 1980.
- Hájek, P. and Havránek, T., *Mechanizing Hypothesis Formation: Mathematical Foundations for a General Theory*, Springer-Verlag, 1978.
- Hájek, P., Havel, I., and Chytil, M., "The GUHA method of automatic hypothesis determination," *Computing*, No. 1, pp. 293-308, March 1966.
- Hayes-Roth, F., "A structural approach to pattern learning and the acquisition of classificatory power," *Proceedings of the First International Joint Conference on Pattern Recognition*, Washington, D. C., pp. 343-355, 1973.
- Hayes-Roth, F. and McDermott, J., "An interference matching technique for inducing abstractions," *Communications of the ACM*, Vol. 21, No. 6, pp. 401-410, 1978.
- Hedrick, C. L., *A Computer Program to Learn Production Systems Using a Semantic Net*, Ph.D. dissertation, Carnegie-Mellon University, July 1974, (Department of Computer Science).

- Hintzman, D. L., *The Psychology of Learning and Memory*, W. H. Freeman and Company, 1978.
- Hoff, B., Michalski, R. S., and Stepp, R., "INDUCE 2 - a program for learning structural descriptions from examples", Technical Report 82-5, Intelligent Systems Group, October 1982.
- Hovland, C. I., "A 'Communication Analysis' of Concept Learning," *Psychological Review*, pp. 461-472, November 1952.
- Hunt, E. B., Marin, J. and Stone, P. T., *Experiments in Induction*, Academic Press, New York, 1966.
- Jouannaud, J. P., and Kodratoff, Y., "An automatic construction of LISP programs by transformations of functions synthesized from their input-output behavior," *International Journal of Policy Analysis and Information Systems*, Vol. 4, No. 4, pp. 331-358, December 1980.
- Kemeny, T. G., "The use of simplicity in induction," *Psychological Review*, Vol. 62, No. 3, pp. 391-408, 1953.
- Kochen, M., "Experimental study of hypothesis formation by computer," *Proc. 1960 London Symp. on Information Theory*, 1960.
- Langley, P. W., Neches, R., Neves, D. and Anzai, Y., "A domain-independent framework for procedure learning," *Journal of Policy Analysis and Information Systems*, Vol. 4, No. 2, pp. 163-197, June 1980.
- Larson, J., *Inductive inference in the variable-valued predicate logic system VL21: methodology and computer implementation*, Ph.D. dissertation, University of Illinois, Urbana, Illinois, May 1977.
- Larson, J. and Michalski, R. S., "Inductive inference of VL decision rules," *Proceedings of the Workshop on Pattern Directed Inference Systems, SIGART Newsletter 63*, pp. 38-44, June 1977.
- Lenat, D. B., *AM: an artificial intelligence approach to discovery in mathematics as heuristic search*, Ph.D. dissertation, Stanford University, Stanford, California, 1976.
- Michalski, R. S., "A Variable-Valued Logic System as Applied to Picture Description and Recognition," *Graphic Languages*, F. Nake and A. Rosenfeld (Ed.), North-Holland Publishing Co., pp. 20-47, 1972.
- Michalski, R. S., "AQVAL/I - Computer implementation of a variable valued logic system VL1 and examples of its application to pattern recognition," *Proceedings of the First International Joint Conference on Pattern Recognition*, Washington, D. C., pp. 3-17, 1973b.
- Michalski, R. S., "Variable-Valued Logic and its Applications to Pattern Recognition and Machine Learning," *Multiple-Valued Logic and Computer Science*, Rine, D. (Ed.), North-Holland, pp. 506-534, 1975a.
- Michalski, R. S., "Synthesis of optimal and quasi-optimal variable-valued logic formulas," *Proceedings of the 1975 International Symposium on Multiple-Valued Logic*, Bloomington, Indiana, pp. 76-87, May 1975b.
- Michalski, R. S., "Pattern recognition as rule-guided inductive inference," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. PAMI-2, No. 4, pp. 349-361, 1980a.

- Michalski, R. S., "Knowledge Acquisition Through Conceptual Clustering: A Theoretical Framework and an Algorithm for Partitioning Data into Conjunctive Concepts," *Policy Analysis and Information Systems*, Vol. 4, No. 3, pp. 219-244, 1980c, (A Special Issue on Knowledge Acquisition and Induction).
- Michalski, R. S. and Chilausky, R. L., "Learning by being told and learning from examples: an experimental comparison of the two methods of knowledge acquisition in the context of developing an expert system for soybean disease diagnosis," *Policy Analysis and Information Systems*, Vol. 4, No. 2, pp. 125-160, June 1980, (Special issue on knowledge acquisition and induction).
- Michalski, R. S. and Larson, J. B., "Selection of most representative training examples and incremental generation of VL_1 hypotheses: the underlying methodology and the description of programs ESEL and AQ11", Technical Report 867, Computer Science Department, University of Illinois, 1978.
- Michalski, R. S., and Negri, P., "An Experiment on Inductive Learning in Chess End Games," *Machine Representation of Knowledge, Machine Intelligence 8*, E. W. Elcock and D. Michie (Ed.), Ellis Horwood, pp. 175-192, 1977.
- Michalski, R. S., Stepp, R., and Diday, E., "A recent advance in data analysis: clustering objects into classes characterized by conjunctive concepts," *Progress in Pattern Recognition*, L. Kanal and A. Rosenfeld (Ed.), North-Holland, Amsterdam, pp. 33-56, 1981.
- Michie, D., "New face of AI", Technical Report 33, University of Edinburgh, 1977.
- Minsky, M., "A framework for representing knowledge," *The Psychology of Computer Vision*, P. H. Winston (Ed.), McGraw-Hill, New York, ch. 6, pp. 211-277, 1975.
- Mitchell, T. M., *Version Spaces: An Approach to Concept Learning*, Ph.D. dissertation, Stanford University, December 1978.
- Moraga, C., "A didactic experiment in pattern recognition", Technical Report AIUD-PR-8101, Dartmund University, 1981.
- Morgan, C. G., "Automated hypothesis generation using extended inductive resolution," *Advance Papers of Fourth International Joint Conference on Artificial Intelligence*, Tbilisi, USSR, pp. 351-356, September 1975.
- Newell, A., Shaw, J. C. and Simon, H. A., "A variety of intelligent learning in a general problem solver," *Self Organizing Systems*, Yovits and Cameron (Eds.), Pergamon Press, New York, 1960.
- Nilsson, N. J., *Principles of Artificial Intelligence*, Tioga Publishing Co., 1980.
- O'Rourke, P., "A comparative study of inductive learning systems AQ11 and ID3", Intelligent Systems Group Report 82-2, Department of Computer Science, University of Illinois at Urbana-Champaign, 1982.
- Pettorossi, A., "An Algorithm for Reducing Memory Requirements in Recursive Programs Using Annotations," *International Workshop on Program Construction*, September 1980.
- Plotkin, G. D., "A further note on inductive generalization," *Machine Intelligence*, Meltzer, B. and Michie, D. (Eds.), Elsevier, Edinburgh, pp. 101-124, 1971.

- Pokorny, D., "Knowledge Acquisition by the GUHA Method," *International Journal of Policy Analysis and Information Systems*, Vol. 4, No. 4, pp. 379-399, 1980, (A special issue on knowledge acquisition and induction).
- Polya, G., *Mathematics and Plausible Reasoning*, Princeton University Press, Princeton, N.J., 1954.
- Popper, K., *The Logic of Scientific Discovery*, Harper and Row, New York, 1968, (2nd edition).
- Post, H. R., "Simplicity of Scientific Theories," *British Journal for the Philosophy of Science*, Vol. 11, No. 41, 1960.
- Quinlan, J. R., "Discovering rules from large collections of examples: a case study," *Expert Systems in the Micro Electronic Age*, Michie, D. (Ed.), Edinburgh University Press, Edinburgh, 1979.
- Russell, B., *History of Western Philosophy*, George Allen and Unwin, London, 1946.
- Sammut, C., *Learning Concepts by Performing Experiments*, Ph.D. dissertation, University of New South Wales, November 1981.
- Shapiro, Ehud Y., "Inductive Inference of Theories From Facts", Research Report 192, Yale University, February 1981.
- Shapiro, A. and Niblett, T., "Automatic Induction of classification rules for a chess endgame," *Advances in Computer Chess, volume 3*, Clarke, M.R.B. (Ed.), Edinburgh University Press, 1982.
- Shaw, D. E., Swartout, W. R. and Green, C. C., "Inferring LISP programs from examples," *Fourth International Joint Conference on Artificial Intelligence*, Tbilisi, USSR, pp. 351-356, September 1975.
- Shortliffe, E., *Computer Based Medical Consultations: MYCIN*, New York: Elsevier, 1976.
- Simon, H. A. and Kotovsky, K., "Human acquisition of concepts for sequential patterns," *Psychological Review*, Vol. 70, pp. 534-546, 1963.
- Simon, H. A. and Lea, G., "Problem solving and rule induction: A unified view," *Knowledge and Cognition*, L. Gregg (Ed.), Lawrence Erlbaum Associates, Hillsdale, N.J., 1974.
- Smith, D. R., "A Survey of the Synthesis of LISP Programs from Examples", Technical Report, Duke University, Bonas, France, September 1980.
- Solomonoff, R. J., "A Formal Theory of Inductive Inference," *Information and Control*, Vol. 7, 1964.
- Soloway, E. M. and Riseman, E. M., "Levels of pattern description in learning," *Fifth International Joint Conference on Artificial Intelligence*, Cambridge, Mass., pp. 801-811, 1977.
- Stepp, R., "The investigation of the UNICLASS inductive program AQ7UNI and User's Guide", Technical Report 949, Department of Computer Science, University of Illinois, Urbana, Illinois, November 1978.
- Stoffel, J. C., "The theory of prime events: data analysis for sample vectors with inherently discrete variables," *Information Processing 74*, North-Holland, Amsterdam, pp. 702-706, 1974.
- Suppes, P., *Introduction to Logic*, Van Nostrand Co., Princeton, 1957.

- Vere, S. A., "Induction of concepts in the predicate calculus," *Proceedings of the Fourth International Joint Conference on Artificial Intelligence, IJCAI, Tbilisi, USSR, 1975.*
- Waterman, D. A., "Generalized learning techniques for automating the learning of heuristics," *Artificial Intelligence, Vol. 1, No. 1-2, pp. 121-170, Spring 1970.*
- Winston, P., *Learning Structural Descriptions from Examples*, Ph.D. dissertation, MIT, September 1970.
- Winston, P. H., *Artificial Intelligence*, Addison-Wesley, 1977.
- Yau, K. C., and Fu, K. S., "Syntactic shape recognition using attributed grammars," *Proceedings of the Eighth Annual EIA Symposium on Automatic Imagery Pattern Recognition, 1978.*
- Zagoruiko, N. G., *Mietody obnaruzhenia zakonomiernostiej (Methods for revealing regularities in data)*, Izd. Nauka, Moscow, 1981.

APPENDIX: ANNOTATED PREDICATE CALCULUS (APC)

This appendix presents definitions of the basic components of the annotated predicate calculus and some rules for equivalence-preserving transformations of APC expressions (rules that are nonexistent in the ordinary calculus).

1. Elementary and Compound Terms—*Terms* can be elementary or compound. An *elementary term* (an *eterm*) is the same as a term in predicate calculus, that is, a constant, a variable, or a function symbol followed by a list of arguments that are eterms. A *compound term* (*cterm*) is a *composite* of elementary terms or is an eterm in which one or more arguments are such composites. The composite of eterms is defined as the *internal conjunction* (&) or *internal disjunction* (\vee) of eterms. (The meaning of these operators is explained later.) The following are examples of compound terms:

$$\text{RED} \vee \text{BLUE} \quad (1)$$

$$\text{height}(\text{BOX}_1 \ \& \ \text{BOX}_2) \quad (2)$$

where RED, BLUE, BOX₁, BOX₂ are constants. Expression (1) and the form in parentheses in (2) are composites. Note that expressions (1) and (2) are not logical expressions that have a truth-status (that is, can be true or false); they are terms to be used only as arguments of predicates. A compound term in which arguments are composites can be transformed (expanded) into a composite of elementary terms. Let *f* be an *n*-argument function whose *n*-1 arguments are represented by list *A*, and let *t*₁ and *t*₂ be elementary terms. The rules for performing such a transformation, that is, term rewriting rules, are:

$$f(t_1 \vee t_2, A) \leftrightarrow f(t_1, A) \vee f(t_2, A) \quad (3)$$

$$f(t_1 \ \& \ t_2, A) \leftrightarrow f(t_1, A) \ \& \ f(t_2, A) \quad (4)$$

Thus, term (2) can be transformed into a composite:

$$\text{height}(\text{BOX}_1) \ \& \ \text{height}(\text{BOX}_2) \quad (5)$$

If list *A* itself contains composites, then it is assumed that the internal disjunction is expanded first, followed by the internal conjunction (that is, the conjunction binds stronger than the disjunction).

2. Elementary and Compound Predicates—Predicates also can be elementary or compound. An *elementary predicate* is the same as a predicate in the predicate calculus, that is, a predicate symbol followed by a list of arguments that are eterms. In a *compound predicate* one or more arguments is a compound term. For example, the following are compound predicates:

$$\text{Went}(\text{Mary} \ \& \ \text{Mother}(\text{Stan}), \text{Movie} \ \vee \ \text{Theater}) \quad (6)$$

$$\text{Inside}(\text{Key}, \text{Drawer}(\text{Desk}_1 \ \vee \ \text{Desk}_2)) \quad (7)$$

The meaning of a compound predicate is defined by rules for transforming it into an expression made of elementary predicates and ordinary “external” logic operators of conjunction (&) and disjunction (\vee). We denote the internal and

external operators identically, because they can be easily distinguished by the context (note that there is no distinction between them in natural language). If an operator connects predicates, then it is an external operator; if it connects terms, then it is an internal operator.

Let t_1 and t_2 be *eterms* and P an n -ary predicate whose last $n-1$ arguments are represented by a list A . We have the following reformulation rules (that is, equivalence-preserving transformations of descriptions):

$$P(t_1 \vee t_2, A) \quad | = \quad P(t_1, A) \vee P(t_2, A) \quad (8)$$

$$P(t_1 \& t_2, A) \quad | = \quad P(t_1, A) \& P(t_2, A) \quad (9)$$

If an argument of a predicate is a compound term that is not a composite of elementary terms, then it is transformed first into a composite by rules (3) and (4). If A contains a composite of terms, then the disjunction is expanded first before conjunction (similarly as in expanding compound terms).

Rules (3), (4), (8) and (9) can be used as bidirectional transformation rules. By applying them forward (from left to right), a compound predicate can be *expanded* into an expression containing only elementary predicates, and by applying them backward, an expression with elementary predicates can be contracted into a compound predicate.

For example, by applying forward rule (8) and then (9), one can expand the compound predicate (6) into

$$\begin{aligned} & \text{Went}(\text{Mary}, \text{movie}) \& \text{Went}(\text{Mother}(\text{Stan}), \text{movie}) \vee \\ & \text{Went}(\text{Mary}, \text{theater}) \& \text{Went}(\text{Mother}(\text{Stan}), \text{theater}) \end{aligned} \quad (10)$$

Comparing logically-equivalent expressions (6) and (10), one can notice that expression (6) is considerably shorter than (10), and in contrast to (10), represents explicitly the fact that Mary & Mother(Stan) went to the same place. Also, the structure of (6) is more similar to the structure of the corresponding natural language expression.

3. Relational Statements—A simple and often used way of describing objects or situations is to state the values of selected attributes applied to these objects or situations. Although such information can be represented by predicates, this is not the most readable or natural way. The APC uses for this purpose a statement:

$$\text{eterm}_i = a \quad (11)$$

stating that eterm_i evaluates to a constant a . Such a statement is called an *atomic relational statement* (or an *atomic selector*). Expression (11) is a special case of a *relational statement* (also called *selector*), defined as:

$$\text{Term}_1 \text{ rel Term}_2 \quad (12)$$

where Term_1 and Term_2 are elementary or compound terms, and *rel* stands for one of the relational symbols: =, \geq , $>$, \leq , $<$.

If Term_1 and Term_2 are both elementary, then expression (12) states that the value of the function represented by Term_1 is in relation *rel* to the value of function represented by Term_2 . For example, the expression:

$$\text{distance}(\text{Boston}, \text{Tampa}) = \text{distance}(\text{Washington}, \text{Dallas}) \quad (13)$$

states that the distance between Boston and Tampa is the same as the distance between Washington and Dallas. If Term_2 is a constant, then it evaluates to itself.

Expression (12) can be represented by a predicate:

$$\text{rel}(\text{Term}_1, \text{Term}_2) \quad (14)$$

If Term_1 and/or Term_2 is compound, then the meaning of expression (12) is defined by expanding it into a form containing only relational statements with elementary terms. The expansion is performed by transforming expression (12) into (14), applying transformation rules (3), (4), (8), and (9), and then converting the elementary predicates into relational statements.

For example, a relational statement:

$$\text{color}(P_1 \vee P_2) = \text{Red} \vee \text{Blue} \quad (15)$$

can be expanded into an expression:

$$(\text{color}(P_1) = \text{Red} \vee \text{Blue}) \vee (\text{color}(P_2) = \text{Red} \vee \text{Blue}) \quad (16)$$

and finally to an expression consisting of only atomic selectors:

$$\begin{aligned} &(\text{color}(P_1) = \text{Red}) \vee (\text{color}(P_1) = \text{Blue}) \vee \\ &(\text{color}(P_2) = \text{Red}) \vee (\text{color}(P_2) = \text{Blue}) \end{aligned} \quad (17)$$

The two selectors in the disjunction (16) are examples of a *referential selector*, defined as a form:

$$\text{Term}_1 \text{ rel Term}_2 \quad (18)$$

where Term_1 (called *referee*) is a nonconstant elementary term and Term_2 (called *reference*) is a constant or the internal disjunction of constants from the domain of Term_1 . If relation *rel* is "=" and Term_2 is the disjunction of some constants, then the referential selector (18) states that the function represented by Term_1 evaluates to one of the constants in Term_2 . The referential selector is very useful for representing concept descriptions.

If the reference of a referential selector contains a sequence of consecutive constants from the domain of a linear descriptor, then the range operator ".." is used to simplify the expression. For example:

$$\text{size}(P) = 2 \vee 3 \vee 4$$

can be written:

$$\text{size}(P) = 2..4$$

The negation of a selector:

$$\sim(\text{Term}_1 = \text{Term}_2)$$

can be equivalently written:

$$\text{Term}_1 \neq \text{Term}_2 \quad (20)$$

An arbitrary predicate $P(t_1, t_2, \dots)$ can be written in the form of a referential selector:

$$P(t_1, t_2, \dots) = \text{True} .$$

Therefore, for the uniformity of terminology, a predicate will be considered a special form of a selector.

To facilitate the interpretation and readability of individual selectors in expressions, they are usually surrounded by square brackets and their conjunction is expressed by concatenating the bracketed forms (see Section 4.7).

APC expressions are created from selectors (relational statements) in the same way as predicate calculus expressions are created from predicates, that is, by using logic connectives (\sim , $\&$, \vee , \Rightarrow , \Leftrightarrow) and quantifiers. One additional useful connective is the *exception operation* (" \backslash "), defined as:

$$S_1 \backslash S_2 \quad | = \quad (\sim S_2 \Rightarrow S_1) \& (S_2 \Rightarrow \sim S_1) \quad (21)$$

where S_1 and S_2 are APC expressions. ($S_1 \backslash S_2$ reads: S_1 *except when* S_2 .) It is easy to see that the exception operator is equivalent to the symmetrical difference.

In addition to ordinary quantifiers there is also a *numerical quantifier*, expressed in the form:

$$\exists(I) v, S[v] \quad (22)$$

where I , the *index set*, denotes a set of integers, and $S[v]$ is an APC expression having v as a free variable.

Sentence (22) evaluates as true if the number of values of v for which expression $S[v]$ is true is an element of the set I . For example, formula:

$$\exists(2..8) v, S[v] \quad (23)$$

states that there are two to eight values of v for which the expression $S[v]$ is true. The following equivalences hold:

$$\exists v, S[v] \text{ is equivalent to } \exists(\geq 1) v, S[v]$$

and

$$\forall v, S[v] \text{ is equivalent to } \exists(k) v, S[v]$$

where k is the number of possible values of variable v .

To state that there are k and only k distinct values for variables v_1, v_2, \dots, v_k for which expression $S[v_1, v_2, \dots, v_k]$ is true, we write:

$$\exists.v_1, v_2, \dots, v_k, S[v_1, \dots, v_k] \quad (24)$$

For example, the expression:

$$\exists.P_0, P_1, P_2 [\text{contains}(P_0, P_1 \& P_2)] \& [\text{color}(P_1 \& P_2) = \text{red}] \Rightarrow [\text{two-red-parts}(P_0)]$$

states that predicate $\text{two-red-parts}(P_0)$ holds if P_0 has two, and only two, distinct parts in it that are red.

Section 4.7 presents an example of the usage of the APC for formulating observational statements and concept descriptions.