

Scientific Model-Building as Search in Matrix Spaces

Raúl E. Valdés-Pérez
School of Comp. Sci. &
Center for Light Microscope
Imaging and Biotechnology
Carnegie Mellon Univ.

Jan M. Zytkow
Dept. of Comp. Sci.
Wichita State Univ.

Herbert A. Simon
Dept. of Psychology
Carnegie Mellon Univ.

Abstract

Many reported discovery systems build discrete models of hidden structure, properties, or processes in the diverse fields of biology, chemistry, and physics. We show that the search spaces underlying many well-known systems are remarkably similar when re-interpreted as search in matrix spaces. A small number of matrix types are used to represent the input data and output models. Most of the constraints can be represented as matrix constraints; most notably, conservation laws and their analogues can be represented as matrix equations. Typically, one or more matrix dimensions grow as these systems consider more complex models after simpler models fail, and we introduce a notation to express this. The novel framework of matrix-space search serves to unify previous systems and suggests how at least two of them can be integrated. Our analysis constitutes an advance toward a generalized account of model-building in science.

Introduction

The discovery of models of atomic and molecular structure, of chemical processes, and of genetic transmission are celebrated events in the history of science. Far from being isolated historical instances, discovery of hidden structure in the form of discrete models is a universal and current task across the natural sciences.

Several discovery systems reported in the AI literature discover models of discrete, hidden structure. These systems include DALTON [Langley *et al.*, 1987], GELL-MANN [Fischer and Zytkow, 1990], MECHEM [Valdes, 1992; 1993 (in press)], MENDEL [Fischer and Zytkow, 1992], BR-3/PAULI [Kocabas, 1991; Valdes, accepted], STAHL [Zytkow and Simon, 1986] and its descendants STAHLp [Rose and Langley, 1986] and REVOLVER [Rose, 1989]. Of these, DALTON, MECHEM, and STAHL perform in chemistry, GELL-MANN and BR-3/PAULI in physics, and MENDEL in biology.

Despite the diversity of scientific domains that these systems treat, there lurk striking similarities in the rep-

resentation of models, problem-solving methods, and domain knowledge used in model construction. Some of these similarities were pointed out elsewhere [Fischer and Zytkow, 1992]. These similarities may eventually allow us to develop a unified discovery system able to search for many types of discrete models. As a prerequisite, we should study existing systems that have already demonstrated a degree of competence on historical or current science. An important theoretical task of comparative analysis, which is relatively scarce in the AI literature, is to identify a unitary core among these systems. Without this, progress is limited to the accumulation of special-purpose programs.

The purpose of this paper is to identify a common representation of discrete models and a systematic analysis of the search spaces and domain constraints using the language of matrices. Our analysis introduces a small set of matrix types that represent the input data, the output models, and the spaces to be searched by the discovery system. Models are proposed by assigning numeric values to entries in a matrix, most assignments being ruled out by the domain constraints. The matrix representation enables the use of powerful methods of matrix algebra and combinatorial algorithms to improve the search for discrete models.

We also introduce a new notation to express how discovery systems carry out the search for models by postulating new entities, processes, and properties. This notation is used later to show how two specific systems that were developed separately can be integrated.

Systems

In this section we will re-interpret six discovery systems and show that they have a surprising degree of similarity. Three types of matrices used in these systems will be highlighted: a reaction matrix \mathcal{R} , a structure matrix \mathcal{S} , and a property matrix \mathcal{P} , defining them in the context of each system. We use the language of matrices and matrix algebra to describe the constraints in these systems. We also show how each system systematically changes the sizes of some few matrices in the course of performing its task.

The emphasis in this paper is on the spaces searched

\mathcal{R}	hydrogen	nitrogen	oxygen	ammonia	water
R_1	[<0]	0	[<0]	0	[>0]
R_2	[<0]	[<0]	0	[>0]	0

$R_1 = \text{react}(\text{hydrogen}, \text{oxygen}) \rightarrow \text{water}$
 $R_2 = \text{react}(\text{hydrogen}, \text{nitrogen}) \rightarrow \text{ammonia}$

Figure 1: Reaction Matrix in DALTON

\mathcal{S}	N	O	H
hydrogen	0	0	[>0]
nitrogen	[>0]	0	0
oxygen	0	[>0]	0
ammonia	[>0]	0	[>0]
water	0	[>0]	[>0]

Figure 2: Structure Matrix in DALTON

by the systems, and not on the detailed ways each system carries out its search, which varies across systems and sometimes even within systems, since several of the programs possess more than one internal search space. One view of problem-solving in science is that it typically proceeds over several spaces which can be quite heterogeneous. Initially proposed by Lea and Simon [Lea and Simon, 1974], this idea has been expanded and applied in the discovery system FAHRENHEIT [Zytrow, 1987], while Klahr and Dunbar [Klahr and Dunbar, 1988] have investigated it as a psychological model.

Some comments on notation follow. Matrices will be represented as tables with two intersecting perpendicular line segments, one to mark the rows, the other the columns. Additional marks are used to show whether a matrix dimension grows, shrinks, or is static: an outgoing (ingoing) arrow means that the dimension grows (shrinks), and a cap means that it is static during problem solving. We will see that most of the systems progressively enlarge their matrix models when smaller models prove inadequate.

DALTON

DALTON's task is to find structural models of chemical reactions and substances in terms of atoms [Langley *et al.*, 1987]. For example, given the following data:

1. two volumes of hydrogen and one volume of oxygen react to form two volumes of water;
2. three volumes of hydrogen and one volume of nitrogen react to form two volumes of ammonia;
3. hydrogen, oxygen, and nitrogen are elementary substances;
4. water consists of hydrogen and oxygen, and ammonia consists of hydrogen and nitrogen;

DALTON uses its bias for simplicity, conservation laws, and the Gay-Lussac law to report correctly that (1) two hydrogen molecules react with one of oxygen to form two water molecules, while three hydrogen molecules combine with one nitrogen molecule to form two ammonia molecules, and that (2) hydrogen, oxygen and nitrogen are diatomic, water is H_2O and ammonia is NH_3 .

In making these inferences, DALTON can be interpreted as filling in two matrices. The first matrix describes inputs and outputs for each reaction; the example discussed in [Langley *et al.*, 1987] has the initial

form shown in Figure 1. The bracketed constraints represent conventional matrix depictions of reactions [Aris and Mah, 1963]: the reactants have negative entries, and the products have positive entries. All non-participating substances have zero entries. In this paper, we will always denote such a reaction matrix by $\mathcal{R}_{r \times s}$, where r is the number of reactions, and s is the number of chemical substances.

A second, structure matrix in DALTON (Figure 2) represents the structure of the chemical substances in terms of atomic elements. Initially, some of the entries are zero to indicate that certain substances do not contain certain atoms. The remaining entries in the matrix are constrained to positive integers. We denote this structure matrix as $\mathcal{S}_{s \times e}$, where s as before is the number of substances, and e is the number of chemical elements involved.

The sizes of the \mathcal{R} and \mathcal{S} matrices are fixed, as indicated in the figures by a "double cap" notation that prevents the matrix from changing size. DALTON does not conjecture new reactions, substances, nor chemical elements, so it never enlarges the two matrices which it receives as input.

DALTON's task is to fill in the reaction matrix and the structure matrix completely with integer entries, subject to the constraints stated above, a criterion of simplicity of entries, and a conservation law on atoms, which is expressed in matrix algebra as follows:

$$\mathcal{R}_{r \times s} \times \mathcal{S}_{s \times e} = \mathbf{0}_{r \times e} \quad (1)$$

This equation implies that each of r reactions must conserve the atoms of all e elements: the product $\mathcal{R} \times \mathcal{S}$ gives the zero matrix $\mathbf{0}$ of dimensions $r \times e$. Conservation means that the net production of atoms of each element is zero for each reaction. Simplicity has a role in choosing the magnitudes of the entries (integers of smaller magnitude are simpler). Equation 1 is the standard way to express linear conservation conditions in sciences such as chemistry.

In our example, DALTON outputs the two matrices in Figure 3 (the output matrix \mathcal{R} is shown transposed to fit on the page). The matrix \mathcal{R} quantifies the qualitative reaction matrix input to DALTON, e.g., three hydrogen molecules enter into the ammonia reaction. The output matrix \mathcal{S} specifies the elementary constituents of each substance. For example, a value of 2 for the matrix entry (*hydrogen*, *H*) in \mathcal{S} means that hydrogen molecules include two atoms of hydrogen. Since all other entries in the hydrogen row are

\mathcal{R}^T	water reaction	ammonia reaction		
hydrogen	-2	-3		
nitrogen	0	-1		
oxygen	-1	0		
ammonia	0	2		
water	2	0		
<hr/>				
\mathcal{S}	N	O	H	
hydrogen	0	0	2	
nitrogen	2	0	0	
oxygen	0	2	0	
ammonia	1	0	3	
water	0	1	2	

Figure 3: Output of DALTON

\mathcal{S}	quark ₁	...	quark _e	
particle ₁	↓			→
⋮				
particle _s				
<hr/>				
\mathcal{P}	property ₁	...	property _p	
quark ₁	↓			
⋮				
quark _e				

Figure 4: Matrix Structure of GELL-MANN

zero, the hydrogen molecule is diatomic, i.e., has structure H_2 .

GELL-MANN

GELL-MANN's task is to propose quark models that account for the known property values of the particles in elementary-particle families [Fischer and Zytlow, 1990]. The models constructed by GELL-MANN are filled-in pairs of matrices shown in Figure 4. The structural \mathcal{S} matrix is analogous to the \mathcal{S} matrix in DALTON: s particles (or "substances") will contain e quarks (or "elements"). The second matrix in GELL-MANN is a property matrix \mathcal{P} which assigns values of p properties to e quarks. The domain constraints on the \mathcal{S} matrix are:

- The matrix entries are non-negative integers.
- The sum of entries over each row equals k , which is the number of quarks contained in each particle.
- The number of k -combinations of the set of e quarks (with infinite repetition number), which by a theorem in elementary combinatorics [Brualdi, 1981] equals $C(e-1+k, k)$, satisfies $s \leq C(e-1+k, k) \leq 3s$, where s is the number of input particles.

Contrary to DALTON, GELL-MANN enlarges the number of columns in the first matrix (and performs the

\mathcal{R}	starting materials	observed products	conjectured substances
reaction ₁	↓		
⋮			
reaction _r			

\mathcal{S}	C	H	N	O ...
starting materials	↓			
observed products				
conjectured substances				

Figure 5: Matrix Structure of MECHEM

number of rows in the second) if it cannot find an acceptable model for the current number of quarks. Adding another column corresponds to postulating one more quark. During its search for an acceptable model, GELL-MANN also increments the value of k , the number of quarks per particle. Each k leads to $C(e-1+k, k)$ possible quark combinations, each represented by one row in the expanded \mathcal{S} matrix. The number of input particles is constant, and equals s .

The quark models proposed by GELL-MANN must also be consistent with the observed property values of the particles. For example, since a proton has a charge of 1, the sum of charges for quarks which constitute the proton must be also 1. This constraint is called an "additivity law" in [Fischer and Zytlow, 1990], and is analogous to laws of conservation. Whereas conservation in DALTON (and generally) is expressed by a constraint of the form $\mathcal{R} \times \mathcal{S} = 0$, additivity in GELL-MANN is expressed as $\mathcal{S}_{s \times e} \times \mathcal{P}_{e \times p} = \mathcal{P}'_{s \times p}$. The matrix \mathcal{P}' contains property values of particles, which are constants given as input to the system. Matrices \mathcal{P} and \mathcal{P}' both contain property values: the first for hidden objects postulated in the model, the second for observable objects given in the input.

Those rows in GELL-MANN's \mathcal{S} matrix corresponding to particles input to the program are tested using the additivity law. However, GELL-MANN also predicts unseen particles by taking advantage of those quark combinations (numbering $C(e-1+k, k) - s$) that were not used to model the known particles. In these cases, the properties of these new particles are predicted by simply pre-multiplying the matrix \mathcal{P} by these $C(e-1+k, k) - s$ rows.

MECHEM

MECHEM's task is to discover the simplest pathway able to explain all the experimental evidence about an aggregate chemical reaction [Valdes, 1992; 1993 (in press)]. MECHEM searches the space of two matrices shown in Figure 5. Some constraints on the \mathcal{R} matrix are:

1. matrix entries admit only integer values,
2. For each row, the sum of the negative integers is -1 or -2 . The sum of the positive integers is 1 or 2

[Each reaction has at most two reactants and two products].

- Each column contains at least one nonzero entry [All substances must occur somewhere in the reaction].
- For each column corresponding to observed or conjectured products, the top-most nonzero entry is positive [Each product must be formed before it can be consumed].

The fourth constraint is used to define a canonical order on reactions in the service of search efficiency [Valdes, 1991]; it is not derived from chemical theory.

New rows and columns can be added in the \mathcal{R} matrix as the program fails with simpler hypotheses, so we see that MECHEM has two dimensions of expansion that guide its search in this matrix space. MECHEM prefers adding new reactions (by incrementing r) over incrementing the number of conjectured substances, so usually the matrix is growing vertically. Three systems considered in this paper (MECHEM, MENDEL, and GELL-MANN) enlarge matrices along two dimensions.

In the \mathcal{S} matrix, the molecular formulas for the starting materials and observed products are known; the program determines the formulas, or matrix entries, for the conjectured substances. This task is common to all systems which fill in entries in the \mathcal{S} matrix. As in DALTON, the conservation conditions can be expressed as the equation $\mathcal{R}_{r \times s} \times \mathcal{S}_{s \times e} = \mathbf{0}_{r \times e}$, in which \mathcal{S} is a structure matrix that contains the molecular formula (involving e chemical elements) of each substance, and $\mathbf{0}_{r \times e}$ is the zero matrix.

In addition to conservation of the elements, MECHEM incorporates other chemical constraints that arise from properties of substances, such as free energy or oxidation number. These constraints can be interpreted as an equation $\mathcal{R}_{r \times s} \times \mathcal{P}_{s \times p} = \mathcal{Z}_{r \times p}$, in which the constraint on the entries of the p columns of \mathcal{Z} vary according to the property. For example, in certain oxidation reactions, the oxidation number should never decrease across a reaction, so all the entries under the oxidation-number column of \mathcal{Z} would be non-negative.

The above are not the only search spaces in MECHEM. For example, to perform its task at a modern level of competence, molecular structures must be inferred for the conjectured substances, not only formulas. The space of molecular structures can also be represented as a matrix, similar to the search space in DENDRAL [Lindsay *et al.*, 1980].

MENDEL

MENDEL's task is to devise genetic explanations for observed inheritance patterns (or "reactions") among phenotypes [Fischer and Zytkow, 1992]. Each phenotype is explained by one or more genotypes. MENDEL searches the pair of matrices \mathcal{R} and \mathcal{S} in Figure 6, in analogy to the matrix \mathcal{R} in DALTON, MECHEM,

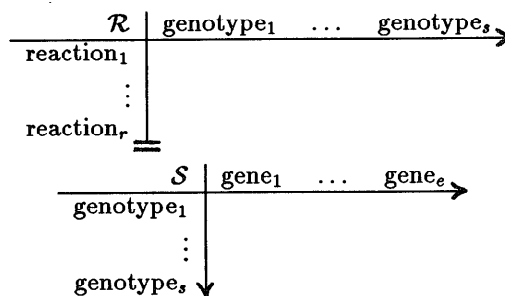


Figure 6: Matrix Structure of MENDEL

and STAHL, and in analogy to \mathcal{S} in DALTON, GELL-MANN, and MECHEM. The domain constraints on \mathcal{S} are identical to GELL-MANN's:

- The matrix entries admit only non-negative integers.
- The sum of entries over each row equals k , which is the number of genes making up a genotype.
- The number $s \stackrel{\text{def}}{=} C(e-1+k, k)$ of possible genotypes having k genes (analogous to the constraint in GELL-MANN) satisfies the constraint $f \leq C(e-1+k, k) \leq 3f$, where f is the fixed number of input phenotypes.

MENDEL enlarges the number of columns in \mathcal{S} if it cannot find explanations of genetic reactions within a specific number of genes. Adding one more column to the matrix corresponds to postulating one more gene. MENDEL, like GELL-MANN, carries out a subordinate search by varying the values of the parameter k , which together with the number e of genes leads to postulating $C(e-1+k, k)$ genotypes; these determine the number of rows in the \mathcal{S} matrix and columns in the \mathcal{R} matrix. MENDEL's search for gene combinations into genotypes is similar to GELL-MANN and DALTON, although several genotypes may be needed to explain one phenotype and several genotype reactions may be needed to explain one phenotype reaction.

The relative number of reactions between genotypes which look identical at the phenotype level is acceptable when it is approximately equal to the observed inheritance statistics that govern mating between phenotypes. Rather than using a predefined conservation principle, MENDEL searches for the right conservation/principle for genetic reactions, and finds out that one gene per parent is preserved in each offspring.

Since the same genotype can occur on both sides of a genetic reaction, and the occurrences should not cancel out, the entries in the \mathcal{R} matrix need to be pairs (n_r, n_p) , where n_r is the number of reactants and n_p is the number of products of a particular genotype.

BR-3/PAULI

PAULI's goal [Valdes, accepted], like its predecessor BR-3 [Kocabas, 1991], is to postulate a small number of

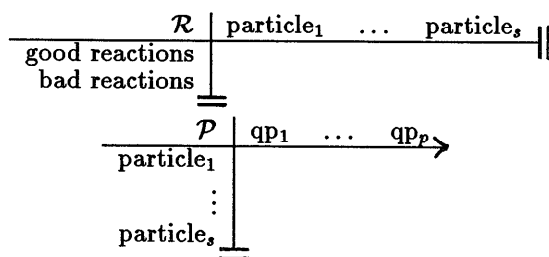


Figure 7: Matrix Structure of BR-3/PAULI

\mathcal{R}	lime	quick lime	fixed air	magnesia alba	calcined magnesia
R ₁	-	+	+	0	0
R ₂	+	-	0	-	+

R₁ = {lime} → {quick lime, fixed air}

R₂ = {quick lime, magnesia alba} → {lime, calcined magnesia}

Figure 8: Reaction Matrix of STAHL

quantum properties, together with values for the properties for each known elementary particle. New properties must explain how certain reactions in physics do not occur and how others occur. The ("good") reactions that occur must conserve each of the postulated properties, while every ("bad") disallowed reaction must disconserve at least one of the properties.

PAULI's matrix search space is shown in Figure 7. A filled-in \mathcal{R} matrix is input to the program. PAULI fills in the \mathcal{P} matrix, and enlarges its number of columns, i.e., the number of quantum properties that it postulates, when simpler models fail.

The only constraint that applies directly to the \mathcal{P} matrix is that the quantum properties of particle/antiparticle pairs should be of equal magnitude and opposite sign. Further constraints on solutions involve both conservation and disconservation. Letting g and b denote the "good" and "bad" reactions respectively, the following matrix equation must be satisfied:

$$\begin{bmatrix} \mathcal{R}_{g \times s} \\ \mathcal{R}_{b \times s} \end{bmatrix} \times P_{s \times p} = \begin{bmatrix} 0_{g \times p} \\ \mathcal{Z}_{b \times p} \end{bmatrix}$$

The first matrix is input to the program, the sub-matrix $0_{g \times p}$ has only zero entries, and the sub-matrix $\mathcal{Z}_{b \times p}$ enforces the disconservation: each row of \mathcal{Z} must contain a nonzero entry. Like GELL-MANN, BR-3 and PAULI could predict many unseen good and bad reactions by combining particles in various ways and testing whether conservation of all properties holds.

STAHL

Unlike other systems, the STAHL program of [Zytkow and Simon, 1986] discovers qualitative models rather than quantitative ones. Consequently, to describe STAHL's search problem we use qualitative matrix entries rather than numbers.

\mathcal{S}	quick lime	fixed air	calcined magnesia
lime	+	+	0
quick lime	+	0	0
fixed air	0	+	0
magnesia alba	0	+	+
calcined magnesia	0	0	+

Figure 9: Structure Matrix of STAHL

We use an example from page 128 of [Zytkow and Simon, 1986] for illustration. The input to the program consists of qualitative input/output facts about chemical reactions shown in Figure 8. A negative entry '-' corresponds to a reactant, a positive entry '+' corresponds to a reaction product, while any non-participating substance receives a zero entry. To represent reaction schemes in which the same substance occurs both as a reactant and a product, pairs of signs can be used, e.g., (-, +).

STAHL's task is to discover the elements and the make-up of substances in terms of these elements, i.e., an \mathcal{S} matrix. In the above example, from the first reaction STAHL notices that lime consists of quick lime and fixed air, and then combining the first and the second, that magnesia alba consists of calcined magnesia and fixed air. In effect the \mathcal{S} matrix in Figure 9 is created. If two rows in the \mathcal{S} matrix have the same entries, STAHL concludes that two substances having different names are in fact identical. In such a case, one row in \mathcal{S} (and the column in \mathcal{R}) can be deleted to give a simpler model; STAHL is the only system in this paper that can be viewed to shrink matrices. The columns of \mathcal{S} can be viewed either as growing and shrinking, or as only shrinking from a maximal possible set of elementary substances. The number of rows in \mathcal{R} grows, since STAHL makes "new" reactions from arithmetic combinations of known ones.

STAHL's \mathcal{R} and \mathcal{S} matrices satisfy a qualitative conservation principle: each element which occurs in a reaction should appear both in its reactants and in its products. This can be expressed identically to DALTON and MECHEM as $\mathcal{R}_{r \times s} \times \mathcal{S}_{s \times e} = 0_{r \times e}$, where matrix multiplication uses qualitative arithmetic following expected rules, for instance $pos \times neg = neg$, $pos + pos = pos$, $pos \times 0 = 0$. The qualitative arithmetic is not associative (e.g., $pos + pos + neg$ could equal pos or 0), but the order of production-rule firing determines how expressions are simplified.

Contradictions can arise when the product $\mathcal{R} \times \mathcal{S}$ has nonzero entries. Such nonzero entries indicate reactions which according to current knowledge (and STAHL's qualitative arithmetic) are unbalanced.

Discussion

The six systems examined in this paper propose discrete underlying models of empirical phenomena across a variety of tasks and sciences. All of the systems

find models of either the structure or properties of substances; this is the main task of DALTON, GELL-MANN, MENDEL, BR-3/PAULI, and STAHL. In addition, DALTON, MECHEM, and MENDEL find models of processes (reactions) in terms of hidden objects. DALTON takes a set of qualitative reactions and specifies them quantitatively, while MECHEM finds a simplest set of reactions (a pathway) from scratch.

All of the systems fill in the entries of one or more matrices. All except DALTON enlarge one or more matrix dimensions, and all including DALTON use constraints expressible as matrix equations of the form $AB = C$ or weaker forms of conservation.¹ The concept of simplicity has a strong presence, as reflected especially by growth in the matrices, which corresponds to entertaining more complex models.

Three matrix types are observed to recur. The most frequent is the reaction matrix \mathcal{R} , which appears in all of the systems except GELL-MANN. Either the structure matrix \mathcal{S} or the property matrix \mathcal{P} appears in all of the systems; all except DALTON and STAHL postulate either new objects or new properties.

Other discovery systems

DENDRAL [Lindsay *et al.*, 1980] and TETRAD-II [Spirtes *et al.*, 1990] discover models of molecular and causal structure, respectively. These models are in the form of graphs, which as is well known can always be represented as adjacency matrices. However, these two systems use different matrix types and different constraints than the ones discussed here, so we have not included them in the analysis of this paper.

AM [Lenat, 1982], GRAFFITI [Fajtlowicz, 1988], and BACON [Langley *et al.*, 1987] are other notable discovery systems that do not seem to fit the present framework. AM and GRAFFITI find plausible mathematical conjectures in elementary number theory and graph theory. BACON finds descriptive, empirical laws in data. These programs make inductive generalizations and introduce new theoretical terms, but do not build discrete models of hidden structure.

What is gained?

It is always possible to view one thing as another thing. A better understanding of a subject is often claimed as a virtue of a new viewpoint. However, since "understanding" is a slippery notion, it is more convincing if the new viewpoint enables new practical accomplishments or unifies seemingly unrelated phenomena. This section discusses what is gained by the matrix representation of discrete models and the matrix-search viewpoint, and culminates by suggesting ways to integrate separately-developed discovery systems.

There are several gains from the interpretation and notation introduced by this paper. First, they provide

¹The order of matrix multiplication in $AB = C$ has no special significance, since the theorem $(AB)^T = B^T A^T$ allows re-writing the former as $B^T A^T = C^T$.

a unifying framework that demonstrates a broad similarity of input/output representation, constraint representation, and elements of search. These similarities raise the question of whether a more general scheme could incorporate these systems as special cases.

A second benefit from the matrix viewpoint is that several constraints can be expressed and satisfied using explicit algebraic techniques, such as Gaussian elimination or linear programming. MECHEM and PAULI do use matrix manipulation to satisfy some constraints. MECHEM converts pathways to matrices in order to solve for the unknown substances by imposing the conservation law of reaction balance. MECHEM also uses matrix algebra to test whether a pathway can explain observed concentrations data. Finally, MECHEM and PAULI both use the simplex algorithm of linear programming to implement some constraints, and the simplex algorithm uses the matrix representation explicitly in the form of tableaux.

A third benefit is that matrix-based heuristics can guide us to find and address other scientific problems that resemble the current ones. One should look for problems that:

- Progressively enlarge classes of objects, structural elements, processes, or properties. Mention of simplicity or Occam's Razor in this connection is a favorable sign.
- Involve integral numbers of combinations of things.
- Involve linear constraints, e.g., conservation or additivity laws.

Examples of possible matches are Feynman diagrams in particle physics (in which the simplest diagrams are called "leading-edge" diagrams), models of ions, and models of atomic nuclei. Finally, the next section uses the matrix-search viewpoint to demonstrate how an integration of GELL-MANN with BR-3/PAULI could be carried out.

Integrating systems

The concept of search in matrix spaces can be applied to show how the task of GELL-MANN can be integrated smoothly with the task of BR-3/PAULI. GELL-MANN's search fills out the two matrices $\frac{\mathcal{S}}{\text{particles}} \mid \frac{\text{quarks}}{\text{quarks}}$ and

$\frac{\mathcal{P}}{\text{quarks}} \mid \frac{\text{properties}}{\text{properties}}$ subject to the constraint

$$\mathcal{S}_{s \times e} \times \mathcal{P}_{e \times p} = \mathcal{P}'_{s \times p}.$$

Given a reaction matrix $\frac{\mathcal{R}}{\text{reactions}} \mid \frac{\text{particles}}{\text{particles}}$, BR-3/PAULI's search fills out a property matrix $\frac{\mathcal{P}''}{\text{particles}} \mid \frac{\text{properties}}{\text{properties}}$ subject to the constraint

$$\begin{bmatrix} \mathcal{R}_{g \times s} \\ \mathcal{R}_{b \times s} \end{bmatrix} \times \mathcal{P}''_{s \times p} = \begin{bmatrix} \mathcal{O}_{g \times p} \\ \mathcal{Z}_{b \times p} \end{bmatrix}.$$

A combined system that carries out the tasks of GELL-MANN and BR-3/PAULI simultaneously would

fill out the \mathcal{R} , \mathcal{S} , and \mathcal{P} matrices and would need to satisfy at least the constraint

$$\begin{bmatrix} \mathcal{R}_{g \times s} \\ \mathcal{R}_{b \times s} \end{bmatrix} \times \mathcal{S}_{s \times e} \times \mathcal{P}_{e \times p} = \begin{bmatrix} \mathbf{0}_{g \times p} \\ \mathcal{Z}_{b \times p} \end{bmatrix}.$$

In the integrated system, *four* distinct matrix dimensions can be enlarged: the two from GELL-MANN, one from BR-3/PAULI, but also a fourth for the reaction dimension, since many new unseen "good" and "bad" reactions can be postulated just as GELL-MANN could postulate unseen particles and their property values. The combined system could accept exactly the \mathcal{R} input to BR-3/PAULI as before, and carry out a substantial theoretical effort by postulating quarks, properties, values, and unseen reactions all within a single system.

Conclusion

We have shown that the search carried out by a number of well-known systems that induce discrete models of hidden structure can be represented by sets of matrices on which constraints are placed. The typical dimensions of the matrices involve reactions (processes), substances (types of objects), and properties of substances. For example, DALTON finds structural models of chemical reactions and substances which can be described by a reaction matrix $\mathcal{R}_{r \times s}$, and a structure matrix $\mathcal{S}_{s \times e}$. $\mathcal{R}_{r \times s}$ describes reactions by the number of molecules of each substance in input and output, and $\mathcal{S}_{s \times e}$ describes the composition of each substance in terms of numbers of atoms of elements. Conservation of atoms is expressed by $\mathcal{R} \times \mathcal{S} = 0$.

The common matrix representation eases comparing these systems, reveals their underlying commonalities and sometimes shows how two systems (e.g., BR-3/PAULI and GELL-MANN) can be integrated into a single one. It also suggests how search algorithms designed for one system could be applied to others.

Hypothesizing new reactions, substances, or properties is accomplished by enlarging a matrix along one or more dimensions. The sizes of the matrices provide an (inverse) measure of a model's simplicity, so that generating small matrices first, then successively enlarging them as required, assures that simpler hypotheses are considered first, and that only as many hidden entities are introduced as are required to account for the data.

Representing model-building as search over a small matrix set does much to reduce the apparent diversity among the various systems, and shows that a few principles are fundamental to the organization and functioning of most of them. Hence, the representation is a significant advance toward a general theory of discrete model-building in scientific discovery.

Acknowledgments. RVP was supported partly by a W.M. Keck Foundation grant for advanced training in computational biology to the University of Pittsburgh, Carnegie Mellon, and the Pittsburgh Supercomputing Center. JMZ contributed to this paper while on sabbatical at Carnegie Mellon.

References

- Aris, R. and Mah, R.H.S. 1963. Independence of chemical reactions. *Ind. Eng. Chem. Fundam.* 2:90-94.
- Brualdi, Richard A. 1981. *Introductory Combinatorics*. North Holland, New York, NY. (Theorem on page 37).
- Fajtlowicz, Siemion 1988. On conjectures of Graffiti. *Discrete Mathematics* 72:113-118.
- Fischer, P. and Zytkow, Jan M. 1990. Discovering quarks and hidden structure. *Proc. of 5th International Symposium on Methodologies for Intelligent Systems*.
- Fischer, P. and Zytkow, Jan M. 1992. Incremental generation and exploration of hidden structure. *Proc. of the ICML-92 Workshop on Machine Discovery*.
- Klahr, D. and Dunbar, K. 1988. Dual space search during scientific reasoning. *Cognitive Science* 12:1-48.
- Kocabas, Sakir 1991. Conflict resolution as discovery in particle physics. *Machine Learning* 6(3):277-309.
- Langley, P.; Simon, H.A.; Bradshaw, G.L.; and Zytkow, J.M. 1987. *Scientific Discovery: Computational Explorations of the Creative Processes*. MIT Press.
- Lea, Glenn and Simon, Herbert A. 1974. Problem solving and rule induction: A unified view. In Gregg, Lee W., editor 1974, *Knowledge and Cognition*.
- Lenat, Douglas B. 1982. AM: Discovery in mathematics as heuristic search. In Davis, R. and Lenat, D.B., editors 1982, *Knowledge-based systems in artificial intelligence*.
- Lindsay, R.K.; Buchanan, B.G.; Feigenbaum, E.A.; and Lederberg, J. 1980. *Applications of Artificial Intelligence for Organic Chemistry: The Dendral Project*.
- Rose, D. and Langley, P. 1986. Chemical discovery as belief revision. *Machine Learning* 1(4):423-451.
- Rose, Donald 1989. Using domain knowledge to aid scientific theory revision. In *Proc. of the 6th International Workshop on Machine Learning*.
- Spirtes, P.; Glymour, C.; and Scheines, R. 1990. Causality from probability. In Tiles, J.E. et al., editors 1990, *Evolving Knowledge in Natural Science and Artificial Intelligence*. Pitman.
- Valdes-Perez, Raul E. Conjecturing hidden entities via simplicity and conservation laws: machine discovery in chemistry. *Artificial Intelligence*. In press.
- Valdes-Perez, Raul E. Discovery of conserved properties in particle physics: A comparison of two models. *Machine Learning*. Accepted for publication.
- Valdes-Perez, Raul E. 1991. A canonical representation of multistep reactions. *Journal of Chemical Information and Computer Sciences* 31(4):554-556.
- Valdes-Perez, Raul E. 1992. Theory-driven discovery of reaction pathways in the MECHEM system. In *Proc. of 10th National Conference on Artificial Intelligence*. 63-69.
- Zytkow, J.M. and Simon, Herbert A. 1986. A theory of historical discovery: the construction of componential models. *Machine Learning* 1(1):107-139.
- Zytkow, J.M. 1987. Combining many searches in the FAHRENHEIT discovery system. In *Proc. of the 4th International Workshop on Machine Learning*. 281-287.