

グラフ理論と化学構造表現

中山 堯・藤原 譲

1. はじめに

グラフとは、頂点の集まりとそれら頂点を結ぶ辺とからなる図形のことである。頂点の位置や辺の長さ・曲がり方はグラフにおいては問題とされない。グラフ理論はこのグラフを扱う理論であるが、歴史的には Euler(1736)によるケーニッヒスベルグの橋の問題、Kirchhoff(1847)による電気回路網の表現、あるいは Cayley(1857)による飽和炭化水素の異性体の数え上げといった実在の問題から抽出・再構成されて発展し、今日では計算機科学、OR、経済学、電子工学、化学、その他広範な応用領域をもっている。化学においては、異性体の数え上げに始まって、反応グラフ、合成設計、速度論、ドキュメンテーション、構造生成、化合物の命名法などに応用されている。グラフの直観的な定義にしたがえば、化学構造はそのままグラフと見做すことができる。すなわち原子を頂点、原子間の結合を辺に対応させると化学構造はグラフとして表現される。ただし、頂点はそれぞれ原子の種類に対応した標識をもち、辺も結合の種類に応じて識別される場合があるため特に化学グラフと呼ばれることがある。Cayley による飽和炭化水素の数え上げは、構造を木と呼ばれる一連のグラフとして抽象化することによって行

なわれたが、現在のグラフによる化学構造表現は膨大な量となった化合物(1981年現在で550万個)の計算機処理のためのデータ構造を与えることが主目的となっている。これまでに種々の化学構造表記法が提案されているが、それらは線形表記法とトポジカル表記法に大別される。

2. 化学構造の線形表記法

化学構造の表記法においてまず注意すべきことは、表現の uniqueness と ambiguity の問題である。これらは次のように定義される。

Uniqueness—ある表記法が unique であることの必要十分条件は、その表記法を任意の化学構造に適用して得られる表記が唯一つに限られることである。

Ambiguity—ある表記法が unambiguous であることの必要十分条件は、その表記法を任意の表記に適用して得られる化学構造(i. e., ある表記の解釈)が唯一つに限られることである。

この2つの特性にしたがえば、化学構造表記法は表1に示す4つのクラスに分類される。どのクラスの表記法を選ぶかは目的によって決められるべきである。たとえば、unique and unambiguous な表現は特定化合物の検索には適しているが化合物の包括的な処理には必ずしも適していない。慣用名は non-unique な表記法であるが、1つの化学構造について複数個の観点からの呼称が

許され処理の便宜を図るのに有利である、また、GREMAS コードのように、ambiguous であるがある種の部分構造検索には強力な表現手段となっているものもある。

表 1 化学構造表記法の分類

class	unique	unambiguous
1	yes	yes
2	yes	no
3	no	yes
4	no	no

表 2 Morgan法にいたるまでに発表されたトポロジカル表記法[18]

開発者	年
Wheland	1949
Mooers	1951
Rey & Kirsh	1957
Meyer & Wenke	1962
Spialter	1962
Dyson, Cossum, Lynch & Morgan	1963
Feldman, Holland & Jacobus	1963
Gould, Gasser & Rian	1965
Gluck	1963, 1965
Morgan	1965

3. 化学構造のトポロジカル表記法

線形表記法は特定の限定された目的のためには有用だが、反面汎用性に乏しく化学構造の大規模データベース構築に対する1つの障害となっている。化学構造データベースに対する最も汎用性の高い利用法は部分構造検索だが、一般的な部分構造検索に対して、線形表記法には明らかな限界が存在する。すなわち線形表記法では化学構造のあらゆる部分構造が一様に表現されず、いわば ad hoc な view を導入していることがアプリケーションに対する汎用性あるいは柔軟性と関連して、化学構造データの独立性を奪う原因となっている。

他方、connection table に代表されるトポロジカルな表記法は、その意味では一様な構造表現であり、すべての構造情報が包含される。原子間の結合関係のみに構造情報を限定すれば行列表現でトポロジカル表記法が実現される。これはグラフの隣接行列による表現に他ならない。

Rush によれば connection table を利用した表記法は Wheland(1949)から Morgan(1965)まで表2のように発表されている。CAS/Morgan法はグラフの頂点の番号付けを Morgan の connectivity value にもとづいて施した結果の connection table を用いるものだが、現在 CAS の Registry System に収録されている化学構造はこの方式を土台としている。しかしそれが表現法として必ずしも十分でないことは部分構造検索のために BASIC グループが fragment file を作成しており、それが CAS ONLINE に利用されていることから明らかである。すなわち connec-

tion table に包含されている構造情報を構造化する余地が大いに残されていると考えられる。そこでその構造化を与える道具としてグラフ理論がとみに注目を集めているわけだが、具体的には unique な表現(隣接行列)を得るためのグラフの頂点の番号付けの問題、グラフの分類のための不変数(特性量)の抽出あるいはグラフの特徴をベースにした構造表現法などが主たるトピックである。

3.1 グラフの頂点の番号付け

グラフの直観的な記述は最初に示したが、少し形式的に定義すれば次のようになる。

グラフとは、

- 1) 頂点の集合 $V = \{v_1, \dots, v_n\}$, および
- 2) V の直積 $V \times V = \{(v, w) | v \in V, w \in V\}$ の部分集合、すなわち辺の集合

$$E = \{e_1, \dots, e_m\}$$

からなる対 $G = (V, E)$ のことである。ここでは辺の重複は許さず、辺の向きも考えない。また (a, a) は辺と見做さない。 v が w に隣接するとは、 (v, w) が E の要素であることをいう。頂点の次数とは、それに隣接する頂点の個数である。隣接行列 A とは、 $n \times n$ の行列 $[a_{ij}]$ で頂点 i と j が隣接しているとき $a_{ij} = 1$, 隣接していないとき $a_{ij} = 0$ としたものである。図1にグラフとその隣接行列を示す。隣接行列 A はグラフ G の一般表現となっているが、行列自体は頂点の番号付けに

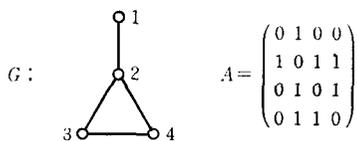


図 1 グラフとその隣接行列

応じて変化し、頂点数 n に対して $n!$ 通りの表現が可能である。したがって、任意の番号付けを許したのでは化学構造表現法に要請される uniqueness という条件に大幅に抵触することになる。頂点の番号付けを決めれば隣接行列も一意に決まるので隣接行列による化学構造表現の問題は、頂点の一意な番号付けの問題として提出されたことになる。

グラフの頂点の番号付けにおける基本方針は、頂点のトポロジカルな等価性(または非等価性)を識別し、それにしたがって頂点の集合を類別することである。これは各頂点に対してグラフのある不変量を計算することによって行なわれる。よく知られているものに、Morgan の connectivity value(cv) があり、下記の手続きで計算される。

S1. 各頂点にそれぞれの次数を割り当て、cv の初期値とし、異なる cv の値の数を k とする。

S2. 新しい cv を次式により計算する：

$$d_j' = \sum_{i=1}^n d_i a_{ij} \quad \text{or} \quad d' = dA$$

ここで d_j' は頂点 j の新しい cv, d_i は頂点 i のもとの cv, a_{ij} は隣接行列の (i, j) 要素である。これは、頂点 j の新しい cv がそれに隣接する頂点のもとの cv の和で与えられることを表わしている。

S3. 新しい cv の異なる値の数を k' とする。 $k' > k$ ならば $k \leftarrow k'$ として、S2 にもどって繰返し、 $k' \leq k$ ならば、もとの (i.e., k を与えた) cv(d) が各頂点の不変数を与える。

頂点の番号付けは、この後で最大の cv をもつ頂点に番号 1 を与え、次に頂点 1 に隣接する頂点に対してその中で cv の降順に 2, 3, ... と番号を

与え、それがすんだら 2, 3, ... の頂点についてその隣接頂点の番号付けを cv にもとづいて同様に繰り返すことによって遂行される。図 2 に cv の計算過程と番号付けの結果を示す。

ところが、cv の定義から類推できるように、正則グラフ(すべての頂点の次数が等しいグラフ)に対してはどの頂点の cv も同一の値となって識別力を発揮できないという限界が存在する。

頂点の一意な番号付けに対するグラフ理論的アプローチを標榜した一連の報告が Randić によってなされている。これは、隣接行列の各行を上から順に取り出し、それを左から右に並べてできる 2 進数が最小になるように番号付けをしようとするものである：隣接行列の各行を 2 進数 b_1, \dots, b_n とみなすと、これらを結合した 2 進数 b_1, \dots, b_n を最小にするためには b_1 が b_1, \dots, b_n の中で最小でなければならない。これを実現するためには、次数最小の頂点を 1 つ選びそれに最小番号 1 を与

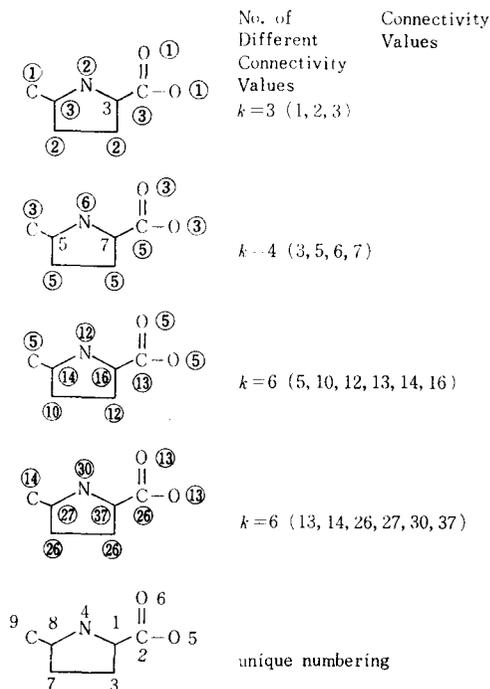


図 2 connectivity value の計算過程とその結果による頂点の番号付け

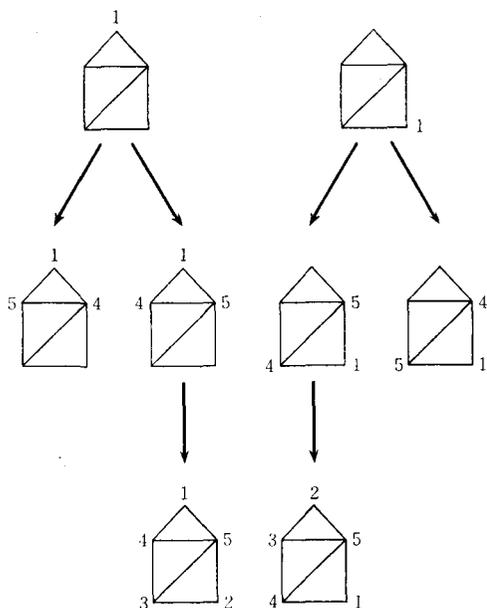


図 3 最小 2 進数法による頂点の番号付け

え、その隣接頂点に最大番号 n 以下を与える。次に有効次数（番号付けのすんだ頂点を無視した場合の次数）が最小の頂点を 1 つ選んで番号 2 を与え、その隣接頂点に残っている番号の中から最大のものを以降を与える。これを繰り返す。図 3 に例を示す。次数あるいは有効次数による頂点の分割は完全ではないから、上記手続きの各ステップごとに 1 つの番号に対して複数の頂点の候補が存在しうる。したがって、複数の番号付けの系列を生成しながら最小のものを残すようにしなければならない。正則グラフの場合、最も primitive な手続きによっても、可能な系列の数が $n \cdot d!$ を越えることはない。 $d \leq 4$ が仮定できれば、これは $n!$ の場合に比較して大きな改良となっている。

ところで、Morgan の connectivity value による方法も、Randić の最小 2 進数による方法も頂点の特性に着目して番号付けを一意に決定しようとするものであるが、頂点の等価性の評価がアルゴリズムに対して必ずしも厳密には反映されていない。グラフの頂点のトポロジカルな等価性は、頂点の集合上に置換群を導入することによって定式化することができる[5]。グラフ G の自己同型

写像 α とは、 G のそれ自身への同型写像、すなわち G の頂点集合上の置換で隣接性を保存するものである。 V の任意の頂点はその置換によってそれと同じ次数の他の頂点に移されるが、このときそれら 2 つの頂点は（トポロジカルに）等価であるという。 G の任意の 2 つの自己同型写像を引き続いて行なった写像は、また G の自己同型写像であるから G の自己同型写像の全体は G の頂点集合上の 1 つの置換群をなす。したがって、この置換群の元をすべて求めてその軌道を求めれば頂点集合を等価性によって識別することができる。対称性の高いグラフは置換群の位数が大きくなりすぎるので（完全グラフならば $n!$ ）、すべての置換を実際に求めるのは非現実的である。そこで、すべての置換を発生させることなく各頂点の類別と不変数を与えるために種々の heuristic な手続きを用意して一意な番号付けを得る方法が内野によって示された[19]。

3.2 グラフの分類

化学構造表現の問題に対するアプローチとしては、前節の主題であった構造の一義的表現（これは特定構造の同定を可能にする）を目標とする立場と、部分構造検索や構造生成などの応用を念頭においたグラフの特徴にもとづく構造の分類を目標とする立場がある。本節では、環構造を中心としてグラフの分類をめざす試みを紹介する。

SSSR (Smallest Set of Smallest Rings). 環構造の記述単位としてまず直観的に思い浮かぶのは、成分としての単環であろう。SSSR はそれをグラフ理論的に言いかえたものである。3 個以上の頂点からなるグラフの 2 重連結成分（ブロック）の頂点はすべて、ある閉路（単環）上の頂点となっている。このブロックは環構造（ring assembly）に対応している。グラフの頂点をすべて含んだ部分グラフが木（tree）になる場合、それを完全木（spanning tree）という。完全木にグラフの残りの辺の 1 つを付加すると 1 つの閉路ができる。ブロックの辺の数を m 、頂点の数を n とすれ

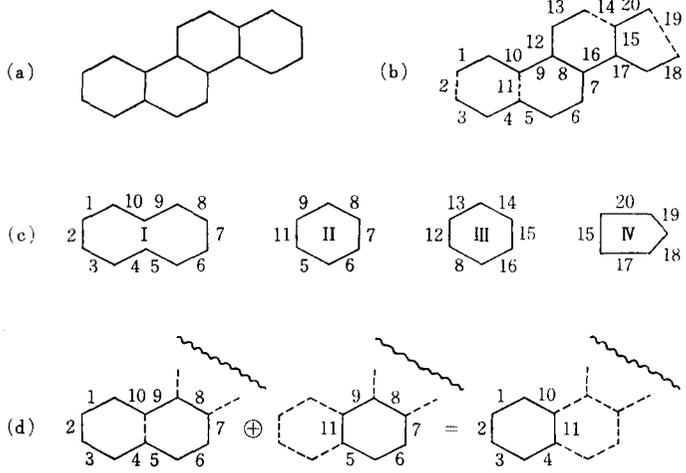


図 4 環構造の基底閉路

ば、完全木の辺の数は $(n-1)$ であるからこの方法により合計 $(m-n+1)$ 個の閉路が生成される。この閉路の集合はその完全木に関する閉路の基本系と呼ばれ、グラフに付随した次のようなベクトル空間の部分空間（閉路部分空間）の基底であることが示される。すなわち、グラフの辺に $1, \dots, m$ と番号を付けると任意の閉路 r に対して、

$$r_i = \begin{cases} 0 & \dots i \notin r \text{ の場合} \\ 1 & \dots i \in r \text{ の場合} \end{cases}$$

と定義することによって、ベクトル

$$r = (r_1, \dots, r_m)$$

を対応させることができる。閉路の和はベクトルとしての和、つまり各成分ごとの排他的論理和として定義する。図 4 は (a) で示されるグラフの完全木が (b) であり、(c) に示される 4 つの閉路 I, II, III, IV なる基底をもつことを表わしている。各基底のベクトル表現は下記のとおりである [20]。

- 閉路 I 11111111110000000000
- 閉路 II 00001111101000000000
- 閉路 III 00000001000111110000
- 閉路 IV 000000000000000101111

他の閉路はこれらの和で与えられる。たとえば、閉路 I と II の和は次の閉路を与える。

$$I \oplus II = 11111000001100000000$$

図 4 (d) にこの様子を示す。ところで Plotkin

によるブロックの SSSR の定義は次のとおりである [3]。

すべての閉路をそれらのサイズで順序づける。最小の閉路は常に SSSR の要素とする。それより大きい閉路は、すでに SSSR の要素となっている閉路と独立なものにかぎり SSSR の要素とする。これを繰り返して SSSR の要素数が $(m-n+1)$ 個となった時点で SSSR が完成する。

m と n はそれぞれブロックの辺と頂点の個数である。図 4 で与えられ

た基底は明らかに SSSR ではない（閉路 I の代りに閉路 I ⊕ II を採用すれば SSSR となる）。基底閉路を求めるアルゴリズムはいくつか報告されているが、完全木を利用するものが多い。グラフの頂点の番号付けが一意に決まれば、番号 1 から出発し各頂点にいたる道を隣接頂点の中の若い番号の中から先に選ぶという構成法によって完全木を一意に生成することができる。したがって、閉路部分空間の基底を一意に決定する手続きを構成できる。こうして、環構造を閉路部分空間の中で記述することが可能となり、グラフの構造的特徴による分類に対する見透しを与える。

トポロジカル・インデクス。SSSR が環構造をベクトル空間に置いて記述しようとする立場であるのに対し、グラフ的な特徴を全体として 1 つの特性量に反映させ、グラフを分類して予備的な screen として機能させようとする立場がある。トポロジカル・インデクスは後者の代表的なものの 1 つである。

隣接行列はグラフを一意に表現するものだが、その特性多項式の係数がグラフの構造（トポロジー）をよく反映することが知られている（ただし特性多項式とグラフが 1 対 1 に対応するわけではない）。

したがって、係数の組をグラフの特性量として

用いることができるが、さらに、それらの和も1つの特性量となりうる。このとき、行列式を展開することなく、グラフの辺に関する非隣接数 $p(G, k)$ によって係数が求まる場合がある[7]。グラフが木である場合の特性多項式 $P_G(x)$ は、

$$P_G(x) = \sum_{k=0}^m (-1)^k p(G, k) x^{n-2k}$$

と表わせる。 $p(G, k)$ はグラフ G において k 本の互いに隣り合わない辺を選ぶ組合せの数である (ただし、 $p(G, 0) = 1$ とする)。すると、トポロジカル・インデックス Z_G は、

$$Z_G = \sum_{k=0}^m p(G, k)$$

と定義される。特性多項式とグラフの対応は1対1でないから、 Z_G もグラフと1対1対応をしないが、グラフの集合に対する screen としては十分に機能しうる事が示されている。

3.3 トポロジカルな化学構造表現

化学構造のトポロジカルな表記法として CAS/Morgan 法が実用されているが、それが必ずしも満足すべき (完成度の高い) 表現法とみなされているわけではないことは前述のとおりである。そこで CAS/Morgan 法の発表以来も実にさまざまな表現法が提案されてきた。筆者らの提案である BCT 表現もこうしたものの1つであるが、部分構造検索によく適合するようにファイル構造を柔軟に編成している点と、それらのファイル・レコードがグラフ理論的にアルゴリズム的に構築される点が特長となっている[2, 10, 11]。

部分構造検索においては検索対象である部分グラフがコード (化学構造表現) から即座に取り出せる形になっていることが検索効率上望ましい。そこで化学構造集合の中で出現頻度の高い部分構造を選んで fragment set を構成し、それによって各化学構造を記述しようとする方法がある。これは言わば計算機によって処理可能な個数の

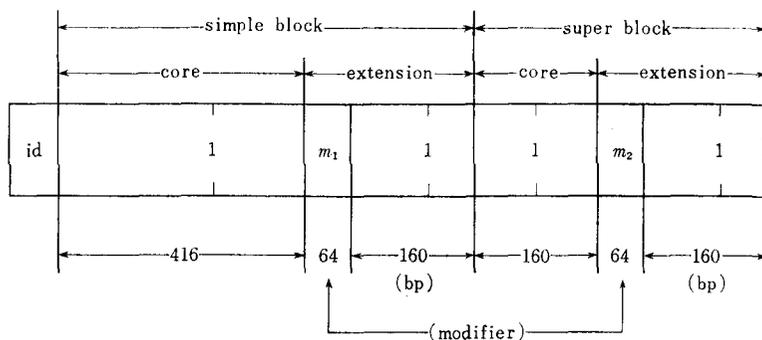


図5 fragment record format

fragment を基底に選んだベクトル空間に化学構造を写像しようとするものであるが、これは線形表記法に他ならない。しかしこの表現にはトポロジカルな情報の落ちがあり、検索における速度と多様性 (一般性) との trade-off が問題となる。

化学構造の BCT 表現は、トポロジカルな表現もすべて保存しながら fragment 情報を表現に対して明示的に付与しようとするものである。BCT 表現においては fragment はブロックによって組織的に与えられる。すなわちグラフ理論におけるブロックの定義 (その頂点を除くとグラフの連結成分が増えるような頂点は切断点と呼ばれ、切断点を含まない極大部分グラフをブロックと定義する) により、各化学構造から algorithmic に fragment が抽出される (その他に、ad hoc な fragment としていくつかの複合ブロックが用意されている)。ブロックまたは複合ブロックを用いた fragment record の形式を図5に示す。このファイルは BCF と呼ばれる。fragment record は、その化学構造内のブロック/複合ブロックの組成を表わすものであるが、core 部と extension 部に分けられる。core 部はそのビット位置がブロック/複合ブロックの識別子に直接対応し、extension 部は modifier とビット位置の組合せが識別子に対応する。core 部に出現頻度の高いものを配置することによって、fragment による screening 効果を高めることができる。

BCT(Block-Cutpoint Tree)は、切断点(Cut-

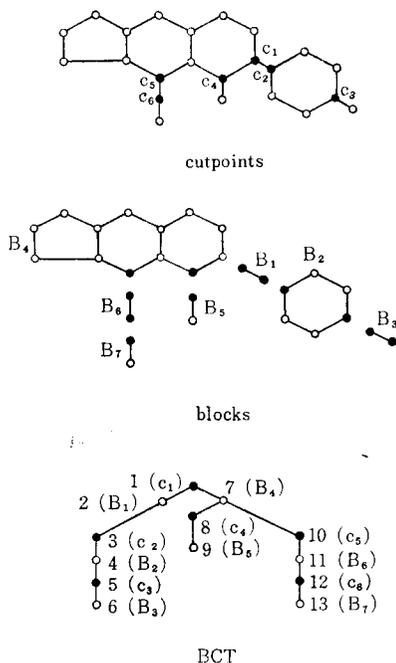


図 6 グラフの Block-Cutpoint Tree

point) とブロック (Block) を頂点集合とする 2 組グラフで、切断点 c がブロック B に含まれるとき (c, B) を辺と定義したものである。図 6 に BCT の例を示す。BCT は木であり、その標準形は容易に決定できる。

つまり、BCT はグラフに関するマクロな view を提供するものであり、ブロックがその記述子となっている。各ブロックの構造は隣接行列で与えられ、ブロック辞書 (BD) に登録される。ブロック辞書は各ブロックの SSSR にもとづいて構造化され、部分構造検索の便宜を図っている。この詳細については [12] を参照されたい。

BCT の組成ブロック間の厳密な結合関係は結合行列を生成することによって表現される (図 7)。これらは BCT と呼ばれるファイルに置かれる。化学構造の BCT 表現に関する全体のファイル構成を図 8 に示す。

BCT 表現がブロックを中間記述子とする階層的ネットワークとして実現されていることが看取されるであろう。また、BCT 表現法は特許における Markush formulas や法律等における自然

BCT code

1	2	3	4	5	6	7	8	9	10	11	12	13
13	5	4	3	2	1	7	2	1	4	3	2	1

	a	b	c	d	c	a	c																						
4	1	7	2	4	2	1	6	6	2	4	1	6	2	1	1	9	6	1	1	2	7	1	4	7	1	13	2	1	1

connectivity matrix

	2	4	6	7	9	11	13
2	0	1	0	2	0	0	0
4	1	0	6	0	0	0	0
6	0	1	0	0	0	0	0
7	11	0	0	0	6	1	0
9	0	0	0	1	0	0	0
11	0	0	0	1	0	0	2
13	0	0	0	0	0	1	0

図 7 グラフの BCT コードとブロック間の結合関係を表わす結合行列

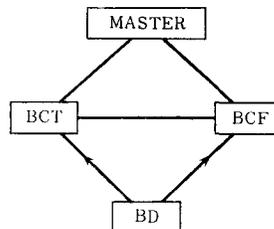


図 8 BCT 表現にもとづく化学構造データベースのファイル構成

言語による規定といった化学構造の包括的な表現に対しても自然に拡張することができるという利点をもっている。

4. おわりに

化学構造表現に対するグラフ理論の応用という観点から現在まで報告されている主要な表記法を概観してみたが、線形表現法のはらむ本質的な限界を補完するものとして登場したトポロジカル表現が、グラフ理論的アプローチによっていくつかの場面ではよい見透しを与えられたといえよう。

グラフの分類の問題では、分類の視点自体がアプリケーションごとに変化する可能性があり、あ

らゆる場合に融通無礙に適應しうる表現を得ることは困難であろうと予想されるが、むしろ種々の場合における heuristic な手法を蓄積して状況に応じた使い分けを可能とするようなシステムが望ましいともいえる。それら種々のアプリケーションの場においてもグラフ理論が強力な道具を与えることが期待される。

参 考 文 献

- [1] Ash, J. E. : Connection Tables and Their Role in a System. *Chemical Information Systems* (ed. J. E. Ash and Hyde), Ellis Horwood, 1975 (以下同)
- [2] Fujiwara, Y. and Nakayama, T. : A Graph Data Base for Storage of Chemical Structures Organized by the Block-Cutpoint Tree Technique. *Anal. Chem. Acta*, Vol. 133 (1981), 647-656
- [3] Gasteiger, J. and Jochum, C. : An Algorithm for the Perception of Synthetically Important Rings. *J. Chem. Inf. Comput. Sci.*, Vol. 19, No. 1 (1979), 43-48
- [4] Golender, V. E. : Graph Potentials Method and Its Application for Chemical Information Processing. *J. Chem. Inf. Comput. Sci.*, Vol. 21, No. 4 (1981), 196-204
- [5] Harary, F. : *Graph Theory*, Addison-Wesley Publishing Co., Reading, Massachusetts, 1969
- [6] 平山健三 : 化学構造の線形表記法. 化学総説 No. 18 (1978), 9-26
- [7] 細矢治夫 : 化学構造表現の数式的手法. 化学総説 No. 18 (1978), 27-46
- [8] リウ, C. L. : 組合せ数学入門Ⅱ (伊理正夫・伊理由美共訳). 共立出版株式会社, 1972
- [9] Morgan, H. L. : The Generation of a Unique Machine Description for Chemical Structures-A Technique Developed at Chemical Abstracts Service. *J. Chem. Doc.*, Vol. 5 (1965), 107-113
- [10] Nakayama, T. and Fujiwara, Y. : BCT Representation of Chemical Structures. *J. Chem. Inf. Comput. Sci.*, Vol. 20, No. 1 (1980), 23-28
- [11] Nakayama, T. and Fujiwara, Y. : Structure Generation on the Basis of BCT Representation of Chemical Structures. *J. Chem. Inf. Comput. Sci.*, Vol. 21, No. 4 (1981), 218-223
- [12] 中山堯, 藤原譲 : BCT 表現にもとづく化学構造データベースと部分構造検索. 第18回情報科学技術研究会発表論文集, 1981, 113-121
- [13] Randić, M. : On the Recognition of Identical Graphs Representing Molecular Topology. *J. Chem. Phys.*, Vol. 60, No. 10 (1974), 3920-3928
- [14] Randić, M. : On Unique Numbering of Atoms and Unique Codes for Molecular Graphs. *J. Chem. Inf. Comput. Sci.*, Vol. 15, No. 2 (1975), 105-108
- [15] Randić, M. : On Canonical Numbering of Atoms in a Molecular and Graph Isomorphism. *J. Chem. Inf. Comput. Sci.*, Vol. 17, No. 3 (1977), 171-180
- [16] Randić, M. : Graph Theoretical Approach to Recognition of Structural Similarity in Molecules. *J. Chem. Inf. Comput. Sci.*, Vol. 19, No. 1 (1979), 31-37
- [17] Randić, M., Brissey, G. M., and Wilkins, C. : Computer Perception of Topological Symmetry via Canonical Numbering of Atoms. *J. Chem. Inf. Comput. Sci.*, Vol. 21, No. 1 (1981), 52-59
- [18] Rush, J. E. : Status of Notation and Topological Systems and Potential Future Trends. *J. Chem. Inf. Comput. Sci.*, Vol. 16, No. 4 (1976), 202-210
- [19] 内野正弘 : 化学構造式の処理. 岩波講座情報科学—23 数と式と文の処理 (伊理正夫編) 岩波書店, 1981, 59-100
- [20] Wipke, W. T. : Use of Ring Assemblies in a Ring Perception Algorithm. *J. Chem. Inf. Comput. Sci.*, Vol. 15, No. 3 (1975), 140-147