

# 発見的手法による知識ベースの検索

志村正道

(東京工業大学 工学部 情報工学科)

## §1. まえがき

近年、知識を利用するエキスパートシステムなどが種々開発され、その実用化も盛んになりつつある。本論文で報告するシステムは植物の名前を同定するための検索システムである。一般に、植物図鑑においては植物名からその植物の特徴を調べるのは容易であるが、その逆は極めて困難である。ここでは、植物の特徴からその植物名を推定する検索システムについて考える。

検索システムはいくつかのキーワードによって、その検索を行なうのが基本的である。すなわち、木構造状の検索項目を木の探索によって検索する方法が基本的である。しかし、実際には意味のない冗長なキーワードがあるとか、検索項目が多い場合などには、このような木の探索方法は効率がよくない。そのため、探索において発見的情報を用いた手法が種々報告されている。ところが、そのようなシステムでは知識ベースとしてテーブルを用いた場合、実際にはキーワードが欠除していたり、あるいはキーワードを変更したりすることが必ずしも容易ではないという欠点が生ずる。したがって、知識をプロダクショナルールの形で表現することを考え、実際に植物を対象とする検索システムを試作した。本論文ではこのシステムについて述べ、またその実験結果について報告する。

## §2. 知識の表現

知識の追加・修正・消去・確認が容易に行なえるような形式を考える。樹木状に知識を格納する方法が一般的であるが、次のような欠点をもっている。すなわち、

- (i) キーワードが欠損している場合探索をすすめることが困難である。
- (ii) データ構造が効率的でなくなる。
- (iii) 知識あるいはデータの拡張が困難である。
- (iv) 多くのデータを扱い難い。

以上の理由によってここではプロダクショナルール(PR)を用いたプロダクショナルシステムを用いることを考える。

いま、植物名を  $P_i$ 、質問事項を  $Q_j$ 、その回答を  $A_{jk}$  とするとき、 $Q_j$  に関する  $P_i$  の知識を

$$P_i = Q_j (A_{jk})$$

と表わす。また複数の知識を表現するには

$$Q_j (A_{jk_1}) \cup Q_j (A_{jk_2}) \cup \dots \cup Q_j (A_{jk_m})$$

ということではなく

$$Q_j (A_{jk_1} \cup A_{jk_2} \cup \dots \cup A_{jk_m})$$

という知識表現を用いる。

また、葉の長さなどの数値をとり扱う場合には min 値と max 値のみを知識として  
もつことにし、 $Q'_j (\text{min-max})$  という形で記憶する。

したがって、本システムで取り扱う知識は次のような二つの形である。

$$P = Q (A)$$

$$P = Q' (\text{min-max})$$

このような知識は実際に内部コードに変換されて格納されるが、本システムに  
おいては図 1 に示すような表現形式を採用している。このような表現形式の特長  
として次のような点がある。例えば、まきひげの  
長さに関する質問事項はまきひげをもつ植物の検索  
索には有効であるが、まきひげをもつ植物は全体  
からみればごく一部であって、それ以外の多くの  
植物に対しては無意味な質問となってしまう。し  
たがって、関係ある植物のみに限ってルールを設  
定しておけばよく、二次元配列を用いる形式に比  
べて効率的な記憶領域の使用ができる。

$P_1$	$Q_1$	$A_{11}$	*
$P_1$	$Q_2$	$A_{12}$	*
$P_1$	$Q_3$	Min	Max
$P_2$	$Q_2$	$A_{21}$	*
⋮	⋮	⋮	⋮
$P_i$	$Q_j$	$A_{ij}$	*
⋮	⋮	⋮	⋮

図 1

### §3. 知識の削除, 追加, 修正, 合成

知識ベースを考える場合に重要なことの一つは知識の削除, 追加, 修正, 合成  
であり, これらの機能が容易に実行できるようなシステムが望ましい。このよう  
なシステムにおいては, とくに知識の量を増やして, 知識ベースを拡張していく  
ということは, 本質的に要求される機能なのである。

上に述べたような機能を考えるには, 実際に次のような問題を解決しておかね  
ばならない。すなわち,

- (i)  $P=Q(A_1)$  なる知識がすでにあるか否を見出さなければならない。
- (ii)  $P=Q(A_1)$  に対し  $P=Q(A_2)$  なるデータが入ったとき, 新しい知識として  
 $P=Q(A_1)$  を捨てて  $P=Q(A_2)$  のみにするか, あるいは  $P=Q(A_1 \cup A_2)$  の  
ように両方も正しいものとすべきかを決定しなければならない。本シス  
テムでは問合せが出来るようになってくる。
- (iii) 新しい植物名が入ったとき, 過去の植物名の領域に新しい分割が入ってく  
る可能性があり, 誤って検索されるとかあるいは検索効率が下がる副作用  
が生じないようにしなければならない。

#### § 4. 知識データの検索

本システムの検索動作について述べる。対象の検索としては、検索される対象のもつ特徴をすべてシステムに入れ、これらの情報によって一度に行なう方法が考えられる。しかし、検索される対象によつてそれぞれ特徴が異なるため効率が必ずしも良くはない。それゆえ、ここで述べるシステムにおいては、検索するのにもっとも適当と考えられる特徴を求め、これについてシステムが検索者に質問を發する。検索者はこの質問に回答するという形をとりながら、逐次特徴に関する情報をシステムに取り入れる。システムは入力された情報すなわち特徴にしたがって検索を行なうが、この動作は探索対象が一意に定まるまで続けられる。

さて、プロダクションシステムにおける認識-作動サイクルは次の3つのステップから構成されている。

- (1) 知識データ（プロダクション）の中から起動条件の満たされているものの集合をつくる。
- (2) その集合の中から、何らかの方法で起動すべきプロダクションを選択する。
- (3) 選択されたプロダクションに従つて、実行を行なう。

本システムにおいては、プロダクションルールの起動は質問の發生ということに相当するが、このとき、質問事項の選択は次のような定め方が考えられる。

- (1) 固定した順序で質問をする。
- (2) 最適な順序で質問をする。
- (3) 最も新しく使われた順序で質問をする。

このうち(2)の最適な順序とは、その質問事項がなるべく多くの植物に関連しており、しかもなるべく多くの分類が可能な質問から選んでいった順序のことである。また、質問回答の結果より、順次条件に適合するプロダクションルールを起動させていき、それに合わない植物を棄却していく。

このような手順をスーパーステップとよぶが、このスーパーステップを実行した場合、検索は次のような状態に至り、たときに終了する。

- (1) 条件にあう植物が一つになった状態
- (2) 質問をすべて出しつくした状態
- (3) 候補植物がなくなった状態

すなわち、(1)の場合は検索を行なつていた対象が決定したわけであるから成功したことになる。(2)の場合は検索すべき入力情報が不足あるいは不備であった場合である。(3)は誤情報あるいは知識の不足のために検索対象が見つからなかった場合である。

#### § 5. 知識ベースのデータ構造

検索システムにおいて用いられる種々のデータのデータ構造を以下に示す。

##### (i) JOHO

JOHOは内部コードに変換された知識データを格納し、6次元(6 Wds)の配列である。各ワードの内容は 1. 植物名 P, 2. 質問事項 Q, 3. 回答 A, 4. 最小値, 最大値, 5. Pのポインタ, 6. Qのポインタである。これらの関係を示

すと図2のようになる。

1	$P_1$	$Q_1$	$A_{1k}$	*	$P_1$ のポインタ	$Q_1$ のポインタ
2	$P_2$	$Q_2$	min	max	$P_2$ のポインタ	$Q_2$ のポインタ
	⋮	⋮	⋮	⋮	⋮	⋮
N	$P_N$	$Q_N$	$A_{Nk}$	*	$N^*$	$N^*$
	⋮	⋮	⋮	⋮	⋮	⋮
$N'$	$P_{N'}$	$Q_{N'}$		*		
$N''$	$P_{N''}$	$Q_{N''}$	$A_{N''k}$	*	$N''$	0

図 2

(ii) MOZIRT

入力された文字列は内部コードに変換されるが、MOZIRTは内部コードと文字列を対応させるために、これらの関係を格納しておく配列である。このとき、図3に示すように、1ワードに4文字があてられ、各文字列は先頭にワード長を記入した後に格納される。

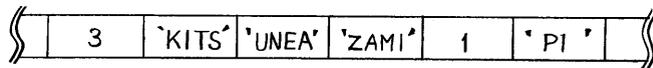


図 3

(iii) NLIB

植物名に関する情報を格納しておく4 Wdsの配列である。NLIBの4 Wdsの内容は 1. MOZIRTに格納されている植物名の文字列を示すポインタ, 2. 質問事項を示すポインタ, 3, 4. 配列JOHOで最初および最後に現れる番地のポインタである。

(iv) FLIB

質問事項に関する情報を格納しておく5 Wdsの配列である。FLIBの5 Wdsの内容は1~4. NLIBと同じ, 5. 質問事項の種類(数値あるいは分岐)である。

(v) XLIB

回答に関する情報を格納しておく3 Wds配列である。XLIBの3 Wdsの内容は 1. MOZIRTのポインタ 2. 回答項目のポインタ 3. 質問事項のポインタである。

以上、5つの配列によって知識ベースが構成されている。すなわち、JOHOはデータそのものを、NLIB、FLIB、XLIBは植物名P、質問事項Q、回答項目Aの索引を、またMOZIRTはこれらに登録された文字列を表わす辞書である。

## § 6. 検索植物名の決定

検索システムでは質問に対する回答結果から候補植物名を決定していくが、この過程について述べる。まず、最初においてはすべての植物が候補となっているが、各候補には荷重が割り当てられ、この荷重の最大なものが最終候補となる。すなわち、検索される植物の特徴が質問事項によって候補植物と一致する場合には荷重を増加し、一致しない場合には大巾に減少させる。また不明の場合には少し減少させる。このような荷重操作により、逐次候補植物の個数は減少していき、前に述べたような終了の状態になったとき終了する。

システムが入力情報をとり入れるために発生させる質問はなるべく最小の質問回数によって該当植物名が決定されるように選ばなければならない。このためには、その特徴によって分類される対象が多くなるようなものを選ぶなければならない。例えば、質問の回答結果に応じて、候補対象が2つのクラス、すなわち、その質問に対応する特徴をもつものともたないものとは分類されるものとしよう。対象の数を  $N$  とし、これを一つの特徴すなわち一回の質問により (a) 1 と  $N-1$  個 (b)  $N/2$  と  $N/2$  個に分類される場合を考えてみる。

(a) の場合は平均

$$\frac{N+1}{2} - \frac{1}{N} \approx \frac{N}{2}$$

(b) の場合は平均

$$\log_2 N$$

回の質問により分類されることとなる。したがって、(b) のようになると多数個づつの分類がなされるような特徴について質問することが望ましい。あるいは最も多くの候補を否定するような質問であることが望ましいとも考えられるので、その質問の有効性を、否定される候補数の期待値で表わすことにしよう。なお、質問の回答として  $m$  個のうちただ一つが入力されるものとし、二つ以上の回答はないものとする。

まず、選択式質問についての期待値は次のように求められる。質問  $Q$  に対して回答  $A_i$  が得られる確率を  $P_i$ 、該当回答がない確率を  $P_0$ 、またこのような質問に対して候補対象  $N$  のうち否定される対象の数をそれぞれ  $\bar{n}_i$ 、 $\bar{n}_0$  とする。このとき、期待値は

$$E = \sum_{i=1}^m P_i \bar{n}_i + P_0 \bar{n}_0$$

ただし、 $m$  は選択分岐数である。ここで、質問  $Q$  に関するデータをもっている対象の個数を  $n$ 、回答  $A_i$  をもつ対象の個数を  $n_i$  とすると  $\bar{n}_i = n - n_i$  となる。

したがって

$$P_i = \frac{n_i}{N}$$

$$P_0 = \frac{N-n}{N}$$

が得られ、期待値は

$$E = \sum_{i=1}^m \frac{n_i}{N} (n - n_i) + \frac{N-n}{N} n$$

となる。

次に数値式回答の期待値を求める。この場合回答は数値となるから、その値を  $v_i$  とし、また、 $\min \leq v_i \leq \max$  を満足する対象の数を  $n_i$  とする。ただし、 $v_i$  は離散値で、 $\min$  から  $\max$  までの値を等確率でとるものとする。このとき  $v_i$  をとる確率  $P_i$  は

$$P_i = \frac{1}{N} \sum_k \frac{1}{\max_k - \min_k + 1} \quad (\min_k \leq v_i \leq \max_k)$$

となり、期待値は

$$E = \sum_{v_i} \frac{n - n_i}{N} \left\{ \sum_k \frac{1}{\max_k - \min_k + 1} \right\} + \frac{N - n}{N} n$$

となる。

以下、この期待値の例を述べる。

例1.  $(N-m)$  種の対象がその質問に無関係で、関係する対象はそれぞれ異なる回答をもっている場合

$$E = \sum_{j=1}^m \frac{m-1}{N} + \frac{N-m}{N} \cdot m = \frac{m}{N} (N-1)$$

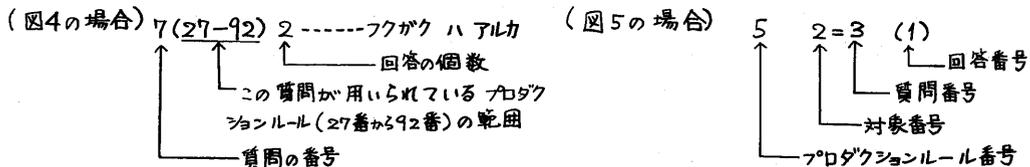
例2.  $l$  種の対象がその質問に無関係で、残り  $N-l$  種が均等に  $m$  個の回答をもっている場合

$$E = \frac{(N-l)^2}{N} \left(1 - \frac{1}{m}\right) + \frac{l}{N} (N-l)$$

実際のシステムにおける実験結果では、平均 7.5 個の質問事項によって候補植物名が決定されたのに対し、固定した一定順序の質問方法によれば 20~40 の質問を行なわねばならなかった。この結果から、上に述べた質問事項の選択はかなり有効であるといえてよからう。

## §7. システムの概要

本システムの知識ベースに格納されている植物は 58 種である。また本システムで準備されている質問項目は図4に示すような 28 種である。図5はプロダクションルールの一例であって 467 個準備されている。図4および5における数字は下のような意味をもっている。



実際に検索を行なった例を図6に示しておく。

本システムにおける付随的ではあるが重要な機能に次のようなものがある。

- (i) 検索中任意の時点で、候補植物名とその荷重を表示することができる。
- (ii) 候補植物の決定に際し、使用したプロダクションルールを表示することができる。
- (iii) 回答が不明の場合にも不明として入力できる。

1( 1-457) 2	-----	マキヒケ <sup>カ</sup> ハ アルカ アル ナイ
2( 2-458) 2	-----	タクヨウ ハ アルカ アル ナイ
3( 3-462) 4	-----	タンヨウ カ フクヨウ カ タンヨウ 5シユツ フクヨウ 3シユツ フクヨウ ウシ <sup>カ</sup> ヨウ フクヨウ
4( 12-463) 4	-----	ハ ノ ツキカタ ハ コンセイ タイセイ ゴ <sup>カ</sup> セイ リンセイ
5( 15-419) 2	-----	ハナ ノ タイシヨウメン ハ ホウシヤカ ウカ
6( 25-414) 3	-----	ハ ノ ヘリ ハ キレコム ノキ <sup>カ</sup> リハ <sup>カ</sup> セ <sup>カ</sup> ンハ <sup>カ</sup> ン
7( 27-92 ) 2	-----	フカ <sup>カ</sup> ク <sup>カ</sup> ハ アルカ アル ナイ
8( 28-265) 4	-----	オシ <sup>カ</sup> ノ カス <sup>カ</sup> ハ タスウ 10 ホ <sup>カ</sup> ン 6 ホ <sup>カ</sup> ン 2 ホ <sup>カ</sup> ン
9( 30-459) 2	-----	ヨウシユウ ハ アルカ アル ナイ
10( 31-460) 2	-----	クキ ハ チウクウ カ ハイ イイ
11( 32-118) 2	-----	クキ ニ ケ ハ アルカ アル ナイ
12( 33-417) 3	-----	カ <sup>カ</sup> ク <sup>カ</sup> ヘン ノ カス <sup>カ</sup> ハ 6 マイ 4 マイ ナイ

図4. 質問の例

1	1 = 1 ( 2)
2	スミレ = マキヒケ <sup>カ</sup> ハ アルカ? (ナイ)
3	1 = 2 ( 1)
4	スミレ = タクヨウ ハ アルカ? (アル)
5	1 = 3 ( 1)
6	スミレ = タンヨウ カ フクヨウ カ? (タンヨウ)
7	2 = 1 ( 2)
8	PIムク <sup>ラ</sup> = マキヒケ <sup>カ</sup> ハ アルカ? (ナイ)
9	2 = 3 ( 1)
10	PIムク <sup>ラ</sup> = タンヨウ カ フクヨウ カ? (タンヨウ)
11	3 = 1 ( 2)
12	アメリカフロウ = マキヒケ <sup>カ</sup> ハ アルカ? (ナイ)
13	3 = 2 ( 1)
14	アメリカフロウ = タンヨウ ハ アルカ? (アル)
15	3 = 3 ( 1)
16	アメリカフロウ = タンヨウ カ フクヨウ カ? (タンヨウ)
17	4 = 1 ( 2)
18	セ <sup>カ</sup> ニハ <sup>カ</sup> ア <sup>カ</sup> イ = マキヒケ <sup>カ</sup> ハ アルカ? (ナイ)
19	4 = 2 ( 1)
20	セ <sup>カ</sup> ニハ <sup>カ</sup> ア <sup>カ</sup> イ = タンヨウ ハ アルカ? (アル)
21	4 = 3 ( 1)
22	セ <sup>カ</sup> ニハ <sup>カ</sup> ア <sup>カ</sup> イ = タンヨウ カ フクヨウ カ? (タンヨウ)
	1 = 4 ( 1)
	スミレ = ハ ノ ツキカタ ハ? (コンセイ)
	3 = 4 ( 2)
	アメリカフロウ = ハ ノ ツキカタ ハ? (タイセイ)
	4 = 4 ( 3)
	セ <sup>カ</sup> ニハ <sup>カ</sup> ア <sup>カ</sup> イ = ハ ノ ツキカタ ハ? (ゴ <sup>カ</sup> セイ)
	4 = 5 ( 1)
	セ <sup>カ</sup> ニハ <sup>カ</sup> ア <sup>カ</sup> イ = ハナ ノ タイシヨウメン ハ? (ホウシヤカ)
	5 = 1 ( 2)
	タチスミレ = マキヒケ <sup>カ</sup> ハ アルカ? (ナイ)
	5 = 2 ( 1)
	タチスミレ = タクヨウ ハ アルカ? (アル)
	5 = 3 ( 1)
	タチスミレ = タンヨウ カ フクヨウ カ? (タンヨウ)
	5 = 4 ( 3)
	タチスミレ = ハ ノ ツキカタ ハ? (ゴ <sup>カ</sup> セイ)
	5 = 5 ( 2)
	タチスミレ = ハナ ノ タイシヨウメン ハ? (サユウカ)
	6 = 1 ( 2)
	ハルタテ <sup>カ</sup> = マキヒケ <sup>カ</sup> ハ アルカ? (ナイ)
	6 = 2 ( 1)
	ハルタテ <sup>カ</sup> = タクヨウ ハ アルカ? (アル)

図5. プロダクショナルールの例

## §8. あとがき

植物の色・形状などの特徴に基づいて、その植物名を決定する植物検索システムを試作した。このシステムはとくに知識の追加・訂正などを容易にするために、プロダクショナルールの形で与えてある。しかし、検索する植物の種類が多く、したがって検索するための特徴が多い場合にはプロダクショナルールによる手法では記憶容量が大きくなる。しかし、その平面欠損データや特種な特徴などがある場合とかあるいは知識の追加・訂正などの機能を有していることを考慮すれば、やはりテーブル方式よりプロダクショナルールによる知識の表現形式の方が優れていると結論できよう。試作モデルでは約60種の植物についてであるが、今後さらに種類を増加するとともにまたそのデータ構造についても重要な課題として考察が必要であろう。

最後に、プログラミングおよび実験については、松田至弘(YHP)、森光彦(ヤマハ)および石塚昭夫(院生)の三君の大きな協力があったことを付記しておく。

```

KENSAKU O HAJIMEMASU.
32.83 26 54          TANYOU KA FUKUYOU KA?
NI TSUITE TSUGINOUCHIKARA ERANDE KUDASAI.
1-->          TANYOU
2-->          5SHUTSU FUKUYOU
3-->          3SHUTSU FUKUYOU
4-->          UJOU FUKUYOU
03350 ?
4
9.29 22 21          YUUSHOU WA ARIKA?
NI TSUITE TSUGINOUCHIKARA ERANDE KUDASAI.
1-->          ARU
2-->          NAI
03350 ?
2
6.44 19 16          IAKUYOU WA ARIKA?
NI TSUITE TSUGINOUCHIKARA ERANDE KUDASAI.
1-->          ARU
2-->          NAI
03350 ?
2
4.75 17 12          ZENTAI NO IRU WA?
NI TSUITE TSUGINOUCHIKARA ERANDE KUDASAI.
1-->          RYOKUSHOKU
2-->          TANSHISHOKU
3-->          TANOOSHOKU
03350 ?
1
4.20 16 10          HANA NO TAISHOUMEN WA?
NI TSUITE TSUGINOUCHIKARA ERANDE KUDASAI.
1-->          HOUSHAKA
2-->          SAYUUKA
03350 ?
2
3.14 8 7            K.KI WA CHUUKU KA?
NI TSUITE TSUGINOUCHIKARA ERANDE KUDASAI.
1-->          YES
2-->          NO
03350 ?
2
2.00 4 5            TOUJOUKA KA?
NI TSUITE TSUGINOUCHIKARA ERANDE KUDASAI.
1-->          YES
2-->          NO
03350 ?
EL
KANOUSEI NO ARU SHOKUBU WA...
PLANT              KOUTEI KAISUU
1                  ITACHINOSASAGE    0
31                 KITSUNEAZAMI     5
37                 KIKEMAN         6
39                 MURASAKIKEMAN    6
40                 SERIBAHIENSOU    6
* IJOH 5 SHU *
2.00 4 5            TOUJOUKA KA?
NI TSUITE TSUGINOUCHIKARA ERANDE KUDASAI.
1-->          YES
2-->          ~NO
03350 ?
2

```

図 6. 検索の例

### 参考文献

1. 松田圭弘, "対話形式の植物検索システムの作成", 東京工大 情報工学科 卒業論文 1980.3
2. 森 光彦, "知識ベースを利用した情報検索システムの作成", 同上, 1981.3
3. 石塚昭夫, "知識ベースと対話形式検索手法の研究", 同上, 1982.3
4. 西村愷彦, "FORTRAN データベース", bit 増刊号
5. 石戸 忠, "実践的植物検索小図鑑1-2", 講談社, 1979
6. 同上, "学習コンピュータ植物図鑑", 文理書院
7. 長沢雅男, "情報検索入門", 森北出版, 1971