



# 法則式発見の機械化

元 田 浩\*

## Computer Assisted Discovery of Law Equations

Key Words : Scientific Discovery, Induction, Data Mining, Knowledge Acquisition

### 1. はじめに

人工知能の一分野に科学的知識の発見という研究分野がある。古来、物理学者が観測データを深く分析し裏に潜む真理を見出して来たように、人間には手におえない多量の実験データ(数値)から、データ間の関係を支配する第一原理の法則を計算機を用いて機械的に発見することは可能であろうか。ここで、あえて第一原理と言う言葉を用いたのは、単にデータにフィットする実験式ではなく対象に生起する現象自体に客観的に内在する法則の発見を意図している。

当研究室では知識発見を中心課題として、法則式発見の問題にも取り組んでおり、測定量の尺度情報の制約しか用いないで、測定データ間に許される一般的な関係式を導出した。また、これに基づき許容関係式の中から測定データを説明できる具体的な関係式を実験的に求めるアルゴリズム開発し、理想気体の方程式やブラックの比熱則など既知の物理法則を再発見し、熱伝達現象や電子回路など数10変数で表現されるかなり複雑な現象の法則同定にも成功した。

### 2. 尺度に関する数学的許容性

#### 2.1 数量の関係制約

我々が扱う数値データとして扱う数量には単位が

ある。単位が分かれば物理学でよく知られている次元解析の手法を用いて、例えば、圧力と温度を加えるというような無意味な演算は排除でき、単位に関し整合性のあるものだけに絞る事ができる。しかし、単位以前に測定量には尺度と言う重要な概念がある。個々の数量は表1に示す尺度を持っており、それらの単位変換に関する群構造制約が複数の数量間に数学的に許容される関係を著しく規定する。例えば2数量間の関係に限定すれば表2に示す関係のみが許される<sup>3)</sup>。このうち、比例尺度に関しては複数の数量の間に成立する著名な定理(Product Theorem)が知られている<sup>1)</sup>。ここで注意すべきは単位が分からなくても尺度は実験的に構成できるし、推定しやすいことである。このことは本手法の適用範囲が物理現象のみに限定されるものでなく、社会、心理現象など他の分野にまで適用可能であることを示唆している。

表1 各尺度の特徴

尺度	基礎的 経験操作	数学的群構造 一意性
間隔	間隔や差の 等値性の決定	一般的線形群 $x' = kx + c$
比例	比の等値 性の決定	相似群 $x' = kx$
絶対	値の等値 性の決定	恒等群 $x' = x$

表2 2数量間の数学的許容関係式

No.	尺度の種類		許容関係
	独立量	従属量	
1	比例	比例	$u(x) = \alpha x^\beta$
2.1	比例	間隔	$u(x) = \alpha x^\beta + \delta$
2.2			$u(x) = \alpha \log x + \beta$
3	間隔	比例	不可能
4	間隔	間隔	$u(x) = \alpha x + \beta$



\* Hiroshi MOTODA  
1943年3月24日生  
1967年東京大学大学院工学系研究科  
原子力工学専攻、修士課程修了  
現在、大阪大学産業科学研究所・知  
能システム科学研究部門・高次推論  
方式研究分野、教授、工学博士、人  
工知能  
TEL 06-6879-8540  
FAX 06-6879-8544  
E-Mail motoda@sanken.osaka-  
u.ac.jp

**定理1 (Product Theorem)** 比例尺度の数量  $x, y, \dots$  で表される数量  $f$  の関係式は  $f=Cx^a y^b z^c \dots$  の形式しか有り得ない。ここで  $C, a, b, c, \dots$  は定数である。

これらの知見を基に、我々は間隔、比例、絶対の3種の異なる尺度が混在する多数量間に許容される関係を表す以下の定理を導出した。

**定理2 (拡張 Product Theorem)**  $R$  を比例尺度の数量の集合、 $I$  を間隔尺度の数量の集合とすると、各  $x_i \in R \cup I$  の間に許される関係は、以下の2式以外には有り得ない\*1。

$$\Pi = \left( \prod_{x_i \in R} |x_i|^{a_i} \right) \left( \prod_{I_k \in C} \left( \sum_{x_j \in I_k} b_{kj} |x_j| + c_k \right)^{a_k} \right)$$

$$\Pi = \sum_{x_i \in R} a_i \log |x_i| + \sum_{I_k \in C_g} a_k \log \left( \sum_{x_j \in I_k} b_{kj} |x_j| + c_k \right) + \sum_{x_\ell \in I_g} b_{g\ell} |x_\ell| + c_g$$

## 2.2 法則式の構造

法則式の構造を決める重要な制約に、次元解析において知られているもう一つの著名な定理(BuckinghamのII-theorem<sup>2)</sup>)がある。この定理も比例尺度のみに関するものであり、我々の手で間隔、比例、絶対の3種の尺度が混在する系に拡張した。

**定理3 (拡張 Buckingham II-theorem)**

$\phi(x, y, z, \dots) = 0$  が系を記述する変数過不足のない完全な方程式で、かつその中の各数量が間隔、比例、絶対尺度の何れかである場合には、 $\phi = 0$  は以下の形式に書き換え可能である。

$$F(\Pi_1, \Pi_2, \dots, \Pi_{n-r-s}) = 0$$

ここで、 $n$  は  $\phi$  の引数の数、 $r$  と  $s$  は  $x, y, z, \dots$  が有する基礎単位及び基礎原点の数である\*2。また全ての  $i$  について、 $\Pi_i$  は絶対尺度量であり、各々定理2の式により求められる。

\*1  $C$  は  $I$  の1つの被覆であり、 $C_{\bar{v}}$  は  $I - I_g$  ( $I_g \subseteq I$ ) の1つの被覆である。左辺の  $\Pi$  は間隔、比例、絶対何れの尺度であってもよい。各係数は定数である。

\*2 基本単位とは、 $\phi$  において他の次元とは独立に数量の尺度を決める単位次元のこと(例：加速度を規定する単位次元は時間と距離)、基礎原点とは、間隔尺度の測定において基準点として人為的に選ばれる原点のことである(例：摂氏温度における0度の定義)。

絶対尺度量  $\Pi_i$  を定義する式  $\rho_i(\Pi_i, x, y, \dots) = 0$  はレジーム、式  $F(\Pi_1, \Pi_2, \dots, \Pi_{n-r-s}) = 0$  はアンサンプルと呼ばれている。レジームはそれ自身であるまとまった意味をもつ現象を意味している。式  $F = 0$  の引数はすべて絶対尺度量であるため、アンサンプルの形式は定理2には従わず任意のものが許される。しかし、多くの例で示されるようにレジームまで規定されればアンサンプルは簡単な式で同定可能である。このように、数量の尺度は法則式の数学的許容性の条件に関して極めて強い制約を与えており、我々の理解にとって受け入れ可能な法則式の形式の多くの部分を決めてしまう。

## 3. 尺度とデータからの構成的法則発見方法

ここまで来れば、全体の式を構成的に求める方法の目処は立つ。幾つかのアルゴリズムを検討した結果、現在の所、以下のアルゴリズムが最も計算量が少なく( $O(n^3)$ )、ノイズにも強いことが確認できている。アルゴリズムを厳密に示す事が目的ではないので、概要のみ示す。

ステップとしては最初にレジームを同定して、次にレジーム間の関係であるアンサンプルを同定する。各数量の尺度情報からレジームの式の形の候補が決まるので網羅的に探索するのが一番単純であるが、個々のレジームにどのような変数が入るのか、そもそもレジーム自身が何個あるのかが分からないので、このようなナイーブな方法では膨大な探索が必要になることは容易に想像がつくであろう。

そこで、まず2つの数量の組み合わせを考え2数量間に成立する関係を検定する。次に1つの数量を共有する2つの2数量の組み合わせに対して無矛盾性を検定し、矛盾がなければ統合する。最初に間隔尺度同士の組み合わせから初めて複合線形式を同定する。次に比例尺度同士および既に得られている複合線形式と比例尺度の組み合わせについて複合積形式を同定する。最後に複合線形式と複合積形式の組み合わせに対して線形対数式を同定する。各ステップではノイズや誤差の影響を考慮してF検定、 $\chi^2$ 検定、ガウス検定などの統計的検定を実施し、早い時期に不適合な組み合わせを棄却する。ここまでのステップで求められた塊はレジームである可能性が強いが単位次元情報を使っていないので、絶対尺度(無次元)量になっているという保証はない(擬似レジームと呼ぶ)。次に擬似レジーム同士の間に線形、積関係などの

簡単な関係が成立するかどうかを同様な手法で検定し、これ以上統合できなくなるまで擬似レジームをまとめて行く。これで1つにまとめられればアンサンプル(つまり法則式)が求まったことになる。1つにまとめられない場合は、より複雑な関数を仮定して全体をフィットするしか方法がなく、仮に求められたとしても第一原理と呼んでいい根拠はない。幸いにして今までに実施した例ではすべてこの段階で法則が求められている。

#### 4. プロトタイプシステムの試作と評価

上述のアルゴリズムを実装した検証用プログラムを構築した。検証用プログラムは実験環境シミュレータが有する対象式を知らずに、シミュレータを実験操作して種々の数量の値の組みを表すデータを取得する。

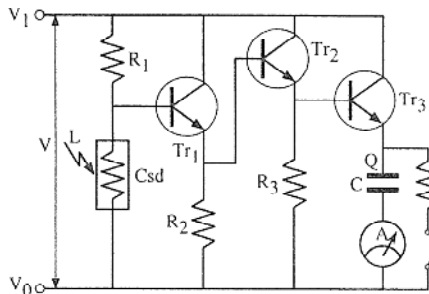


図1 光測定の電子回路

一例として図1に示される3石トランジスタからなる一定時間内の光量増加率を測定する回路について、回路方程式を同定する問題を示す。この系は以下の18個の数量からなる、

$$RQ(\text{比例尺度}) = \{L, r, R_1, R_2, R_3, h_{ie1}, h_{ie2}, h_{ie3}, Q, C, X, K, B\}, IQ(\text{間隔尺度}) = \{V_1, V_2\}, AQ(\text{絶対尺度}) = \{h_{fe1}, h_{fe2}, h_{fe3}\}$$

ここで、 $L, r$ は光量と光素子Csdの感度、 $X, K, B$ は表示電流計針の位置、針パネの定数、磁石磁場の強さを表す。また、 $h_{ie1}, h_{fe1}$ は、それぞれ*i*番のトランジスタのベース入力インピーダンスと電流増幅率を表す。さらに $V_0, V_1$ は電源の正極、負極の電位、 $R_i$ は各抵抗の抵抗値、 $C, Q$ はコンデンサの容量、蓄積電荷を表す。

最初に実験データから間隔尺度である $V_1$ と $V_2$ の組み合わせが検定され $\Pi_1 = V_1 - V_0$ が得られた。次に

$RQ$ 内の数量及び $\Pi_1$ の間の2つの数量間の組み合わせが検定され11組の2数量間関係の集合が得られた。無矛盾性の検定を経て積形式からなる6個の2~4数量からなる擬似レジーム $\Pi_i$ にまとめられた。該当する対数関係式は見つからず、さらにこれらの $\Pi_i$ と $AQ$ 内の無次元量同士の2数量の関係として線形並びに積関係を想定し、4個の中間変数 $\theta_i$ を得た。これから1数量を共有する3つの線形式が得られ、最終的にこれらをまとめて1つの多重線形式が得られた。その結果、正しい回路方程式(1)が得られた。

$$\left( \frac{R_3 h_{fe2}}{R_3 h_{fe2} + h_{ie2}} \frac{R_2 h_{fe1}}{R_2 h_{fe1} + h_{ie1}} \frac{rL^2}{rL^2 + R_1} \right) (V_1 - V_0) - \frac{Q}{C} - \frac{Kh_{ie3}X}{Bh_{fe3}} = 0 \quad (1)$$

中間段階では各 $\Pi$ や $\theta$ は無次元にはならないが、最終的には本来の法則式を反映したモデルが得られることが確認された。

この他に多くの例で性能を検証した。ノイズに対するロバスト性を示す結果を3に示す。最右欄の数値が示すようにボトムアップの探索の効果により、極めてノイズに対して頑強な性能を示している\*3。

表3 計算量とノイズロバスト性に関する評価

例	n	TC	NL(S)
Fechner	2	0.34	±45%
開放感	4	1.06	±40%
理想気体	4	1.00	±40%
運動量保存	8	6.14	±35%
Coulomb	5	1.63	±35%
Stoke's	5	1.59	±35%
運動エネルギー	8	6.19	±30%
電子回路	18	21.9	±20%

n: 属性量の数, TC: CPU計算時間, NL(S): ノイズ許容限界。

#### 5. 今後の展望

以上のように提唱する方式では、10変数を超える比較的規模の大きい問題に関しても、その挙動を表す法則式を現実的な条件で発見可能である。10変数を超える規模の対象事象をまとめて一度にモデル化することは、通常、科学者や技術者にとっても困難な場合が多く、本研究のような体系的手法を採用するメリットは大きい。誌面の都合で法則が1本の式

\*3 各数量に印可可能な、正しい法則が得られるノイズ量の限界を%で示している。

で書ける場合しか紹介しなかったが、連立方程式の場合にも適用可能であり、60変数、26本の連立方程式で書けるプラントの動特性も正しく求められている。

現在の手法は、必要な時に必要なデータが実験によって入手可能であることを前提としている。社会現象など、パラメータを変えて歴史を再現することが不可能な現象に対して(つまり与えられたデータだけから法則を発見できるかどうか)も、適用可能にするためには大きな困難が横たわっている。データが十分にある場合に対してはある程度の見通しは得られているが、ノイズに対する頑強性などはまだ明らかでない。今後、これらの問題に対してこの手法がどこまで適用可能であるか、その限界を明らかにすると共に、単位次元のよく分からない未知の分野で新しい法則を発見して行きたいか考えている。

## 謝 辞

本研究は筆者の所属する高次推論研究分野の鷺尾隆助教授との共同研究によるものである。

## 参 考 文 献

- 1) Bridgman, P. W. (1922). *Dimensional Analysis*. New Haven, CT : Yale University Press.
- 2) Buckingham, E. (1914). On physically similar systems ; Illustrations of the use of dimensional equations. *Physical Review*, IV(4), 345-376.
- 3) Luce, R. D. (1959). On the Possible Psychological Laws. *The Psychological Review*, 66(2), 81-95.

