

## 線図形の類似度とその計算法

田中栄一 粟野宏昭 増田澄男

神戸大学工学部

Cox 等は凸多角形間の類似度の尺度として、1つの距離関数を定義し、その計算法を提案している。本文は、この距離関数を一般線図形に適用したときの計算法について述べている。ここで線図形とは、線分が連結していても、そうでなくてもよく、また、2次元でも3次元でもよい。計算時の枝刈りのとき、最急降下法を用いることができる。この距離測度を化合物の形の比較に応用したとき、100例の計算で、力づくの方法に比べ、平均36.4%の計算時間で済んだ。

## A PROXIMITY MEASURE OF LINE DRAWINGS AND ITS COMPUTING METHOD

Eiichi Tanaka, Hiroaki Awano and Sumio Masuda

Kobe University, Faculty of Engineering  
1-1, Rokkodai, Nada, Kobe, 657 Japan

Cox et al. proposed a distance function between convex polygons and its computing method. This paper describes a distance function between general line drawings and its computing method. A general line drawing in this paper is not only a connected one but also nonconnected one, and is not only a 2-D one but also a 3-D one. This metric can be applied to comparison of the shapes of chemical compounds. The method of steepest descent can be used in the process of pruning. The computing time was reduced on average to 36.4% of that of a brute force method in computation of distances of 100 pairs of chemical compounds.

# 1. まえがき

3次元物体の類似度を測ることはパターン認識の基本的な問題の1つであり、多くの応用がある。たとえば、分子の匂いはその形と強い関係があり [1]、Amoore 等 [2] は PAPA (確率的パターン認識機械) を用いて分子のシルエット間の類似度を測った。しかし、この方法は3次元物体には適切ではない。Marsili 等 [3] は次のような非常に単純な類似度を定義した。2つの分子を重ね合わせるとき、共通部分の体積を  $T_c$ 、両者の体積を  $T_u$  とする。分子間の類似度は  $S = T_c/T_u$  の最大値とし、会話形システム DRACO で  $S$  の最大値を求めた。分子  $A, B$  の体積を  $T_a, T_b$  として、Meyer 等 [4] は  $S = T_c/(T_a \cdot T_b)^{0.5}$  の最大値を類似度としたが、自動的な計算方法ではなかった。

最近 Cox 等 [5] は凸多角形間の距離関数を定義し、多角形の平行移動に関しては、距離関数が単峰であることを示した。従って、多角形の回転がないときには、最急降下法で距離を求めることができる。Bloch-Boulanger 等 [6] はこの距離を凸多面体に適用している。凸多角形は線図形の特殊な場合に過ぎない。Cox の距離を一般の線図形に適用できるように修正すると、距離関数が多峰になる。しかし、多くの応用では正確な類似度は必要ではなく、近似値でよいことが多い。

本文は、一般線図形の類似度尺度として、1つの距離関数を定義する。これは Cox の距離と基本的に等しい。この距離関数に基づいて距離の近似値を求める際に、枝刈り時に最急降下法が適用できることについて述べる。香料分子の形状の類似度を計算したときの、枝刈りの効果についても述べる。なお、本文は文献 [7] に基づいている。

## 2. 線図形間の類似度

2つの線図形  $A, B$  を考える。 $A, B$  はそれぞれ  $m, n$  個の頂点 (extreme points) を持つものとし、次のように表す。

$$A = (Va, Ea), \quad B = (Vb, Eb) \quad (1)$$

ここで、 $V$  は頂点の集合、 $E$  は辺の集合であり次のように表す。

$$Va = \{a_1, a_2, \dots, a_m\}, \quad Vb = \{b_1, b_2, \dots, b_n\}, \\ Ea = \{(a_i, a_j) | a_i \text{ と } a_j \text{ を結ぶ線があるとき}\}, \\ Eb = \{(b_i, b_j) | b_i \text{ と } b_j \text{ を結ぶ線があるとき}\}.$$

$e$  を  $B$  の辺とし、 $d(a_i, e)$  を  $a_i$  から  $e$  への距離とする。 $A, B$  を凸多角形とし、 $A$  と  $B$  の距離を次のように定義する。

$$\bar{b}(a_i : B) = \begin{cases} \min_{e \in Eb} \{d(a_i, e)\}, & \text{if } a_i \notin \text{Int}(B), \\ 0, & \text{そうでないとき.} \end{cases} \\ \bar{D}(A, B) = \sum_{i=1}^m \bar{d}(a_i, B)^2, \\ \bar{D}(B, A) = \sum_{i=1}^n \bar{d}(b_i, A)^2. \\ \bar{D}(A : B) = \bar{D}(A, B) + \bar{D}(B, A). \quad (2)$$

ここで、 $\text{Int}(B)$  は  $B$  の内部を表す。 $A$  の重心を座標軸の原点に置く。 $B$  の重心の位置を  $(x, y)$ 、その回りの回転角を  $\theta$  とすると  $B$  を  $B(x, y, \theta)$  と表せる。Cox 等 [5] 及び Bloch-Boulanger [6] は次のことを証明した。

$A$  と  $B$  が凸多角形あるいは凸多面体のとき、距離関数  $\bar{D}(A : B(x, y, \theta))$  は、 $\theta$  が固定されているとき、 $(x, y)$  に関して単峰である。

しかし、 $\bar{D}$  は  $(x, y)$  が固定されているとき、 $\theta$  に関して単峰ではない。従って、 $\theta$  を適当な刻み巾で、 $0 \leq \theta < 360$  の範囲で変化させ、それぞれの  $\theta$  に対して、 $\bar{D}$  を最急降下法で計算し、すべての場合の最小値を求めればよい。本文では、一般線図形  $A, B$  間の距離を次のように定める。

$$d(a_i, B) = \min_{e \in Eb} \{d(a_i, e)\}, \\ d(b_i, A) = \min_{e \in Ea} \{d(b_i, e)\}. \\ D(A, B) = \sum_{i=1}^m d(a_i, B)^2, \\ D(B, A) = \sum_{i=1}^n d(b_i, A)^2. \\ D(A : B) = D(A, B) + D(B, A). \quad (3)$$

(2) の距離関数は凸多角形及び凸多面体を対象に定義されているため、図形の内部 ( $\text{Int}$ ) を考えるこ

とができるが、本文では一般線図形を対象にしている  
 ので、(3)では図形の内部を定義していない。一般  
 線図形に対しては、 $D(A : B(x, y, \theta))$  は固定した  $\theta$   
 のときも  $(x, y)$  に関して単峰ではない。

### 3. 枝刈り法

線図形  $A$  の閉包を  $P(A)$  と書く (図1)。  $P(A)$  と  
 $P(B)$  の共通領域を  $R(A, B)$  とし、  $P(A) = P(A) -$   
 $R(A, B)$ ,  $R(B) = P(B) - R(A, B)$  とする。図2の  
 例では、  $R(A) = R_1(A) \cup R_2(A)$ ,  $R(B) = R_1(B) \cup$   
 $R_2(B)$  である。  $D(A : B)$  は次のように書ける。

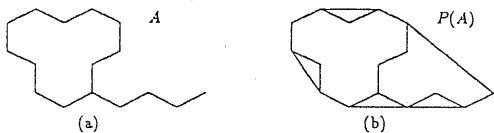


図1. (a) 化合物  $A$  の線図形,  
 (b) 線図形  $A$  を覆う凸包  $P(A)$

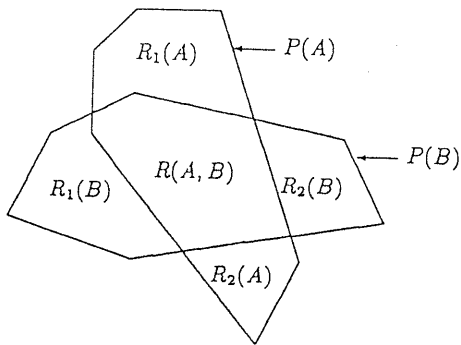


図2. 2つの凸包  $P(A)$  と  $P(B)$  の重ね合わせ

$$\begin{aligned}
 D(A : B) &= \sum_{a_i \in P(A)} d(a_i, B) + \sum_{b_i \in P(B)} d(b_i, A) \\
 &= \sum_{a_i \in R(A)} d(a_i, B) + \sum_{a_i \in R(A, B)} d(a_i, B) \\
 &+ \sum_{b_i \in R(B)} d(b_i, A) + \sum_{b_i \in R(A, B)} d(b_i, A) \\
 &\geq \sum_{a_i \in R(A)} d(a_i, B) + \sum_{b_i \in R(B)} d(b_i, A).
 \end{aligned} \tag{4}$$

また、次の関係がある。

$$\begin{aligned}
 d(a_i, B) &\geq d(a_i, P(B)), \text{ if } a_i \in R(A), \\
 d(b_i, A) &\geq d(b_i, P(A)), \text{ if } b_i \in R(B).
 \end{aligned} \tag{5}$$

$\tilde{D}(A : B)$  を次のように置く。

$$\begin{aligned}
 \tilde{D}(A : B) &= \sum_{a_i \in R(A)} d(a_i, P(B)) \\
 &+ \sum_{b_i \in R(B)} d(b_i, P(A)).
 \end{aligned} \tag{6}$$

このとき次の関係がある。

$$D(A : B(x, y, \theta)) \geq \tilde{D}(A : B(x, y, \theta)). \tag{7}$$

ここで  $d(a_i, P(B))$  は、 $\theta$  を固定したとき  $(x, y)$  関  
 して単峰であることに注意しよう。単峰関数の  
 和は単峰であるから、 $\tilde{D}(A : B)$  は同じ条件下で  
 単峰である。従って、 $\tilde{D}(A : B)$  は最急降下法で  
 計算できる。  $D(A : B(x, y, \theta))$  の計算時に、まず  
 $\tilde{D}(A : B(x, y, \theta))$  を計算し、その値がこれまでの  $D$   
 の最小値より大きいときは、(7)の関係があるので、  
 $D(A : B(x, y, \theta))$  を計算しなくてもよい。本文での  
 線図形は連結図形でも非連結図形でもよく、また2  
 次元図形でも3次元図形でもよい。

### 4. 計算法

距離関数  $D(A : B(x, y, \theta))$  は、固定した  $\theta$  に  
 対しても多峰関数であるから、非常に多くの  $(x, y)$   
 について  $D(A : B(x, y, \theta))$  を計算しなければならない。  
 本文では、 $\theta$  を固定したとき、多くの指定し  
 た  $(x, y)$  について  $D$  を計算する代わりに、複数の  
 初期値  $(x, y)$  から出発して、 $D$  を最急降下法で計算  
 することを、力ずくの方法 (brute force method) と  
 呼ぶことにする。前節で述べたように、固定した  $\theta$   
 に対して、 $\tilde{D}$  は任意の1つの初期値  $(x, y)$  から出  
 発して、最急降下法で計算できる。従って、計算途  
 中では  $(x, y)$  は固定していないので、本節では記号  
 $D(A : B(x, y, \theta))$  及び  $\tilde{D}(A : B(x, y, \theta))$  の代わり  
 に、記号  $D(A, B(\theta))$ ,  $\tilde{D}(A, B(\theta))$  を用いる。また、  
 $n \cdot \Delta\theta = 360^\circ$  とする。

Procedure  $Dist(A, B, \theta, dmin)$

begin

$\theta$  に対して、 $(0, 0)$  の近くの幾つかの初期位置  $(x, y)$   
 から、最急降下法で  $D(A : B(\theta))$  を計算する。

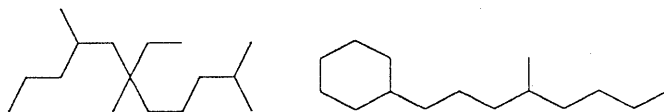


図3. 2つの線図形

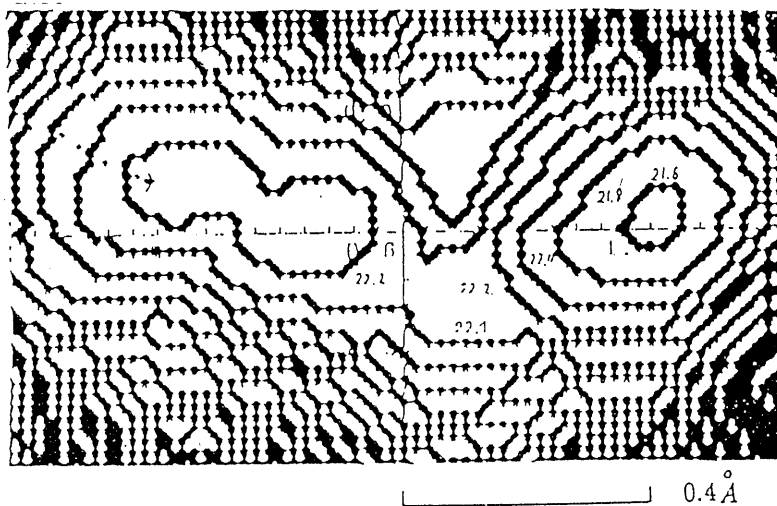


図4. 2つの線図形間の距離の等高線

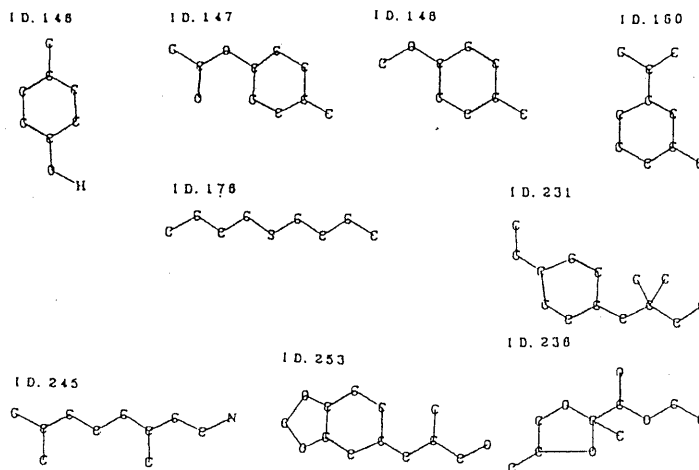


図5. 分子の形 (殆どの水素分子は省略している)

表 1. 線図形間の距離

ID.A/ID.B	147	148	160	176	231	236	245	253
146	3.7	0.5	5.0	11.4	10.9	7.1	13.8	11.8
147		3.7	8.2	10.4	4.8	6.5	10.5	4.7
148			5.3	10.8	9.7	6.4	12.5	10.4
160				17.8	17.5	13.9	16.6	13.6
176					11.1	8.7	3.2	9.5
231						6.0	10.0	5.8
236							8.6	7.4
245								9.2

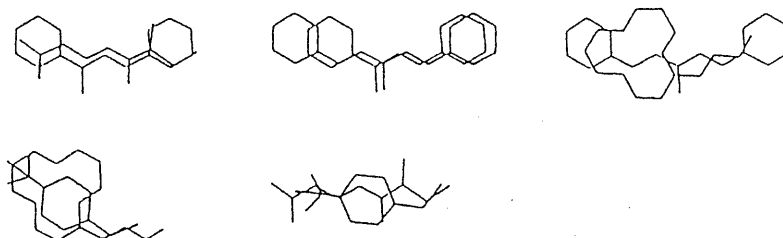


図 6. 距離を最小にする重なり位置の例

$dmin := \min\{D(A : B(\theta))\}$ .  
end.

Procedure *Distmin*

begin

$Dist(A, B, 0, dmin)$ .

$Dmin := dmin$ .

for  $i := 1$  to  $n$  do

begin

$\theta := \theta + \Delta\theta$ .

最急降下法で  $\tilde{D}(A : B(\theta))$  を計算する。

if  $\tilde{D}(A : B(\theta)) < Dmin$ ,

$Dist(A, B, \theta, dmin)$ .

if  $dmin < Dmin$ ,

$Dmin := dmin$ .

end

$D(A : B) := Dmin$ .

end.

## 5. 実験結果

図 3 の線図形間の距離の等高線を図 4 に示す。この図では少なくとも 2 つの極小値があることがわかるが、その差は少ない。分子間の距離の計算では次のことが観察できた。

- (1)  $r$  を隣接原子間の平均距離とすると、線図形間の距離の最小値は  $[-r \leq x \leq r, -r \leq y \leq r]$  の範囲にある。
- (2) 上記の範囲内にある複数の極小値の差は極めて小さい。

図 5 の線図形間の距離を表 1 に示す。初期値を  $(0,0), (1,1), (1,-1), (-1,1), (-1,-1)$  (単位は Å) として  $Dist(A, B, \theta, dmin)$  で計算することを、ここで力づくの方法とする。また、 $n = 180, \Delta\theta = 2^\circ$  とした。 $t_1, t_2$  を力づくの方法及び本方法での平均計算時間とする。100 例の実験では、 $t_1 = 1278.7, t_2 = 465.9$  (秒) であった。本方法で計算時間が 36.4% に短縮されたことになる。計算機は WS NEWS (SONY) を用いた。図 6 に最もよく合った位置の例を示す。

## 6. あとがき

Cox の距離関数を少し修正して一般線図形間の距離関数として用いたときの計算法について述べた。この距離関数は一般線図形に対して多峰であるが、計算の枝刈りに最急降下法が使える。100 対の香料分子間の距離の計算では、力ずくの方法に比べて、36.4% の計算時間で済んだ。

謝辞 化合物データを頂いた花王文理科学研究所  
美濃所長、原稿作成にご協力いただいた田中圭子氏  
に深謝する。

### 文献

1. J.E. Amoores: 匂い-その分子構造-, 恒星社厚生閣 (1970).
2. J.E. Amoores, G. Palmieri and E. Wanke :  
Molecular shape and odour : Pattern analysis by PAPA, Nature, Vol.216, pp.1084-1087 (1967).
3. M.Marsili, P. Floersheim and A.S. Dreiding :  
Generation and comparison of space-filling molecular models, Computers and Chemistry, Vol.7, No.4, pp.175-181 (1983).
4. A.Y. Meyer and W.G. Richards : Similarity of molecular shape, J. Computer-Aided Molecular Design, Vol.5, pp.427-439 (1991).
5. P. Cox, H. Maitre, M. Minoux and C. Ribeiro :  
Optimal matching of convex polygons, Pattern Recognition Letters, Vol.9, pp.327-334 (1989).
6. I. Bloch-Boulanger, H. Maitre and M. Minoux :  
Optimal matching of 3-D convex polyhedra, Telecom Paris 89C003 (1989).
7. E. Tanaka, H. Awano and S. Masuda : A proximity measure of line drawings for comparison of chemical compounds, Proc. CAIP'93, pp.291-298 (1993).