

## 化学物質発ガン性データベースの開発及び サポートベクターマシンによる発ガン性予測

田辺 和俊<sup>1</sup>、鈴木 孝弘<sup>2</sup>、貝原 巳樹雄<sup>3</sup>、小野寺 夏生<sup>1</sup>

<sup>1</sup>筑波大学大学院図書館情報メディア研究科(〒305-8550 つくば市春日 1-2)

<sup>2</sup>東洋大学経済学部経済学科(〒112-8606 東京都文京区白山 5-28-20)

<sup>3</sup>一関工業高等専門学校物質化学工学科(〒021-8511 一関市萩荘字高梨)

### 【緒言】

ガンは日本人の死因の第1位であるが、その原因には物理的(放射線、紫外線等)や生物的(ウイルス、ピロリ菌等)要因以外の化学物質による発ガンが最も大きい。しかし、環境中に存在する化学物質の中で発ガン性が判明しているものは僅少であり、動物愛護の観点から動物実験のスクリーニングとして、化学物質の発ガン性を構造から予測する手法が注目されている。現在、構造の類似した同族体については発ガン強度を高い精度で予測できる場合もあるが、任意の構造の化学物質の発ガン性を高精度で予測することは困難である。

この原因の1つは、非同族体の発ガン性を統一的に予測するモデルを開発するために必要な信頼性の高い発ガン性データの不足である。動物実験に基づく発ガン性データは、IARC、NTP等の機関で収集、公開しているが、発ガン性のランクに統一性がなく混乱している。

もう1つは、非同族体の発ガン性を構造のみから予測できるかという問題である。Hansch-FujitaによるQSARは元来、同族体を対象としたものであるが、化学物質の発ガン機構は複雑であり、しかも化学物質により発ガン機構への関与の仕方が異なるため、非同族体の発ガン性を構造のみから統一的に予測することは簡単ではない。我々はこのようなきわめて複雑な予測問題に有効なニューラルネットワーク(ANN)を用いて非同族体の発ガン性予測を検討したが、多数の局所解の存在のために最適解が得られず、予測精度を確定できないという問題があった[1]。

そこで本研究では、各種発ガン性DBに収録さ

れている多数の化学物質についての発ガン性を総合評価して大規模な発ガン性DBを構築し、近年非線形解析法として注目されているサポートベクターマシン(SVM)[2]を用いてこの発ガン性データを解析することにより、非同族体の発ガン性予測の可能性を検討した。

### 【発ガン性データベースの開発】

動物実験による発ガン性の実測データは、信頼性が高いIARC、EU、EPA、NTP、ACGIH、JSOHの6種のDB[3]から収集した。しかし、それぞれのDBにより発ガン性のランクやその信頼度に関する表現が混乱している[4]。そこで、これらの発ガン性のランクを統一するために、PRTR-MSDSのランク付け[5]を参考に統一基準を設定し、これに基づき表1のように発ガン性の信頼度をA(+++)、B(++), C(+), D(+/-), E(-)の5段階にランク付けした。総物質数は1508である。

表1 発ガン性データの統一ランク

ランク	A	B	C	D	E
発ガン性の信頼度	+++	++	+	+/-	-
IARC	1	2A,2B		3	
EU	1	2	3		
EPA	A	<b>B1,B2</b>	B1,B2,C	D	
NTP	A	<b>B</b>	B		-
ACGIH	A1	<b>A2,A3</b>	A2,A3	A4	
JSOH	1	<b>2A,2B</b>	2A,2B		
物質数	167	407	186	631	117
設定値	0.9	0.7	0.5	0.3	0.1

注：ランクBにおける太字は複数の機関で格付けされている場合を表す。

## 【SVMによる発ガン性予測】

以上の発ガン性データを有する化学物質の中から、QSAR解析が可能な化学物質として、化学構造が不明確なもの、金属原子等を含むもの、混合物、高分子等を除く908種の化学物質について、平面構造から立体構造を生成し、構造を最適化した後、CACHe Project-Leader(富士通)に入力して、表2に示す66種の記述子を作成した。

SVMによる解析にはLIBSVM [6] を用い、発ガン性データは表1に示した値に設定し、SVMによる回帰(SVR)を行った。最適条件の探索や予測精度の評価のために、全ての化学物質をランクごとに分子量順にソートしてセットIとセットIIに交互に振り分け、cross-validation-testを行った。

SVMではANNと異なり局所解の問題はなく、解は一義的に決まる。ただし、SVMではカーネルの選択の問題があり、問題によって最適のカーネルを選択する必要がある。LIBSVMではlinear、polynomial、radial basis function(RBF)及びsigmoidの4種類のカーネルが用意されているので、これらを試行して最適のカーネルを探索した。

また、ANNと同様、SVMでも過学習の問題があり、パラメータの最適設定が必要である。LIBSVMのSVRでは2個のパラメータ $g$ (gamma)と $c$ (cost)が調節可能であり、これらのパラメータの数値を大きくすると過学習状態に陥る。そこで、これら2個のパラメータをグリッドサーチして検証セットの予測誤差が最小になる最適条件を探索した。

Cross-validation-testによりモデルの最適化を行った結果、得られた予測値のクロス集計表を表3に示す。設定値と同一のランクに予測できた物質

表2 解析に用いた記述子

分類	記述子	個数
構造的	各種の原子・結合・原子団の個数、環の大きさ・個数等	41
電子的	電荷、双極子モーメント、配座・立体エネルギー、HOMO・LUMOのエネルギー等	11
立体的	分子容、表面積、分子屈折、カップ形状指数等	6
その他	分子結合指数、logP、分子量等	8

表3 発ガン性予測結果(%)のクロス集計表

予測値	設定値				
	0.9	0.7	0.5	0.3	0.1
1.0~0.8	15	1	0	0	0
0.8~0.6	45	28	11	6	4
0.6~0.4	35	63	69	67	58
0.4~0.2	5	8	19	25	32
0.2~0.0	0	0	1	2	6
計	100	100	100	100	100

の割合は全物質中のわずか29%に過ぎず、また動物実験のスクリーニングという視点からの的中率、すなわち予測値が0.4以上になった物質の中で実際に発ガン陽性の物質の割合も66%に過ぎない。しかし、誤判定の中で特に問題なのは発ガン陽性を陰性と判定するfalse negativeであるが、A~Cランクの物質の中で発ガン陰性と予測された物質は1%以下である。

以上の結果から、今回の発ガン性データに対するSVMモデルの予測性能は充分ではないが、有効な記述子を追加することで、実用的な予測レベルまで向上すると考えられる。SVMは、予測性能はANNより若干劣るものの、局所解の問題がなく最適解が一義的に決まり、ANNより圧倒的に短い時間で処理可能という利点がある。SVMを用いて大規模な発ガン性データを解析することにより、不特定の化学物質の発ガン性を高い精度で予測するシステムの開発が期待できる。

## 参考文献

- [1] 田辺和俊, 大森紀人, 小野修一郎, 鈴木孝弘, 松本高利, 長嶋雲兵, 上坂博亨, *J. Comp. Chem. Jpn.*, **4**, 89-100 (2005)
- [2] V. N. Vapnik, *The Nature of Statistical Learning Theory*, Springer (1995)
- [3] 日本化学物質安全・情報センター, 化学物質の発がん性評価とその分類基準 (2007)
- [4] 田辺和俊, 大森紀人, 小野修一郎, 松本高利, 長嶋雲兵, 上坂博亨, 鈴木孝弘, *情報知識学会誌*, **16**, 63-84 (2006)
- [5] 浦野紘平, PRTR-MSDS 対象化学物質の毒性ランクと物性情報, 化学工業日報社 (2001)
- [6] C. C. Chang, C. J. Lin, <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>