

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5448447号  
(P5448447)

(45) 発行日 平成26年3月19日(2014.3.19)

(24) 登録日 平成26年1月10日(2014.1.10)

(51) Int.Cl. F I  
**G06F 19/16 (2011.01)** G O 6 F 19/16  
**G06F 19/24 (2011.01)** G O 6 F 19/24

請求項の数 2 (全 32 頁)

(21) 出願番号	特願2008-517917 (P2008-517917)	(73) 特許権者	504132272
(86) (22) 出願日	平成19年5月25日 (2007.5.25)		国立大学法人京都大学
(86) 国際出願番号	PCT/JP2007/060736		京都府京都市左京区吉田本町36番地1
(87) 国際公開番号	W02007/139037	(74) 代理人	100115200
(87) 国際公開日	平成19年12月6日 (2007.12.6)		弁理士 山口 修之
審査請求日	平成22年5月24日 (2010.5.24)	(74) 代理人	100158366
(31) 優先権主張番号	特願2006-147433 (P2006-147433)		弁理士 井戸 篤史
(32) 優先日	平成18年5月26日 (2006.5.26)	(72) 発明者	奥野 恭史
(33) 優先権主張国	日本国 (JP)		京都府京都市左京区吉田下阿達町46-2 9 京都大学大学院薬学研究科内
前置審査		(72) 発明者	種石 慶
			京都府京都市左京区吉田下阿達町46-2 9 京都大学大学院薬学研究科内

最終頁に続く

(54) 【発明の名称】ケミカルゲノム情報に基づく、タンパク質-化合物相互作用の予測と化合物ライブラリーの合理的設計

(57) 【特許請求の範囲】

【請求項1】

(A)

第1の化学物質は化合物であり、

第2の化学物質は核酸、タンパク質又はそれらの複合体であり、

第1の特徴量は、前記第1の化学物質のconstitutional

descriptors、topological descriptors、walk and path counts、connectivity indices

、information indices、2D autocorrelations、edge adjacency indices、burden eigen

value descriptors、topological

charge indices、eigenvalue-based indices、functional group counts、atom-centered

fragments、及びmolecular propertiesから選ばれる複数を組み合わせた記述子がベクトル

として表現され、

第2の特徴量は、前記第2の化学物質のアミノ酸配列からスペクトラム法により計算され

た記述子のみがベクトルとして表現され、

互いに相互作用することが既知である前記第1の化学物質と前記第2の化学物質とについ

て、

該第1の化学物質の第1の特徴量と該第2の化学物質の第2の特徴量とを組み合わせた特

徴ベクトルから、請求項1に記載のデータ処理方法によってを、サポートベクターマシン

によって写像変換することにより、特徴空間に学習モデルを構築する工程と、

(B)

10

20

相互作用の予測対象となる第1の化学物質の第1の特徴量と  
相互作用の予測対象となる第2の化学物質の第2の特徴量とを組み合わせた特徴ベクトルを

写像変換することにより前記特徴空間にマッピングする工程と、

(C)

前記特徴ベクトルの前記特徴空間内の位置によって、

前記相互作用の予測対象となる第1の化学物質と前記相互作用の予測対象となる第2化学物質とが

相互作用するか否かが判定される工程と、

からなるデータ処理方法。

10

【請求項2】

請求項1に記載のデータ処理方法を計算機に実行させるためのデータ処理プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、2種類の化学物質のデータベースを用いて効率よくスクリーニングを行い、合理的設計を行うためのデータ処理方法に関する。詳細には、化合物の化学特性などの化学物質情報と遺伝子の配列情報など生物学的情報との両情報を用いて、化合物-タンパク質相互作用予測するデータ処理方法に関する。

20

【背景技術】

【0002】

ヒトゲノムが公開されている現在、ゲノム情報を利用した医薬応用が注目されている。ゲノム情報に基づく医薬品開発には、タンパク質(遺伝子)と化合物の相互作用を解明することが不可欠であるが、これら相互作用を実験的に解明するには莫大な労力が必要である。タンパク質と化合物との相互作用を予測する従来技術としては、一般に、立体構造モデルを用いた予測システムが知られている(非特許文献1~4)。このシステムは、生体高分子の立体構造情報に基づき、リガンドとの安定な複合体構造およびその結合の強さを推定する方法であり、ドッキングスタディと呼ばれている。

【0003】

30

特表2002-530727号公表公報(特許文献1)は、ケミカル空間の高活性領域を特定し、ライブラリ構築を行うことを記載している。しかし、この方法は、ケミカル情報(構造活性相関情報、ファルマコフォア情報など)のみを用いてケミカル空間を定義しているに過ぎない。

【0004】

また発明者自身の従来法(非特許文献5)は、化合物群とタンパク質群とについて、別々にクラスター解析などの統計処理をした後、化合物処理データとタンパク質処理データとを融合し、2次元マップに表示することにより、タンパク質と化合物の相互作用ペアを推測するものである。

【0005】

40

2004年に米国NIHは国家プロジェクトとしてケミカルゲノミクスプロジェクトを開始した。以来、米国や欧州を中心にゲノム情報の化学分野への応用が世界中で取り組まれている。従って、少なくとも米国など先進国において、効率的な予測方法に対する需要が存在する。

【0006】

【特許文献1】特表2002-530727号公表公報

【非特許文献1】Yoshifumi Fukunishi, Yoshiaki Mikami, and Haruki Nakamura "Th efillingpotential method:A methodforestimating the free energy surfaceforprotein -liganddocking" J.Phys. Chem. B.(2003) 107, 13201-13210.

【非特許文献2】Shoichet, B.K., D.L. Bodian, and I.D. Kuntz, "Molecular docki

50

ng usingshapedescriptors.” J.Comp. Chem., 1992. 13(3)、380-397.

【非特許文献3】 Jones G、 WillettP、 Glen RC、 Leach AR、 Taylor R. “ Development and validation of a geneticalgorithm for flexible docking.” J Mol Biol.1997.267(3):727-748.

【非特許文献4】 Rarey M、 KramerB、 Lengauer T. “ Time-efficient docking of flexibleligands intoactivesites ofproteins.” Proc Int Conf IntellSyst Mol Biol.1995; 3:300-308.

【非特許文献5】 Okuno、 Y.、 Yang、 J.、 Taneishi、 K.、 Yabuuchi、 H.、 Tsujimoto、 G、 “ GLIDA:GPCR-LiganddatabaseforChemical Genomic Drug Discovery ” Nucleic AcidsResearch、 34、 D673-677、 2006

10

【発明の開示】

【発明が解決しようとする課題】

【0007】

タンパク質の立体構造モデルを用いた予測手法は、X線結晶解析などで信頼できる立体構造を得ている場合でなければ科学的根拠を有する予測が成り立たないこと、さらに立体構造を用いた計算は負荷が高く計算時間が膨大にかかるため、無限のバリエーションを有する化合物群やタンパク質群の組み合わせを網羅的に計算することが不可能であることなどの問題点を有していた。

【0008】

また、化合物の化学物質情報のみに基づく手法では、予測性能が低く効率化が求められていた。本発明は化合物のスクリーニングを合理的にかつ効率的に行うことを目的とする。

20

【課題を解決するための手段】

【0009】

本発明のデータ処理方法は、第1の化学物質は化合物であり、第2の化学物質は核酸、タンパク質又はそれらの複合体であり、第1の特徴量は、前記第1の化学物質の1種類以上の化学物質情報からなるベクトルとして表現され、第2の特徴量は、前記第2の化学物質の1種類以上の生物学的情報からなるベクトルとして表現されるものであって、第1の特徴量を表す第1空間と第2の特徴量を表す第2空間との相関が高くなるように、多変量解析手法又は機械学習法によって、第1の特徴量および前記第2の特徴量を写像変換するデータ処理方法である。

30

【0010】

また、本発明の他のデータ処理方法は、第1の化学物質と第2の化学物質とが相互作用するか否かが判定されるデータ処理方法であり、以下の工程(A)~(C)からなる。(A)互いに相互作用することが既知である第1の化学物質と第2の化学物質とについて、該第1の化学物質の第1の特徴量および該第2の化学物質の第2の特徴量を、第1の特徴量を表す第1空間と第2の特徴量を表す第2空間との相関が高くなるように、多変量解析手法又は機械学習法によって、第1の特徴量および第2の特徴量を写像変換する工程、(B)

相互作用の予測対象となる第1の化学物質の第1の特徴量を写像変換することによって該第1の化学物質を、(A)の写像変換後の第1の特徴量を表す第1空間にマッピングし、相互作用の予測対象となる第2の化学物質の第2の特徴量を写像変換することによって該第2の化学物質を、(A)の写像変換後の第2の特徴量を表す第2空間にマッピングする工程、(C)相互作用の予測対象となる第1の化学物質の、(B)の写像変換後の第1の特徴量の座標位置と相互作用の予測対象となる第2の化学物質の、(B)の写像変換後の第2の特徴量の座標位置とによって、該第1の化学物質と該第2の化学物質とが相互作用するか否かが判定される工程。

40

【0011】

また、本発明の他のデータ処理方法は、クエリ側の予測対象となる化学物質と相互作用すると予測される化学物質を選出するものであって、以下の工程(A)、(B)、(C1)、(D1)(E1)又は(A)、(B)、(C2)、(D2)、(E2)からなる。(A

50

）互いに相互作用することが既知である第1の化学物質と第2の化学物質とについて、該第1の化学物質の第1の特徴量および該第2の化学物質の第2の特徴量を、第1の特徴量を表す第1空間と第2の特徴量を表す第2空間との相関が高くなるように、多変量解析手法又は機械学習法によって、第1の特徴量および第2の特徴量を写像変換する工程、(B)相互作用の予測対象となる第1の化学物質の第1の特徴量を写像変換することによって該第1の化学物質を、(A)の写像変換後の第1の特徴量を表す第1空間にマッピングし、相互作用の予測対象となる第2の化学物質の第2の特徴量を写像変換することによって該第2の化学物質を、(A)の写像変換後の第2の特徴量を表す第2空間にマッピングする工程、(C1)相互作用の予測対象のうち、クエリ側が第1の化学物質である場合、第1空間において、相互作用の予測対象となる第1の化学物質における第1の特徴量が占有する領域を目的領域として算出する工程と、(D1)第2空間において、第1空間の目的領域に対応する標的領域を算出する工程と、(E1)第2空間にマッピングされた第2の化学物質のうち、該標的領域内に存在する第2の化学物質が選出される工程と、又は、(C2)相互作用の予測対象のうち、クエリ側が第2の化学物質である場合、第2空間において、相互作用の予測対象となる第2の化学物質における第2の特徴量が占有する領域を標的領域として算出する工程と、(D2)第1空間において、第2空間の標的領域に対応する目的領域を算出する工程と、(E2)第1空間にマッピングされた第1の化学物質のうち、該目的領域内に存在する第1の化学物質が選出される工程。

10

**【0012】**

また、本発明のデータ処理方法は、第1の化学物質は化合物であり、第2の化学物質は核酸、タンパク質又はそれらの複合体であり、第1の特徴量は、前記第1の化学物質の1種類以上の化学物質情報からなるベクトルとして表現され、第2の特徴量は、前記第2の化学物質の1種類以上の生物学的情報からなるベクトルとして表現され、第1の特徴量と第2の特徴量とを、多変量解析手法又は機械学習法によって写像変換することにより、特徴空間を構築するデータ処理方法である。写像変換される特徴量は、第1の特徴量と第2の特徴量を連結した特徴量であってもよい。

20

**【0013】**

本発明の他のデータ処理方法は、上記データ処理方法で構築した特徴空間を用い、相互作用の予測対象となる第1の化学物質の第1の特徴量と相互作用の予測対象となる第2の化学物質の第2の特徴量とを多変量解析手法又は機械学習法によって写像変換することにより特徴空間にマッピングし、第1の特徴量の特徴空間内の位置と第2の特徴量の特徴空間内の位置とによって相互作用の予測対象となる第1の化学物質と相互作用の予測対象となる第2化学物質とが相互作用するか否かが判定されるデータ処理方法である。また、特徴空間にマッピングされる特徴量は、相互作用の予測対象となる第1の化学物質の第1の特徴量と相互作用の予測対象となる第2の化学物質の第2の特徴量とを連結した特徴量であってもよい。

30

**【0014】**

また、本発明のライブラリは、上記データ処理方法によって得られた判定結果や選出された化学物質から設計される。また、本発明のデータ処理プログラムは、上記データ処理方法を計算機に実行させるためのプログラムである。また、本発明の化学物質の製造方法は、上記データ処理方法を用いて相互作用の予測のクエリ側の化学物質と相互作用すると予測された計算対象側の化学物質を合成するものである。さらに、本発明の化合物は、上記データ処理方法を用いることによって生産された化合物である。

40

**【発明の効果】****【0015】**

本発明により、化合物の化学特性などの化学物質情報とゲノム情報などの生物学的情報との両情報を用いて、化合物-タンパク質の相互作用を予測する計算手法を開発することができた。この手法では、従来型の化合物の化学特性情報のみを用いた予測とは異なり、遺伝子の配列情報を加えることによって予測の精度向上をはかることに成功した。

**【0016】**

50

従って、この計算手法は具体的に次の2つに適用することができる。

- 1) 化合物ライブラリの生物活性に基づく合理的設計が実現できる。
- 2) 従来法よりも性能の良いリード化合物探索ができる。

【0017】

以下に、本発明の好ましい実施形態を示すが、当業者は本発明の説明および当該分野における周知慣用技術からその実施形態などを適宜実施することができ、本発明が奏する作用および効果を容易に理解することが認識されるべきである。従って、本発明のこれらおよび他の利点は、必要に応じて添付の図面等を参照して、以下の詳細な説明を読みかつ理解すれば、当業者には明白になることが理解される。

【発明を実施するための最良の形態】

【0018】

以下、本発明を説明する。本明細書の全体にわたり、単数形の表現は、特に言及しない限り、その複数形の意味をも含むことが理解されるべきである。従って、単数形の冠詞または形容詞（例えば、英語の場合は「a」、「an」、「the」など）は、特に言及しない限り、その複数形の意味をも含むことが理解されるべきである。また、本明細書において使用される用語は、特に言及しない限り、当該分野で通常用いられる意味で用いられることが理解されるべきである。したがって、他に定義されない限り、本明細書中で使用される全ての専門用語および科学技術用語は、本発明の属する分野の当業者によって一般的に理解されるのと同じ意味を有する。矛盾する場合、本明細書（定義を含めて）が優先する。

【0019】

本明細書において「化学物質」とは、物質という一般用語の中で、とくに化学的な立場で物質を取り扱う場合の用語であり、任意の一定の分子構造をもつ物質をいう。

【0020】

本明細書において「空間座標」とは、対象物の空間内の位置を特定するための指標である。「空間」とは、「スペース」と交換可能に使用され、集合の別名であり、なんらかの位相や幾何学的構造を想定するものをいう。空間は、例えば、化学物質情報や生物学的情報により定義される。空間としては、化合物記述子や化学特性などの化学物質情報で定義されるケミカル空間や、発現情報、パスウェイ情報、機能情報、および生物活性などの生物学的情報で定義されるバイオ空間（Current Opinion in Chemical Biology 2005、9:296-303）を例示することができる。

【0021】

本明細書において「特性」とは、ある化学物質が有する特別の性質をいう。特性としては、例えば、融点、沸点、比重などの物理特性、反応性、酸性、アルカリ性などの化学特性、タンパク質構造、酵素活性、レセプターとの結合能、サイトカイン能、細胞との相互作用力などの生物学的特性を挙げることができるがそれらに限定されない。

【0022】

本明細書において、第1特性と第2特性との間に「単純な関連が見られない」とは、両者の間に直感的な相関が見られないことをいう。第1特性を示すベクトルを  $a$ 、 $b$  とし、第2特性を示すベクトルを  $x$ 、 $y$  とし、 $a$  と  $x$ 、 $a$  と  $y$ 、 $b$  と  $x$ 、 $b$  と  $y$  のペアには相関が直感的な相関が見られないとする。しかし、線形結合ベクトル  $m * a + n * b$  と  $M * x + N * y$  ( $m$ 、 $n$ 、 $M$ 、 $N$  は0で無い係数) との間に相関が見られる場合がある。つまり、第1特性と第2特性との間に単純な関連が見られないものであっても、適当な変換を加えた  $F(a, b, \dots)$  と  $F'(x, y, \dots)$  の間に相関が見られる場合がある。

【0023】

本明細書において「化合物」とは、化学変化によって2種またはそれ以上の元素の単体に分けることができる純粋物質をいう。2種以上の元素の原子の化学結合によって生じた純粋物質といってもよい。通常、各元素の組成比は一般に定比例の法則に従って一定であるが、不定比化合物のように組成比がある範囲で連続的に変化しても安定な結晶をつくる

10

20

30

40

50

ものもまた、本明細書において化合物の範疇に入れる。

【0024】

本明細書において「生体物質」とは、生物に関連する任意の物質を言う。生体物質もまた、化学物質の一種として捕らえることができる。本明細書において「生体」とは、生物学的な有機体をいい、動物、植物、菌類、ウイルスなどを含むがそれらに限定されない。生体物質には、タンパク質、ポリペプチド、オリゴペプチド、ペプチド、ポリヌクレオチド、オリゴヌクレオチド、ヌクレオチド、核酸（例えば、cDNA、ゲノムDNAのようなDNA、mRNAのようなRNAを含む）、ポリサッカライド、オリゴサッカライド、脂質、これらの複合分子（糖脂質、糖タンパク質、リポタンパク質など）などが包含されるがそれらに限定されない。通常、生体物質は、核酸、タンパク質、脂質、糖、プロテオリピッド、リポプロテイン、糖タンパク質およびプロテオグリカンなどであり得る。好ましくは、生体物質は、核酸（DNAまたはRNA）またはタンパク質である。

10

【0025】

本明細書において「レセプター」とは、細胞上または核内などに存在し、外界からの因子または細胞内の因子に対する結合能を有し、その結合によりシグナルが伝達される分子をいう。レセプターは通常タンパク質の形態をとる。レセプターの結合パートナーは、通常リガンドという。

【0026】

本明細書において「アゴニスト」とは、あるリガンドのレセプターに結合し、その物質のもつ作用と同等又は類似の作用を現わす因子をいう。また、「アンタゴニスト」とは、あるリガンドのレセプターへの結合に拮抗的に働き、それ自身はそのレセプターを介した生理作用を現わさない因子をいう。拮抗薬、遮断剤（ブロッカー）、阻害剤（インヒビター）などもこのアンタゴニストに包含される。

20

【0027】

本明細書において「相関が最大になるように」「座標を変換」とするとは、ある空間における各要素と別の空間における各要素との関係が、全体としてみたときに最大限に相関している状態を言う。このような定義は種々の計算手法によって達成することができる。

【0028】

本明細書において「相関」とは、数理統計学や生物統計学において、一般に2つまたはそれ以上の変数のあいだの関連性をいい、2つの確率変数の間の直線的な共変関係をいう。数量分類学において、対象とする二つの操作的分類単位（OTU）間の類似係数（類似の程度）によって表現することができる。従って、相関分析とは、対のデータに基づいて相関の有無を検証し、あるいは相関の大きさを推定したりする統計方法のことである。

30

【0029】

本明細書において「相関係数」とは、2つの確率変数X、Yの間の関連を示す、一次変換で不変な量をいう。XとYの共分散をXおよびYのおのおのの分散の積の平方根で割った値である。「相関関数」には、空間の相関関数と時間の相関関数がある。2点の物理量A(r)とA(r')の積の平均値A(r)A(r')を空間の相関関数という。これが2点間の距離に対して指数関数的に変化するとき、これを

【0030】

【数101】

$$\overline{A(r)A(r')} \propto \exp\left(-\frac{|r-r'|}{\xi}\right)$$

40

【0031】

において、相関距離  $\xi$  が定義される。時間の相関関数に対しても同様に相関時間  $\tau$  が定義される。  $\tau$  は臨界点で発散する。

【0032】

本明細書において「相関距離」とは、空間の各点に存在する確率変数A(r)の相関関数  $\overline{A(r)A(r')}$  が2点間距離  $r = |r - r'|$  の増大とともに絶対値として減

50

少していく場合、その減衰の目安となる距離をいう。相関関数が  $\exp(-r/)$  の形である場合、 が相関距離である。減衰が指数関数的でない場合に一義的に定義することはできないが、多くの場合、適当な長さの尺度として定まる。

### 【0033】

本明細書において相関を最大にする方法として、正準相関分析 (CCA)、カーネル正準相関分析 (kernel CCA)、サポートベクターマシン (SVM) 法等の多変量解析手法及び機械学習法がある。

### 【0034】

本明細書において「化学物質情報」とは、化学物質の化合物としての情報をいう。より詳細に定義すると、化学構造自体や化学構造から計算処理によって算出された各種記述子 (文献 Handbook of Chemoinformatics: From Data to Knowledge Gasteiger, Johann (EDT) / Publisher: John Wiley Published 2003/10)、化学構造や記述子より計算推定される化学特性、さらには化合物を計測して得られる化学特性が含まれる。通常、各々の情報は数値または数値列 (ベクトル) として表現される。各々の化合物は、適宜選択した各種化学物質情報の数値列を連結した数値列 (ベクトル) として表現される。

### 【0035】

通常、化学物質情報は、化合物記述子によって定義される。このような記述子としては、例えば、一次元記述子、二次元記述子および三次元記述子からなる群より選択されるものを挙げるができる。「記述子」とは、ディスクリプタともいい、ある情報を記述するための方法およびそれにより表された表現物をいう。記述子は、電波、磁波、音、光、色、画像、数字、文字など又はそれらの組み合わせによって表現することができる。分子構造を特徴づける記述子群の特定は、数多くの化合物を解析するプロセスで重要な工程である。数多くの記述子が提案されているが、分子構造へのアプローチに応じてこれらを分類することができる (M. Hassan et al., Molecular Diversity, 1996, 2, 64; M. J. McGregor et al., J. Chem. Inf. Comput. Sci., 1999, 39, 569; R. D. Brown, Perspectives in Drug Discovery and Design, 1997, 7/8, 31 参照。以上は先に本明細書に参考文献として援用する。R. D. Brown et al., J. Chem. Inf. Comput. Sci., 1996, 36, 572; R. D. Brown et al., J. Chem. Inf. Comput. Sci., 1996, 37, 1; D. E. Patterson et al., J. Med. Chem., 1996, 39, 3049; S. K. Kearsley et al., J. Chem. Inf. Comput. Sci., 1996, 36, 118 参照。以上を本明細書に参考文献として援用する)。1次元 (1D) 特性は、分子量や  $c\log P$  等の全体的な分子特性を表す。2次元特性 (2D) には、分子の機能性や結合性が含まれる。2D記述子の例としては、MDLサブストラクチャーキー (MDL Information Systems Inc., 14600 Catalina St., San Leandro, CA 94577) (M. J. McGregor et al., J. Chem. Inf. Comput. Sci., 1997, 37, 443 参照。これを本明細書に参考文献として援用する) や MSI50記述子 (Molecular Simulations Inc., 9685 Scranton Road, San Diego, CA 92121-3752) 等が挙げられる。例えば、薬剤化合物に対する要件を特定する際に有用な、周知の5つの法則 (rule of five) は、1次元記述子及び2次元記述子から導かれる (C. A. Lipinski et al., Advanced Drug Delivery Reviews, 1997, 23, 3 参照。これを本明細書に参考文献として援用する)。

### 【0036】

3次元記述子 (3D) の算出には、適度なエネルギーを有する少なくとも1つの3次元

10

20

30

40

50

構造体が必要である。更に、複数のコンホメーション（立体配座）からの寄与を考慮にいられて、3次元記述子を算出してもよい。また、リガンド結合において重要な特徴に基づいて、あるいは、その他の重要な所望の特徴に応じて、記述子を選択するようにしてもよい。あるいは、数多くの化合物群の解析に多数の記述子を用いる場合には、主成分分析（PCA）や部分最小2乗法（PLS）等の統計手法により最少数の重要な記述子群を求めればよい。

#### 【0037】

本明細書において「生物学的情報」とは、遺伝子やタンパク質などの生体物質である化学物質の生物学的特性にかかわる情報をいう。生物学的情報とは、例えば、遺伝子の塩基配列情報、一次構造情報、二次構造情報、三次構造情報、四次構造情報、立体構造情報、発現情報、パスイエイ情報、機能情報、および生物活性情報などが挙げられる。本発明においては、この生物学的情報は、計算処理や計測によって数値化した数値を各要素として持つ数値列（ベクトル）として表現される。

10

#### 【0038】

生体物質の特性は、その独自の記述子で生物学的特性を指標に表現することができる。従って、本明細書において、ある状態に関する「指標」とは、その状態を表すための目印となる関数をいう。本明細書では、例えば、生物または細胞であれば、その生物または細胞内の種々の物理的指標（電位、生体内温度、移動速度・距離、局在化率、扁平率、伸長率、回転速度など）、化学的指標（ゲノム量、特定の遺伝子の転写産物（例えば、mRNA、翻訳タンパク質、翻訳後修飾されたタンパク質、イオン濃度、pHなどの量、代謝産物の量、イオン量など）、生物学的指標（個体差、進化速度、薬物応答など）など、あるいはその生物または細胞の環境、例えば、温度、湿度（例えば、絶対湿度、相対湿度など）、pH、塩濃度（例えば、塩全体の濃度または特定の塩の濃度）、栄養（例えば、ビタミン量、脂質量、タンパク質量、炭水化物量、金属イオン濃度など）、金属（例えば、金属全体の量または特定の金属（例えば、重金属、軽金属など）の濃度など）、ガス（例えば、ガス全体の量または特定のガス（例えば、酸素、二酸化炭素、水素など）の量）、有機溶媒（例えば、有機溶媒全体の量または特定の有機溶媒（例えば、エタノール、DMSO、メタノールなど）の量）、圧力（例えば、局所圧または全体の圧（気圧、水圧）など、粘性、流速（例えば、培地中に生物が存在する場合のその培地の流速、膜流動など）、光度（ある特定波長の光量など）、光波長（例えば、可視光のほか紫外線、赤外線なども含み得る）、電磁波、放射線、重力、張力、音波、対象となる生物とは異なる他の生物（例えば、寄生虫、病原菌、細菌、ウイルスなど）、化学薬品（例えば、医薬品、食品添加物、農薬、肥料、環境ホルモンなど）、抗生物質、天然物、精神的ストレス、物理的ストレスなどのような指標に対する反応性または耐性を、そのような状態に関する「指標」として使用することができる。

20

30

#### 【0039】

本明細書において「遺伝子」とは、遺伝形質を規定する因子をいう。通常染色体上に一定の順序に配列している。タンパク質の一次構造を規定する遺伝子を構造遺伝子といい、その発現を左右する遺伝子を調節遺伝子という。本明細書では、特定の状況において、「遺伝子」は、「ポリヌクレオチド」、「オリゴヌクレオチド」および「核酸」ならびに／あるいは「タンパク質」、「ポリペプチド」、「オリゴペプチド」および「ペプチド」をさすことがある。

40

#### 【0040】

本明細書に置いて「一次構造」とは、特定のペプチドのアミノ酸配列をいう。「二次構造」とは、ポリペプチド内の局所的に配置された三次元構造をいう。これらの構造はドメインとして一般に公知である。ドメインは、ポリペプチドの緻密単位を形成し、主に50～350アミノ酸長であるポリペプチドの部分である。代表的なドメインは、 $\alpha$ -シート及び $\beta$ -ヘリックスである。「三次構造」とは、ポリペプチドモノマーの完全な三次元構造をいう。「四次構造」とは、独立した三次単位の非共有結合により形成される三次元構造をいう。

50



## 【0041】

本明細書において用いられる数値処理は、例えば、生命システム解析のための数学、コロナ社、清水和幸(1999)などにおいて記載される周知技術を適用することができる。

## 【0042】

本明細書において「ライブラリ」とは、スクリーニングをするための化合物または生体物質などの化学物質の一定の集合をいう。ライブラリは、同様の性質を有する化合物の集合であっても、ランダムな化合物の集合であってもよい。好ましくは、同様の性質を有すると予測される化合物の集合が使用されるが、それに限定されない。

## 【0043】

本明細書において「相互作用」とは、2以上の分子が存在する場合、ある分子と別の分子との間の作用をいう。そのような相互作用としては、水素結合、ファンデルワールス力、イオン性相互作用、非イオン性相互作用、受容体リガンド相互作用、静電的相互作用およびホスト-ゲスト相互作用が挙げられるがそれらに限定されない。分子間の相互作用を定量化する手法の1つとしては、分子間相互作用の熱力学的あるいは速度論的な定量評価法が挙げられるが、これらに限定されるものではない。分子間相互作用の熱力学的あるいは速度論的な定量評価法としては、例えば、熱測定(calorimetry)、表面プラズモン共鳴法、超遠心分析法などが挙げられるがこれらに限定されない。分子間相互作用は、複合体を形成した状態と解離した状態の熱力学量の変化を示す指標により表現可能であり、例えば、結合定数、解離定数、結合/解離に伴う標準化学ポテンシャル変化、エンタルピー変化、イオン結合数変化などにより表現可能である。

## 【0044】

本明細書において、「相互作用情報」は、例えば、分子間における結合の有無、結合活性、薬理活性などにより表現されるが、これらに限定されるものではない。「薬理活性」は、一般に、薬物の効力を示す指標である。例えば、生理活性の50%阻害効果を示す濃度であるIC50や生理活性の50%亢進効果を示す濃度のEC50として与えられ、また、「結合の有無」および「結合活性」についても、例えば、解離定数Kdとして与えられる。

## 【0045】

本明細書において「情報処理」または「処理」とは、情報をより取り扱いやすくするため、一つの形式から他の形式へ変換または統合することをいう。また、「データ処理」とは、データを一つの形式から他の形式へと変換または統合する過程をいう。

## 【0046】

本明細書において「パターン認識」とは、自然情報処理の1つである。例えば、形態、図形、物体、画像、音声、生理的現象のような単純な数量として与えられない情報を識別し認定することをいう。このような諸情報を、「パターン情報」又は、単に「パターン」という。パターン認識を行う識別器としては、例えば、SVM(サポートベクターマシン)、ベイズ分類、ニューラルネットワークなど、機械学習により大量のデータから識別パラメータを構成する手法を用いることが可能である。本明細書において、「相互作用パターン」は、上述の識別器により統計モデルとして定義づけられる情報をいう。

## 【0047】

本明細書において「モデル」とは、ある特定の条件に従う数学的、物理的な系をいい、物理、生物などの自然科学における系を理解するために用いられる。特に、統計的な解析の対象となる場合、そのモデルを「統計モデル」という。本明細書において使用される場合、ある現象の「モデル化」とは、データの背後にある現象の解明と予測、制御や新たな知識発見のために「モデル」を導入することをいう。

## 【0048】

本明細書においては、集合Xの要素が任意に与えられたとき、関数fによって集合Yの要素がひとつ対応づけられていることを「集合Xから集合Yへの写像」と表現する。また、関数fによって集合Xから集合Yへと移行することを「写像変換」という。本発明においては、

10

20

30

40

50

集合は化学物質の集合を意味する。

【0049】

また、本明細書における「マッピング」とは、トレーニングデータをCCAやPCAやカーネル化などの計算によって写像変換することによって、算出される重み係数行列やカーネル関数をテストデータにかけることでなされる。

【0050】

本明細書におけるライブラリは、例えば、コンビナトリアル・ケミストリ技術、醗酵方法、植物および細胞抽出手順などの手法により、作製又は入手することができる。コンビナトリアル・ライブラリを作成する方法は、当該技術分野で周知である。例えば、E. R. Felder, *Chimia* 1994, 48, 512-541; Gallopら, *J. Med. Chem.* 1994, 37, 1233-1251; R. A. Houghten, *Trends Genet.* 1993, 9, 235-239; Houghtenら, *Nature* 1991, 354, 84-86; Lamら, *Nature* 1991, 354, 82-84; Carellら, *Chem. Biol.* 1995, 3, 171-183; Maddenら, *Perspectives in Drug Discovery and Design* 2, 269-282; Cwirllaら, *Biochemistry* 1990, 87, 6378-6382; Brennerら, *Proc. Natl. Acad. Sci. USA* 1992, 89, 5381-5383; Gordonら, *J. Med. Chem.* 1994, 37, 1385-1401; Leblら, *Biopolymers* 1995, 37, 177-198; およびそれらで引用された参考文献を参照のこと。これらの参考文献は、その全体を、本明細書中で参考として援用する。

【0051】

コンビナトリアル・ケミストリとハイスループット・スクリーニングの最近の発展に伴い、数多くの化合物に対する実験的アプローチが可能になった。例えば、D. K. Agrafiotis et al., *Molecular Diversity*, 1999, 4, 1; U. Eichler et al., *Drugs of the Future*, 1999, 24, 177; A. K. Ghose et al., *J. Comb. Chem.*, 1, 1999, 55; E. J. Martin et al., *J. Comb. Chem.*, 1999, 1, 32; P. R. Menard et al., *J. Chem. Inf. Comput. Sci.*, 1998, 38, 1204; R. A. Lewis et al., *J. Chem. Inf. Comput. Sci.*, 1997, 37, 599; M. Hassan et al., *Molecular Diversity*, 1996, 2, 64; M. J. McGregor et al., *J. Chem. Inf. Comput. Sci.*, 1999, 39, 569; R. D. Brown, *Perspectives in Drug Discovery and Design*, 1997, 7/8, 31を参照のこと。以上を本明細書に参考文献として援用する。このため、数多くの化合物に関する演算特性を解析する技術が、薬剤開発において、ますます重要になってきている。特定ライブラリ、すなわち、標的ライブラリの構築並びにプライマリ・ライブラリの構築という2つの適用例では、数多くの化合物に関する演算特性の解析により、薬剤設計にとって特に重要な情報を提供することができる。

【0052】

以下に本発明を実施するための実施形態の説明を記載するが、この実施形態は本発明を実施するための単なる例示であり、本発明の範囲はそのような好ましい実施形態に限定されないことが理解されるべきである。

【0053】

まず、本発明の実施形態について、図を用いて説明する。本発明者らの手法は、図1上段の化学物質情報のみを用いていた従来手法に、図1下段のバイオインフォマティクス技術を融合させたものであり、これはゲノム情報等に代表される生物学的情報を考慮に入れた新しい相互作用予測方法である。

【0054】

10

20

30

40

50

リード化合物の探索やライブラリ設計を計算で行う際には、個々の化合物の類似度を反映するため、相対的な位置を示す座標空間が必要である。例えば、図1上段のケミカル空間上の丸印はそれぞれ異なる化合物を表しており、特性の似ている化合物は相対的に近い位置関係になるように配置されている。これら化合物の位置から構成される座標空間をケミカル空間という。

【0055】

同様に、タンパク質又はその遺伝子について、類似関係を相対的な位置として表現したものが図1下段のバイオ空間である。図1下段のバイオ空間内の四角印がタンパク質又は遺伝子を表す。さらに、個々の化合物とタンパク質の結合をリンク(図1、中央矢印)することによって、ケミカル空間とバイオ空間を融合したモデルを作ることができる。

10

【0056】

図2は、本発明の概念図であり、本発明の実施態様の一例を示したものである。ケミカル空間とバイオ空間との融合モデルから、クエリとなる化合物についての標的タンパクを予測することができる。まず、1)クエリの化合物(星印)の化学構造からその化合物がケミカル空間座標にマッピングされる(図2上段)。2)ケミカル空間にマッピングされた該化合物の近隣にある化合物からバイオ空間へのリンク情報をたどること(図2中央・矢印)により、その未知化合物が関係するバイオ空間内のエリア(図2下段・円内)を指定することができる。3)このエリア内のタンパク質群が、この活性未知の化合物が相互作用する可能性のあるタンパク群と推定される。

【0057】

20

図3は、本発明の概念図であり、本発明の別の実施態様を示したものである。ケミカル空間とバイオ空間との融合モデルがあれば、図3中央の矢印で示すように、クエリのタンパク質から、該タンパク質に相互作用するリード化合物群を予測することも可能となる。

【0058】

さらに、ケミカル空間とバイオ空間とを融合したモデルは化合物ライブラリの合理的設計にも適用できる。図4に示すように、広大なケミカル空間のうち、限定されたバイオ空間に対応するエリア(biologically relevant chemical space)内の化合物群が、バイオ空間を形成するタンパク群と相互作用する可能性が高いと考えられる。そのため、生物活性を有する化合物ライブラリを合理的に設計することができる。また、バイオ空間のタンパク質を、例えばGPCRファミリーなどに限定するとFocused libraryを設計すること

30

【0059】

ケミカル空間とバイオ空間とを融合したモデルの構築方法について、本手法の特徴を示す。図5を参照しながら説明する。従来法(図5上段)は、化学物質情報のみを用い、化合物の化学特性ができる限り多様になるように、ケミカル空間座標を定義していた。しかし、化合物の多様性と生物活性との直接の因果関係がないという点が問題であった。そこで、本手法は、化学物質情報と生物学的情報との両情報を用いて、ケミカル空間とバイオ空間との両空間の相関が高くなるように互いの空間座標を定義することにした。これはバイオ空間との関連を考慮してケミカル空間座標を定義すると、化合物の生物活性を反映した空間座標を構築できるためである。

40

【0060】

次に、本発明の実施形態をより詳細に説明する。第1の化学物質は化合物であり、第2の化学物質は、核酸、タンパク質又はそれらの複合体等の生体物質である。また、第1の空間は、ケミカル空間と呼ばれ、第2の空間は、バイオ空間と呼ばれる。化合物としては、任意の化合物のライブラリを用いることができる。例えば、コンビナトリアル・ライブラリ、ある企業が有する任意の化合物ライブラリ、公的機関が運営するデータベース、コンソーシアムによって運営されている化合物ライブラリ・データベースなどを挙げるができるがそれらに限定されない。

【0061】

第1の特徴量は化学物質情報からなるベクトルとして表現され、第2の特徴量は生物学的

50

情報からなるベクトルとして表現される。これら第1の特徴量と第2の特徴量とは互いに単純な関連が見られないものであり得る。

【0062】

本発明の空間座標の作成にはどのような情報を用いてもよい。1つの実施形態では、本発明における化合物の空間座標は化学物質情報から定義される。ここで、化学物質情報は、化合物記述子によって定義される。化合物記述子は、一次元記述子、二次元記述子および三次元記述子からなる群より選択される。一次元記述子は、化学組成を記述することを特徴とし、二次元記述子は、化学トポロジーを記述することを特徴とし、三次元記述子は、三次元形状および官能性からなる群より選択される特徴を記述することを特徴とする。

【0063】

また、化合物記述子はファルマコフォアであってもよい。ファルマコフォアは、少なくとも3つの空間的に離れたファルマコフォア中心を含む。各ファルマコフォア中心は、(i)空間位置と、(ii)ある化学特性を特定する所定のファルマコフォア型と、を含む。基本セットのファルマコフォア型には、少なくとも、水素結合受容体、水素結合供与体、負電荷中心、正電荷中心、疎水性中心、芳香族中心、ならびに他のいずれのファルマコフォアの型にも入らないデフォルトカテゴリが含まれる。また、空間位置を、隣接するファルマコフォア中心間の隔絶距離あるいは隔絶距離範囲として与えることによって、ファルマコフォアをより詳細に記述することができる。

【0064】

化学物質は、化合物記述子、化学構造や記述子より計算推定される化学特性、さらには化合物を計測して得られる化学特性の数値を各要素として持つ数値列(ベクトル)として表現される。よって、各々の化学物質は第1空間座標上のベクトルとして位置が特定される。

【0065】

本発明における生体物質の空間座標は、生物学的情報によって定義される。ここで、生物学的情報は、配列情報、二次構造、三次構造、四次構造、立体構造情報、発現情報、パスウェイ情報、および機能情報からなる群より選択される少なくとも1種類の情報を用いることができる。生体物質は、生物学的情報を計算処理や計測によって数値化した数値を各要素として持つ数値列(ベクトル)として表現される。よって、各々の生体物質は第2空間座標上のベクトルとして位置が特定される。

【0066】

好ましい実施形態では、第1空間と第2空間と間の相関が高くなるように、より好ましくは相関が最大になるように、第1空間の特徴量及び第2空間の特徴量が写像変換される。ここで、第1空間と第2空間との相関が高くなるようにするには、正準相関(CCA)、カーネル正準相関(kernel CCA)、サポートベクターマシン(SVM)法などの多変量解析や機械学習法またはそれらの等価方法により行うことができる。

【0067】

正準相関分析(CCA)とは、2種類の異なるデータセット、例えば、化合物とタンパク質が与えられたときに、そのデータセット間の相関関係を解析する多変量解析手法の一種である。CCAは、両データセット間の相関を最もよく表すように写像変換し、それによって2つのデータ間の相関を解析するものである。具体的手順を以下に示す。

【0068】

正準相関分析の具体的手順は、以下のような物を例示することができる(T. W. Anderson. An Introduction to Multivariate Statistical Analysis. Wiley & Sons, 1984.、H. Hotelling. Relations between two sets of variates. Biometrika, 28: 321-377、1936.)。

【0069】

行列X、Yの行には化学物質のエントリが、列には化学物質情報が並ぶ2種の異なるデ

10

20

30

40

50

ータ（例えば化合物とタンパク質）を行列  $X$ 、 $Y$ （第 1 空間が行列  $X$ 、第 2 空間が行列  $Y$ ）と表現したとき、

【 0 0 7 0 】

【数 1 1 1】

$$X = n \times p \text{ 行列} \quad Y = n \times q \text{ 行列} \quad n \geq p \geq q \geq 1$$

$$f = Xa \quad g = Yb$$

$$\sigma_f^2 = \frac{1}{n-1} f'f \quad \sigma_{fg} = \frac{1}{n-1} f'g \quad \sigma_g^2 = \frac{1}{n-1} g'g$$

第 1 空間と第 2 空間の相関を最大にするために、

10

相関係数

【 0 0 7 1 】

【数 1 1 2】

$$r_{fg} = \frac{\sigma_{fg}}{\sigma_f \sigma_g}$$

【 0 0 7 2 】

を最大にする係数ベクトル  $a$ 、 $b$  の組を探し、

ここで

20

【 0 0 7 3 】

【数 1 1 3】

$$\sigma_f^2 = \sigma_g^2 = 1$$

【 0 0 7 4 】

、の条件付のとき、

【 0 0 7 5 】

【数 1 1 4】

$$r_{fg} = \sigma_{fg}$$

30

【 0 0 7 6 】

を最大にするとき、

【 0 0 7 7 】

【数 1 1 5】

$$r_{fg}$$

【 0 0 7 8 】

を正準相関、

【 0 0 7 9 】

【数 1 1 6】

40

$$f, g$$

【 0 0 8 0 】

を正準変量と呼ぶ。より具体的には、正準相関解析において、 $X$ と $Y$ の特異値分解を行い、

【 0 0 8 1 】

【数 1 1 7】

$$X = U_X D_X V_X' \quad Y = U_Y D_Y V_Y'$$

$$U_X' U_Y = U D V'$$

50

U、D、Vを算出し、そのU、D、Vを用い、

【0082】

【数118】

$$A = \sqrt{n-1} V_X D_X^{-1} U \quad B = \sqrt{n-1} W_Y D_Y^{-1} V$$

$$F = \sqrt{n-1} U_X U \quad G = \sqrt{n-1} U_Y V$$

を求め、

ただし、A、B、F、Gは、

【0083】

【数119】

$$f_i = X a_i \quad g_i = Y b_i \quad i=1, \dots, q$$

$$A = [a_1, \dots, a_q] \quad B = [b_1, \dots, b_q]$$

$$F = [f_1, \dots, f_q] \quad G = [g_1, \dots, g_q]$$

$$F = XA \quad G = YB$$

であり、ここで、i=1から順番に相関の高いもの

【0084】

【数120】

$r_{fg}$

【0085】

を得ることができる。

【0086】

カーネル正準相関分析 (kernel CCA) とは、通常正準相関分析にカーネル法を導入した手法であり、線形モデルに基づく正準相関分析に対して、非線形モデルに基づく相関分析がカーネル正準相関分析では可能である。上記正準相関分析で対象にした第1空間のXと第2空間のYとをそれぞれヒルベルト空間に写像したX'とY'について、正準相関分析を行う方法である (S. Akaho, A kernel method for canonical correlation analysis, International Meeting of Psychometric Society (IMPS), 2001)。

【0087】

サポートベクターマシン (SVM) 法とは、教師付き識別問題を解くための機械学習アルゴリズムである (文献 B. E. Boser, I. M. Guyon, and V. N. Vapnik. A training algorithm for optimal margin classifiers. In D. Haussler, editor, 5th Annual ACM Workshop on COLT, pages 144-152)。SVMでは、データを2種類に分類するために各データ点との距離が最大となる分離平面 (超平面) を求めるマージン最大化という考え方をを用いる特徴を有する。さらに、カーネル関数を用いてパターンを有限もしくは無限次元の特徴空間へ写像し、特徴空間上で線形分離を行う方法を取ることによって非線形分離問題にも優れた性能を示すという特徴も有する。ここで、教師付き分類問題を、タンパク質と化合物の結合予測に適用すると、タンパク質と化合物が結合するというクラスと結合しないというクラスを分類する識別器を作ることになる。この場合、文献や実験などで得られる既知のタンパク質と化合物の結合データを教師データとして用いることができる。

【0088】

10

20

30

40

50

多変量解析手法とは、複数の変数（項目、属性、次元数）を持つデータ（多変量データ）を利用し、その変数間の相互の関係性をとらえるために使われる統計的手法の総称である。重回帰分析や判別分析、正準相関分析、主成分・因子分析、クラスター分析、多次元尺度法、フェース分析、数量化分析、コンジョイント分析などの手法がある。複雑なデータが持つ傾向や特徴を「要約」したり、結果に影響する相関関係を明らかにして「原因発見」や「推定・予測」を行ったり、また因果関係のモデル化などに有効である。CCA、カーネルCCAもまた、この多変量解析に該当する。

【0089】

機械学習法とは、人工知能における研究課題の一つで、人間が自然に行っている学習能力と同様の機能をコンピュータで実現させるための技術・手法のことである。ある程度の数のサンプルデータ集合を対象に解析を行い、そのデータから有用な規則、ルール、知識表現、判断基準などを抽出する。

10

【0090】

1つの実施形態では、本発明は、ある化学物質と相互作用すると予測される化学物質を選出するデータ処理方法を提供する。

【0091】

この方法において用いられている「目的領域」とは、ケミカル空間内に算出される領域を指す。「標的領域」とはバイオ空間内に算出される領域を指す。

【0092】

相互作用の予測対象のうちクエリ側がタンパク質である場合、すなわち標的タンパク質に作用する化合物を予測する場合について詳細に説明する。まず、タンパク質について任意の活性を選択する。標的タンパク質の特性を選択すると、それに対応する空間内の任意の標的領域を算出することができる。

20

【0093】

標的タンパク質又は標的タンパク質群が与えられている場合、その標的タンパク質と配列や構造などが相同なタンパク質群を選出し、それらが占有する空間領域として標的領域を算出することができる。また、タンパク質に定義されている機能（例えば、遺伝子オントロジー（gene ontology）など）に基づいて標的領域を算出する場合は、標的タンパク質又は標的タンパク質群と同等の機能が定義されているタンパク質群を選出し、それらが占有する空間領域として標的領域を算出することができる。

30

【0094】

さらには、タンパク質に定義されている遺伝子発現パターン、パスウェイ位置情報、生物活性情報（例えば、マイクロアレイデータ、反応経路、薬理活性など）に基づいて標的領域を算出する場合は、標的タンパク質又は標的タンパク質群と同等の遺伝子発現パターン、パスウェイ位置情報、生物活性情報が定義されているタンパク質群を選出し、それらが占有する空間領域として標的領域を算出することができる。

【0095】

算出された標的領域から所定の距離以下に存在するケミカル空間の目的領域を算出する工程を説明する。まず標的領域内の各々のタンパク質に対し結合し得る化合物群が特定される。特定された化合物各々について、ケミカル空間内で所定の距離以下に存在する目的領域を算出する。ここでの距離とは、ユークリッド距離、マンハッタン距離などの距離の公理を満たすものから、相関係数やカーネルなどの類似度指標も含む。

40

【0096】

ケミカル空間の目的領域に存在する化合物が選出される工程では、算出された目的領域に対応する化合物が選出される。これは、目的領域が算出されていれば、自動的な計算によって選出することも可能である。

【0097】

また、他の実施形態では、相互作用の予測対象のうちクエリ側が化合物である場合、すなわち目的の化合物に作用するタンパク質を予測する方法を提供する。

【0098】

50

1つの実施形態では、本発明は、サンプルデータを用いて第1空間と第2空間とを相関させるようトレーニングすること工程をさらに包含する。本明細書において「トレーニング」とは、装置の使用のための訓練に要する計算機操作で、取付け操作、操作卓操作、変換操作、印刷操作のような活動や、必要なデモンストレーションを行なうのに使われる操作をいう。本明細書において「トレーニングデータ」とは、操作の始めにロボットのコンピュータへ入力される練習データをいう。

【0099】

1つの実施形態において、トレーニングは、直交行列  $A = C_{x \times x}^{-1/2} U$  と  $B = C_{y \times y}^{-1/2} V$  を生成する（ここで、 $\det(A) = \det(B) = 1$  かつ数131で示される数式の通りである）。

【0100】

【数131】

$$C_{xx} = E[(X - m_x)(X - m_x)^T], C_{yy} = E[(Y - m_y)(Y - m_y)^T], C_{xy} = E[(X - m_x)(Y - m_y)^T], K = C_{xx}^{-1/2} \cdot C_{xy} \cdot C_{yy}^{-1/2} = U \cdot S \cdot V^T$$

【0101】

第1モダリティの第1空間を表す  $A X$  と第2モダリティの第2空間を表す  $B Y$  との間の相関は最大となり、これにより、該第1モダリティから該第2モダリティへの特徴の移転が可能となることを特徴としてもよい。

【0102】

別の実施形態では、前記トレーニングにより、行列  $A$  と行列  $B$  を生成する。第1モダリティの第1空間を表わす  $X A$  と第2モダリティの第1空間を表わす  $Y B$  との間の相関は最大となり、これにより、第1モダリティから第2モダリティへの特徴の移転が可能となる。特徴の移転は、行列  $X$ 、 $Y$  の行には化学物質のエントリが、列には化学物質情報が並ぶ2種の異なるデータ（例えば化合物とタンパク質）を行列  $X$ 、 $Y$ （第1空間が行列  $X$ 、第2空間が行列  $Y$ ）と表現し、正準相関解析を行って両空間の相関を最大にすることができる。

【0103】

本明細書において「モダリティ」とは、特徴的属性をいう。1つの実施形態において、第1空間を表す  $A X$  のクエリは、第2空間を表す  $B Y$  の前記クエリの結果のみが与えられると、 $B Y$  は  $A X$  と最大の相関を有することから特定可能である。

【0104】

また、本発明は、第1の化学物質の第1の特徴量と第2の化学物質の第2の特徴量とを多変量解析手法又は機械学習法によって1つの特徴空間へ写像変換するデータ処理方法を提供する。機械学習法としては  $SVM$  法を用いることができる。

【0105】

また、具体的な実施形態では、機械学習法において、1) 相互作用予測の対象となる第1化学物質の特徴量と第2化学物質の特徴量とを、第1の化学物質の第1の特徴量の空間座標のデータベースおよび第2の化学物質の第2の特徴量の空間座標のデータベースによって構築された特徴空間にマッピングする工程；2) 該問い合わせペアが、空間エリア内に存在する場合、第1化学物質と第2化学物質とが結合すると判定し、空間エリア内に存在しない場合、第1化学物質と第2化学物質とが結合しないと判定する工程を包含する。

【0106】

特徴空間にマッピングする特徴量は、第1の化学物質の特徴量と第2の化学物質の特徴量とを連結した特徴量であってもよい。また、第1の化学物質の特徴量及び第2の化学物質の特徴量はベクトルで表現されるが、連結された特徴量はベクトル同士の内積であるカーネルを含む。

【0107】

1つの実施形態では、本発明の方法は、さらに、選出された化学物質をインシリコで生産する工程を包含する。インシリコでの生産方法は、本明細書において別の場所において

10

20

30

40

50



記載されており、周知の技術を用いることができる。

【0108】

別の実施形態では、本発明の方法は、さらに、ウェットで選出された化学物質を生産する工程を包含する。ウェットでの生産方法は、本明細書において別の場所において記載されており、周知の技術を用いることができる。ウェットでの生産の代表例としては、コンビナトリアル・ケミストリを用いることがあり得る。ウェットでの生産は、遺伝子組み換え技術を用いて達成されてもよい。

【0109】

1つの具体的な実施形態では、本発明の方法は、さらに、前記第1空間の化学物質の選出の後、該第1空間の化学物質の前記第2の特徴量を測定して、実際に所望の活性を有する化学物質を選出する工程をさらに包含する。

10

【0110】

本発明において、第1の化学物質と第2の化学物質とが相互作用するか否かが判定される工程には、相互作用する確率を算出する工程が含まれる。また、相互作用する確率はスコアとして算出されてもよい。ここでいうスコアとは、第1空間内の任意の化学物質Aと第2空間内の任意の化学物質Bの結合のしやすさ（結合予測の統計的有意性）を示すものである。

【0111】

例えば、化学物質Aと化学物質Bとの結合スコアは次のように定義できる。第2空間内の全ての化学物質数をN、そのうち化学物質Aと結合する個数がLあったとする。ここで、Bに近接する第2空間内の化学物質K個を考えた場合、そのうち化学物質Aと結合する個数がHであったとする。その際の化学物質AからBの結合スコアは $\log(H/K)/(L/N)$ というようなオッズスコアとして定義できる。また逆に、第1空間内の全ての化学物質数をn、そのうち化学物質Bと結合する個数がlあったとする。ここで、Aに近接する第1空間内の化学物質k個を考えた場合、そのうち化学物質Bと結合する個数がhであったとする。その際の化学物質BからAの結合スコアは $\log(h/k)/(l/n)$ というようなオッズスコアとして定義できる。従って、化学物質AからBへのスコアとBからAへのスコアから、化学物質AとBの総合スコアを $\log(H/K)/(L/N) + \log(h/k)/(l/n)$ と定義することができる。これらのスコアはCCAやkernel CCAで算出された第1空間と第2空間の相関モデルから算出される。また、SVM法を用いた場合は、第1空間内の化学物質群と第2空間内の化学物質群との間の既知の結合ペアと結合しないペアを分離する超平面からの距離に相当するものからスコアを換算することができる。

20

30

【0112】

「相互作用情報」には、解離定数Kd、50%阻害効果濃度IC50、50%亢進効果濃度EC50などが挙げられる。医薬品開発の場合、結合の有無、結合活性、薬理活性の基準として、Kd、IC50、EC50が、マイクロモルオーダー、ナノモルオーダーであると望ましい。また、相互作用する確率のスコアは、特徴空間内でクラス分類（例えば、結合する化合物とタンパク質ペアのクラスと結合しない化合物とタンパク質ペアのクラスなど）をする境界面（超平面）からの予測対象までの距離で表される。境界面から遠距離にあるほど、相互作用する確率が高くなる。

40

【0113】

本発明は、本発明のデータ処理方法により得られた結果から設計されたライブラリを提供する。また、本発明は、本発明の方法を用いることによって生産された化学物質を提供する。

【0114】

また、本発明は、本発明のデータ処理方法をコンピュータに実行させるプログラムを提供する。本発明のプログラムは、コンピュータ読み出し可能な記録媒体に格納される。記録媒体としては、プログラムを記録することができる限り、任意の形態（例えば、フレキシブルディスク、MO、CD-ROM、CD-R、DVD-ROMのような任意のタイプ

50

)を使用することができることが理解される。さらに、本発明のデータ処理方法に用いられるデータ構造物は、記録媒体に格納される。具体的には、第1の化学物質の空間座標のデータベースにより定義される第1空間と第2の化学物質の空間座標のデータベースにより定義される第2空間とを備えるデータ構造物や、本発明の方法により構築された特徴空間を備えるデータ構造物などが含まれる。

【0115】

本発明は、特定の標的タンパク質に活性を持つ化合物を探索する目的に適用可能なほか、特定の化合物が与えられたときにその化合物に作用される複数のタンパク質の推定にもつながり、薬物の副作用に関する知見を提供する。また、特定の化合物に作用できる人工タンパク質（遺伝子改変タンパク質）の創製などへの応用も可能であると考えられる。

10

【0116】

また、本発明のデータ処理方法を用いて第1の化学物質と第2の化学物質との間の相互作用を予測するデータ処理装置が提供される。また、別の実施形態において、特定の化学物質と相互作用すると予測される化学物質を選出するデータ処理装置が提供される。さらに、別の実施形態において化学物質ライブラリを設計するデータ処理装置が提供される。これらの処理装置は、その処理装置に備え付けられた演算装置上で本発明の方法が実行されることにより実施される。

【0117】

本明細書において引用された、科学文献、特許、特許出願などの参考文献は、その全体が、各々具体的に記載されたのと同じ程度に本明細書において参考として援用される。

20

【0118】

以上、本発明を、理解の容易のために好ましい実施形態を示して説明してきた。以下に、実施例に基づいて本発明を説明するが、上述の説明および以下の実施例は、例示の目的のみに提供され、本発明を限定する目的で提供したのではない。従って、本発明の範囲は、本明細書に具体的に記載された実施形態にも実施例にも限定されず、特許請求の範囲によってのみ限定される。

【実施例】

【0119】

以下に実施例を示して本発明をさらに詳しく説明するが、この発明は以下の例に限定されるものではない。

30

【0120】

(実施例1：CCAを用いたバイオ空間およびケミカル空間でのスクリーニング)

従来法(PCA)と本手法(CCA)の性能評価をするために、それぞれで構築した「ケミカル空間とバイオ空間の融合モデル」を用いてインシリコスクリーニング(In silico screening)を行った。

【0121】

図6は、予測性能を評価する方法の一つであるROC曲線である。このグラフは曲線が上に位置すれば位置するほど、予測性能が良いことを表すものであり、本手法の曲線が従来法の曲線より上に位置することから、本手法の予測性能が、従来法の予測性能に比べて高いことがわかる。以下にその具体的手順を示す。

40

【0122】

(プロトコール)

1、融合モデルの構築に用いる既知のタンパク質と化合物のデータは、Drug Bank データベース(<http://redpoll.pharmacy.ualberta.ca/drugbank/>)2005年8月リリース版から取得した。

2、全ての化合物エントリのmol fileについて、DragonXソフトウェアを用いて、937個の化合物記述子を算出した。ここで、計算された化合物数は3079個である。さらに、CCA計算を行うにあたり、属性となる記述子のプロファイルは独立していなければならないため、相関係数0.8以上の相関性を持つ記述子は、情報量の高い記述子300個に縮約した。

50

3、全てのタンパク質エントリの `fasta file` について、`mismatch string kernel` を生成する手法と同様の手法により、ミスマッチを考慮した連続する2アミノ酸の組成比からなる400次元(アミノ酸20種\*20種)のプロファイルを算出した。ここで、プロファイル化されたタンパク質数は3476個である。また、上記2で算出された化合物との結合数は8006個であった。

4、予測性能の評価には、5分割交差検証法(5 fold cross validation)を用いた。すなわち、上記2と3で作成した8006個の化合物-タンパク質結合データを無作為に4:1に分類し、80%の結合データをトレーニングデータとしてCCA計算やPCA計算をし、バイオ空間とケミカル空間のそれぞれの座標を構築した。残りの20%の結合データはテストデータとして用いた。ここで、テストデータの負例(結合しないデータ)は、正例である結合データを構成する化合物とタンパク質の組合せで結合しない組合せを発生させ、正例と同数無作為に選出した。このように作成したテストデータを、トレーニングデータによって構築したバイオ空間とケミカル空間とにそれぞれマッピングした。マッピングとは、トレーニングデータをCCAやPCA計算することによって、算出される重み係数行列(PCAの場合は主成分得点係数行列)をテストデータ行列にかけることによってなされる。バイオ空間とケミカル空間とにそれぞれマッピングされたタンパク質と化合物について、各々についてスコアを算出した。

5、上記4のようにテストデータを予測したとき、実際に結合するデータを結合すると予測できたものの比率を真陽性率、実際には結合しないデータを結合すると予測してしまったものの比率を偽陽性率と呼ぶ。ここで、特定の予測スコア(閾値)以上の値を持つデータは陽性とみなし、特定スコア以下の値を持つデータは陰性とみなす。上記CCAとPCAにおいて予測した化合物-タンパク質結合スコアに基づいて、スコアの閾値を動かし、それに伴う偽陽性率と真陽性率を(x, y)としてプロットし、ROC曲線を作成した(図6)。

#### 【0123】

(実施例2:化合物-GPCR相互作用予測手法)

化合物-GPCR相互作用予測手法の開発を以下の手順で行った。なお、本実施例にあたり用いたデータセットおよび解析方法等に関する詳細を記載した文献等については、本文中に参照番号を付し、本実施例の末尾にその参考文献一覧を添付した。これらの文献は、本明細書中で参考として援用される。

#### 【0124】

(#1 化合物-タンパク質相互作用情報の収集)

GLIDA (GPCR-Ligand Database) [1]、DrugBank [2]、IUPHAR Receptor database [3]、PDSP Kidatabase [4] から、相互作用する化合物-GPCRの組み合わせ5207例(化合物:866、GPCR:317)を収集した。ただし、ここではヒト、マウス、ラットのGPCRを用い、GPCRの定義はGPCRDB [5] に従った。また、化合物については、続く記述子(descriptor)の計算に構造情報が必要であるため、mol(sdf)形式のファイルが提供されているGLIDAおよびPubChem Compound [6] に登録されている化合物を用いた。

#### 【0125】

(#2 記述子の計算)

化合物およびタンパク質を特徴ベクトルとして表現するために、以下の方法によりそれぞれの記述子を計算した。

・化学記述子(chemical descriptor)

収集した化合物の構造から、化合物の構造・物性に関する記述子をDRAGONX ver. 1.2 [7] により計算した。この研究では、カテゴリー1-10(constitutional descriptors、topological descriptors、walk and path counts、connectivity indices、information indices、2D autocorrel

10

20

30

40

50

ations、edge adjacency indices、Burden eigenvalue descriptors、topological charge indicesおよびeigenvalue-based indices)、カテゴリー17-18(functional group countsおよびatom-centered fragments)、カテゴリー20(molecular properties)の計929記述子を計算した。なお、分子の三次元座標に依存する記述子(カテゴリー11-14)、官能基や原子タイプの数を取る記述子(カテゴリー15および16)、電荷記述子(カテゴリー19)は、ここでは用いなかった。続いて、これらの記述子のうち、すべての化合物で同一の値として計算出力されるものを取り除き、結果として残った797種類の記述子を以下で用いた。

10

・タンパク質記述子(protein descriptor)

ミスマッチを許容したスペクトラム法[8]により計算した。この方法は、タンパク質配列を固定長kのアミノ酸配列に分解し、この中に現れる、最大m個のミスマッチまで許容した長さkのアミノ酸配列パターンの頻度を数えることにより計算される。発明者らは(k, m)を(2, 1)に設定した。したがって、計算される記述子は、1アミノ酸のミスマッチを許容した2連アミノ酸202種類となる。

【0126】

(#3 サポートベクターマシン(SVM)による学習モデルの構築)

SVMは、Vapnikら[9]により提案された学習アルゴリズムであり、その高い汎化能力から各方面において多用されている。SVMは、2つの異なるグループの特徴ベクトルを最大マージンで分離するような超平面を構築する。ここで、最大マージンとは、分離した超平面から各サンプル間までの最短距離を指す。

20

【0127】

発明者らは、化合物-タンパク質相互作用の有無を分離する超平面を求めるために、正例(相互作用するパターン)および負例(相互作用しないパターン)に対応する化学記述子、タンパク質記述子をそれぞれ組み合わせて特徴ベクトルを構築し、SVMを用いて学習モデルを構築した。ただし、負例については、相互作用しないパターンの情報が得られないため、2つの記述子をランダムに組み合わせて正例と同数を生成した。ここで、SVMライブラリとして、「Sequential Minimal Optimization」アルゴリズム[10, 11]を採用しているlibsvmプログラム[12]のコードを用いた。SVMモデルが得られると、新しいベクトル(化合物-タンパク質ペア)が、相互作用有/相互作用無のどちらのクラスに属するかを予測することができる。さらに、判別だけでなく、サンプルのスコアリングを行う方法も報告されている[13]。これは、分離面に近いサンプルは、分離面から遠いサンプルよりも、誤って分類される確率が高いであろうという考えに基づいている。発明者らは、化合物-タンパク質間の相互作用予測において、この方法により、その可能性のスコア化、および順位付けを行った。

30

【0128】

(#4 化合物構造類似性によるリガンド予測)

発明者らは、モデルの比較対象となるリガンド予測方法として、化学記述子から計算される化合物の類似性を用いた。この類似性は、一般的な化合物探索の方法であり、リード化合物を発見する手助けになると言われている[14]。この研究では、上述の「#2」で計算した797種類の化学記述子を主成分分析し、主成分座標上で既知リガンドと隣接する化合物から順にスコア付けを行った。化合物A-タンパク質Bペアのスコアは、主成分空間において化合物Aから見て最近傍にあるタンパク質B既知リガンドとの類似度で表現される。主成分は、累積寄与率80%までのもの(30主成分)を用いた。また、類似度の尺度として相関係数(Pearson correlation coefficient)を用いた。

40

【0129】

(#5 交差検証法(cross validation)によるモデルの評価)

n分割交差検証法(n-fold cross-validation)を用いて学習

50

モデルの予測性能を評価した。この評価法では、最初に全学習データセットが  $n$  個の等サイズなサブセットに分割される。続いて、それぞれのサブセットについて、残りの  $n - 1$  サブセットで学習して作られた分類器を用いて予測する。そして、この操作は、すべてのサブセットが一度だけ予測されるように繰り返されて評価される。予測性能の尺度としては、以下の式で計算される正確度 (Accuracy) を用いた。正確度は以下の式で表される。

【0130】

$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$$

ここで、TP は真陽性、TN は真陰性、FP は偽陽性、FN は偽陰性を表す。

【0131】

ランダムな組み合わせで作られる負例によるスコア変動を考慮し、負例を交換しながら 10 回の異なるデータセットを生成して 5 分割交差検証法 (5-fold cross-validation) を繰り返し行い、正確度の平均値により発明者らのモデルを評価した。続いて、計算された相互作用予測スコアから ROC 分析を行った。ここで、各評価において、化合物の構造類似性に基づいたリガンド予測法 (#4) を比較対象とした。

【0132】

(#6 ヒト  $\beta$ 2 アドレナリン受容体 ( $\beta$ 2 AR) のリガンド予測)

発明者らが収集した化合物 - GPCR 相互作用情報は、今までの研究により「強く結合する」と知られているもののみであり、その他の大部分の化合物 - GPCR 相互作用は不明である。発明者らの疑問は、リガンド探索において、予想に反して相互作用すると予測された化合物が本当に相互作用しないかどうかということである。そこで、発明者らは、インビトロ (in vitro) 結合阻害実験により、相互作用予測スコアと相互作用の有無との関連性を確認した。

【0133】

そのために、ヒト  $\beta$ 2 AR を標的タンパク質とし、作成した学習モデルを用いてリガンド予測を行った。この受容体は、喘息治療の標的として治療薬の開発が進められている生理学的に重要な遺伝子である。リガンド予測の対象化合物は、GPCR との相互作用が知られている上記 866 化合物 (ただし、モデル構築時にヒト  $\beta$ 2 AR との相互作用を学習した化合物は除く) とした。これらの化合物の化学記述子に対して  $\beta$ 2 AR のタンパク質記述子 (protein descriptor) を組み合わせ、予測用データセットとした。負例組み合わせによるスコア変動を考慮し、負例を交換しながら学習と予測の試行を 30 回繰り返し、各リガンドについて、得られたスコアの最大値を最終的な予測スコアとした。

【0134】

次に、相互作用予測スコア上位 50 (Top 50) の化合物について、さらなる調査・実験を行った。まず、文献・特許調査 (SciFinder、PubMed) により、 $\beta$ 2 AR との相互作用に関する報告が存在しないか確認した。

【0135】

続いて、相互作用情報を確認できなかった化合物のうち、入手可能な化合物について、インビトロ (in vitro) 結合阻害実験による検証を行った。この実験では、ヒト  $\beta$ 2 AR 強制発現細胞株から膜画分を調製し、放射性  $\beta$ 2 AR リガンドである [ $^{125}\text{I}$ ] - シアノピンドロールに対する競合的な阻害効果を確認した。

【0136】

ところで、相互作用しないという情報の欠如により、発明者らのモデルでは、ランダムに発生させた化合物 - タンパク質ペアを負例 (相互作用なし) パターンとして採用している。このため、相互作用予測スコアの低い化合物が本当に相互作用しないか、ということを確認する必要がある。そこで、発明者らは、予測スコア下位 50 (Bottom 50) の化合物についても、上位 50 (Top 50) と同様の文献調査・検証実験を行った。

【0137】

10

20

30

40

50

(結果)

(交差検証法による新規リガンド予測モデルの評価)

まず、手始めに今回開発した方法と従来法との比較検討を行った。公共データベースから収集したGPCR-リガンド相互作用情報を用い、化合物-タンパク質相互作用パターンの特徴ベクトルをSVM分類器の入力とし、学習モデルを構築した。負例を交換しながら5分割交差検証法を10回試行した結果、発明者らが開発したモデルの予測性能(accuracy)は91.3%±0.3%だった。対照として、化合物類似度に基づいた従来法についても同様に5分割検証法を行ったところ、予測性能は81.9±0.3%だった。また、ROC曲線からも、発明者らの開発したモデルの予測性能が高いことが判明した(図7)。

10

【0138】

(ヒト2ARリガンド予測)

次に、新規手法をヒト2ARの新規リガンド予測に適用し、その有効性を実験により検証した。また、新規手法でのみ予測されるリガンドが従来法では検出できないような新規骨格を持つ化合物を含むかどうか調べた。構築したモデルを用いて、866種類のGPCRリガンドについて2ARとの相互作用予測スコアを算出した。

【0139】

新規モデルが予測した2ARリガンド候補Top50の化合物のうち、文献・特許調査により14種の化合物が2ARとの相互作用に関する報告を確認した(図8(B-1)左)。さらに、残りの相互作用不明な化合物のうち、入手可能な21種類についてインビトロ(in vitro)結合阻害実験を行ったところ、17種類の化合物が相互作用( $10^{-5} \text{ M} < \text{IC}_{50} < 10^{-3} \text{ M}$ )を示した(図8(B-1)右)。実験のヒット率は81%(17/21)にのぼり、ここにおいても高い予測的中率が示された。

20

【0140】

一方、下位50(Bottom 50)の化合物については、2ARリガンドとして報告されているものは文献および特許調査では確認されなかった(図8(B-2)左)。さらに、残りのうち入手可能な9化合物についてインビトロ(in vitro)結合実験を行ったところ、2個の化合物が同程度の強さの相互作用を示したが、残りの7化合物は相互作用を示さなかった(図8(B-1)右)。

【0141】

これらの予測結果を従来法によるものと比較した図が図9である。実験で相互作用を確認した化合物の半数近くは、化合物の構造類似性に基づく従来の方法ではスコアが低かった。実際に、これらの化合物は、典型的な2AR作動薬の構造(カテコラミン骨格、イソプレナリン誘導体)および2AR拮抗薬の構造(アリルアルキルアミン誘導体)とは異なる多様な骨格(図9左)を持っており、化合物の構造類似性に基づく従来の方法では発見できないリガンド群であるといえる。すなわち、相互作用情報に基づく新しいモデルは、多様な構造を持つ化合物が同一タンパク質に作用するという関係を正しく予測することができたといえる。また、これらの化合物の中には、ニューロペプチド受容体アンタゴニストなど、従来はペプチド受容体に作用する化合物として知られていたものも含まれていたが、遠縁にあたる2ARとも相互作用することが実験により確認された。

30

40

【0142】

(実施例3:化合物ライブラリの設計)

本手法を用いて、化合物ライブラリの設計、すなわち特定の化合物に対する標的遺伝子の予測を行った。図10に設計手法を示す。

【0143】

予測は、米国NCBI/PubChemデータベース内の化合物6,391,005件を用いて、それらの化合物が標的とし得るタンパク候補を予測し、化合物ライブラリを構築した。ここで、予測の基準座標となるケミカル空間とバイオ空間との融合モデルの作成には、薬物とその標的タンパク質のデータを蓄積したカナダのDrugBankデータベースを用いた。

50

## 【0144】

図11は、PubChem化合物の生物活性予測の結果である。各列は化合物と標的タンパクの結合可能性の信頼性を表すスコアごとに分かれている。すなわちスコアが高ければ高いほど、その化合物の生物活性の信頼性は高いと考えられる。ここで言うスコアは、第1空間内の任意の化学物質Aと第2空間内の任意の化学物質Bの結合のしやすさ（結合予測の統計的有意性）を示すものと考えられる。

## 【0145】

例えば、化学物質Aと化学物質Bとの結合スコアは次のように定義できる。第2空間内の全ての化学物質数をN、そのうち化学物質Aと結合する個数がLあったとする。ここで、Bに近接する第2空間内の化学物質K個を考えた場合、そのうち化学物質Aと結合する個数がHであったとする。その際の化学物質AからBの結合スコアは $\log(H/K)/(L/N)$ というようなオッズスコアとして定義できる。また逆に、第1空間内の全ての化学物質数をn、そのうち化学物質Bと結合する個数がlあったとする。ここで、Aに近接する第1空間内の化学物質k個を考えた場合、そのうち化学物質Bと結合する個数がhであったとする。その際の化学物質BからAの結合スコアは $\log(h/k)/(l/n)$ というようなオッズスコアとして定義できる。図11では化学物質AからBへのスコアとBからAへのスコアとから、化学物質AとBの結合可能性の総合スコアを $\log((H/K)/(L/N)) + \log((h/k)/(l/n)) + 20$ と定義した。

## 【0146】

また、各行の項目は標的タンパクの機能（遺伝子オントロジー（gene ontology）に基づく）ごとの分類を表している。表中の数値は、該当する部分に対応する（予測された）化合物の数である。例えば、receptor activityに関するタンパクを標的とし、スコア値27以上の信頼性を示す化合物は、198個予測されたことになる。

## 【0147】

図12では、図11と同様に、PubChem化合物の生物活性予測の結果であるが、標的タンパクの機能分類を異なる基準で行ったものである。図12の見方は図11と同様である。

## 【図面の簡単な説明】

## 【0148】

【図1】本発明の概念図である。

【図2】本発明の別の概念図である。

【図3】本発明の別の概念図である。

【図4】本発明の別の概念図である。

【図5】ケミカル空間とバイオ空間とを融合したモデルの構築方法についての概念図である。

【図6】従来法（PCA）と本手法との性能評価のため、DrugBankデータベースから取得したデータを基に、インシリコスクリーニング（In silico screening）を行った結果のROC曲線である。

【図7】従来法（化合物構造類似性による相互作用予測法）と本手法との性能評価のため、収集した相互作用する化合物GPCRのデータを基に、相互作用予測を行った結果のROC曲線である。

【図8】2ARリガンド予測結果の検証の結果を示す。予測スコア上位50位（Top 50）（B-1）と下位50（Bottom 50）（B-2）の化合物に対する調査・実験結果を示す。それぞれ、左側が、文献調査および実験検証で判明した化合物の内訳であり、右側が、 $[^{125}I]$ -シアノピンドロールに対する結合阻害曲線、縦軸が阻害された割合、横軸が各化合物の濃度を示す。

【図9】新規モデルおよび従来法による2ARリガンド予測結果の比較を示す。各点が化合物であり、縦軸は新規モデル、横軸は従来手法による相互作用予測スコアを示している。

10

20

30

40

50

【図10】本発明の化合物ライブラリ設計の概念図である。

【図11】PubChem化合物の生物活性予測の結果である。

【図12】標的タンパクの機能分類を異なる基準で行ったPubChem化合物の生物活性予測の結果である。

【図13】本発明におけるCCAを適用した際の計算フローである。

【図14】本発明におけるSVMを適用した際の計算フローである。

【表1】

【参考文献】

- [1] Okuno Y, Yang J, Taneishi K, Yabuuchi H & Tsujimoto G. GLIDA: GPCR - ligand database for chemical genomic drug discovery. *Nucleic Acids Res.* 34,D673 - 7(2006).
- [2] Fredholm B.B., Fleming W.W., Vanhoutte P.M. & Godfraind T. The role of pharmacology in drug discovery. *Nat. Rev. Drug Discov.* 1,237 - 8(2002).
- [3] Wishart D.S., Knox C, Guo A.C., Shrivastava S, Hassanali M, Stothard P, Chang Z & Woolsey J. DrugBank: a comprehensive resource for in silico drug discovery and exploration. *Nucleic Acids Res.* 34,D668 - 72(2006).
- [4] Roth B.L., Kroeze W.K., Patel S & Lopez E. The Multiplicity of Serotonin Receptors: Uselessly diverse molecules or an embarrassment of riches? *The Neuroscientist* 6,252 - 262(2000).
- [5] Horn F, Bettler E, Oliveira L, Campagne F, Cohen FE, & Vriend G. GPCRDB information system for G protein - coupled receptors. *Nucleic Acids Res.* 31, 254 - 7(2003).
- [6] Wheeler D.L., Barrett T, Benson D.A., Bryant S.H., Canese K, Chetverin V, Church D.M., et al. Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res.* 34,D173 - 80(2006).
- [7] DRAGON software is available at [http://www.taletc.mi.it/main\\_net.htm](http://www.taletc.mi.it/main_net.htm)
- [8] Leslie C.S., Eskin E, Cohen A, Weston J & Noble WS. Mismatch string kernels for discriminative protein classification. *Bioinformatics* 20,467 - 76(2004).
- [9] Vapnik V.N. *The Nature of Statistical Learning Theory* Springer: New York, (1995).
- [10] Platt J. Fast Training of Support Vector Machines using Sequential Minimal Optimization. Microsoft Research Technical Report MSRTR - 98 - 14(1998).
- [11] Keerthi S.S., Shevade S.K., Bhattacharyya C, & Murthy R.R.K. Improvements to Platt's SMO Algorithm for SVM Classifier Design. *Neural Comput.* 13,637 - 649(2001).
- [12] Chang C.C. & Lin C.J. LIBSVM: a library for support vector machines. Software is available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>
- [13] Platt J. Probabilistic outputs for support vector machines and comparison to regularized likelihood methods. *Advances in Large Margin Classifiers* (pp.61 - 74). (1999)
- [14] Oprea T.I. Chemical space navigation in lead discovery. *Curr. Opin. Chem. Biol.* 6,384 - 9(2002)

10

20

【産業上の利用可能性】

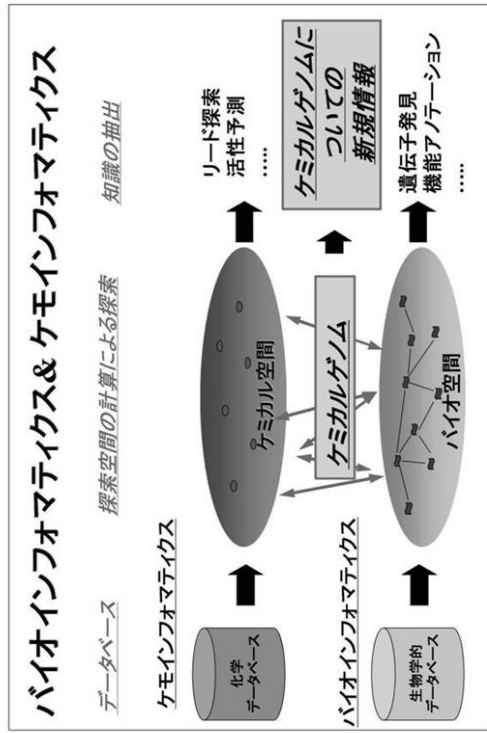
【0149】

本技術によって、特に創薬の分野では新薬開発のコストを大幅に下げ、また研究開発サイクルも短縮することができる。これにより、従来よりも短い期間でより良い薬品を市場に送り出すことができる。また、製薬コストに占める研究開発費の割合を下げることで、社会的には医療費負担の低減という貢献が期待できる。

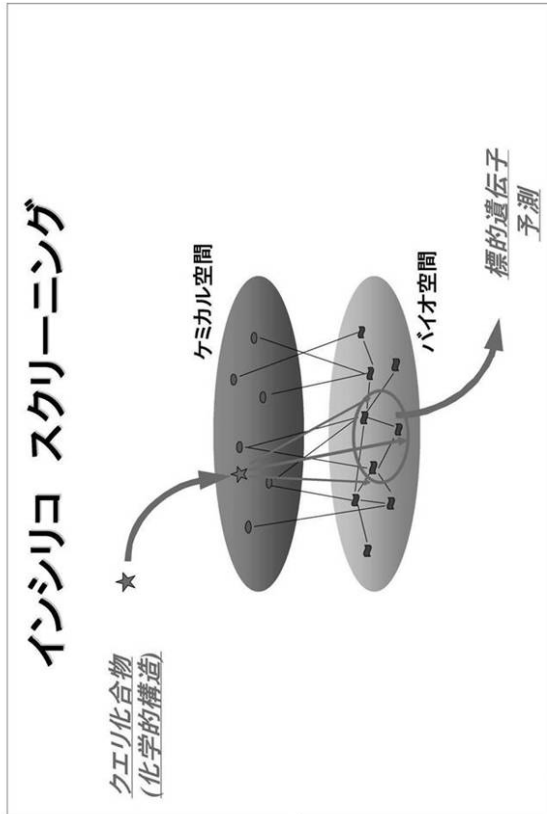
30



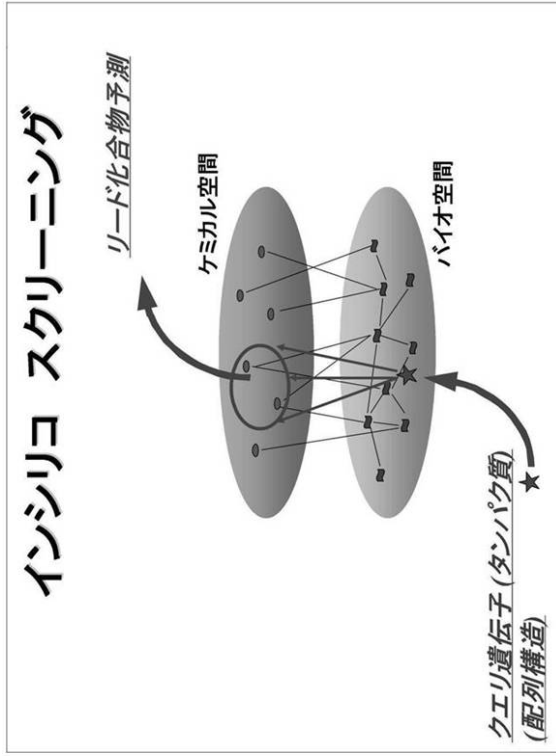
【 図 1 】



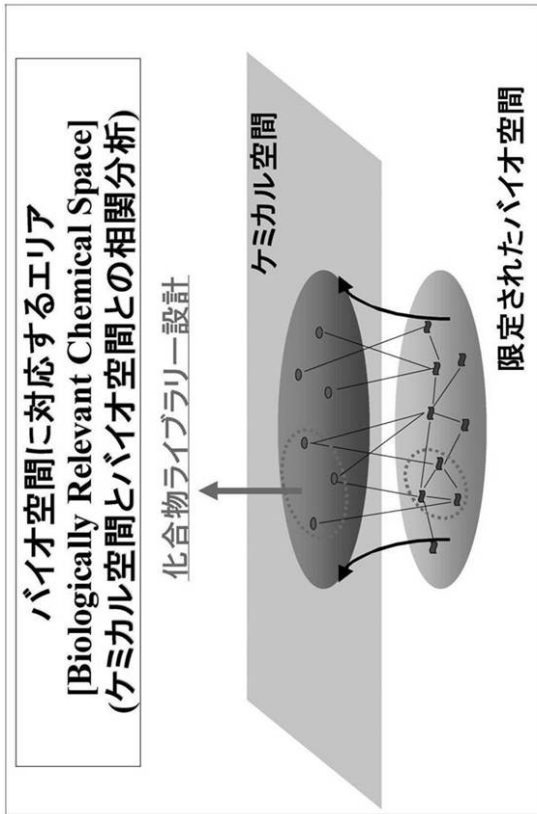
【 図 2 】



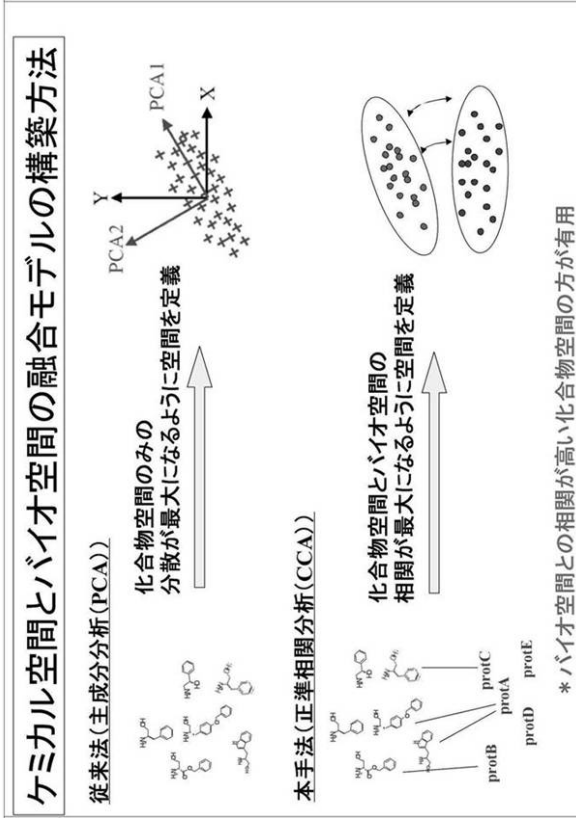
【 図 3 】



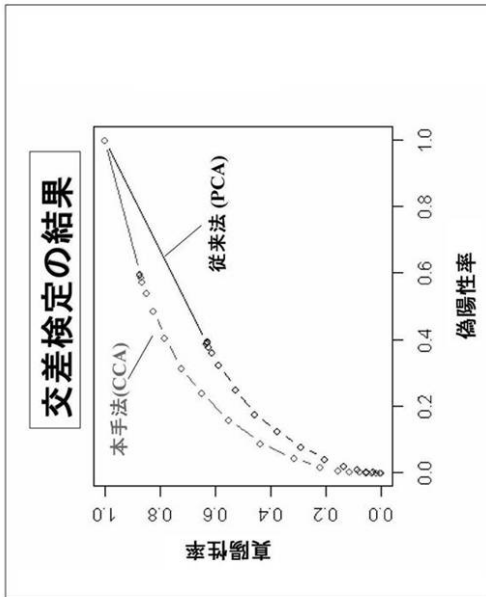
【 図 4 】



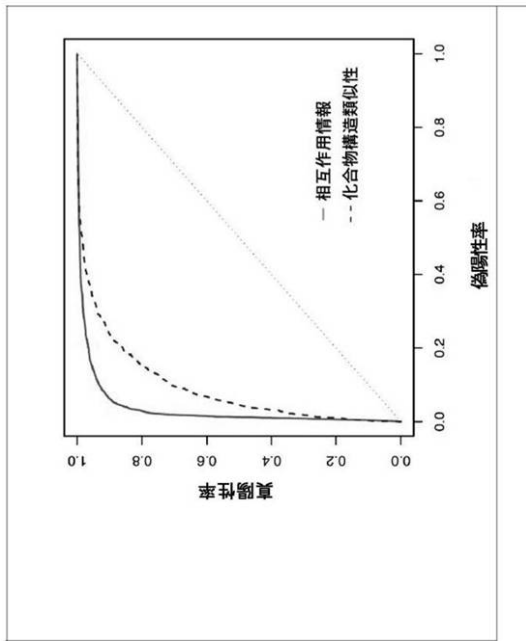
【 図 5 】



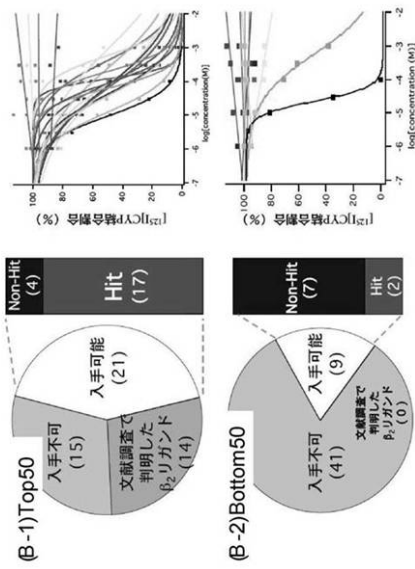
【 図 6 】



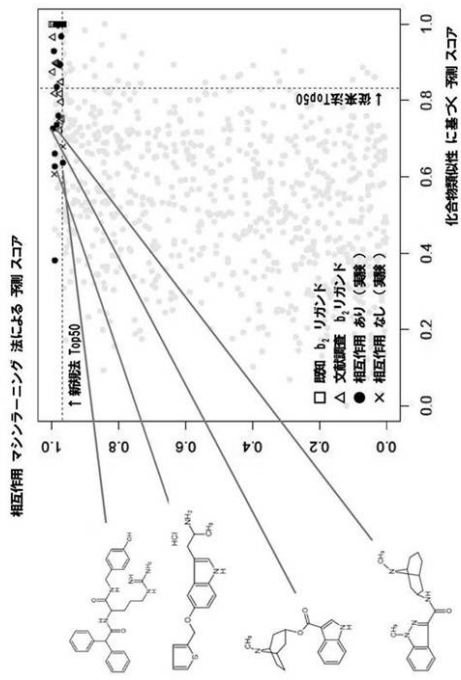
【 図 7 】



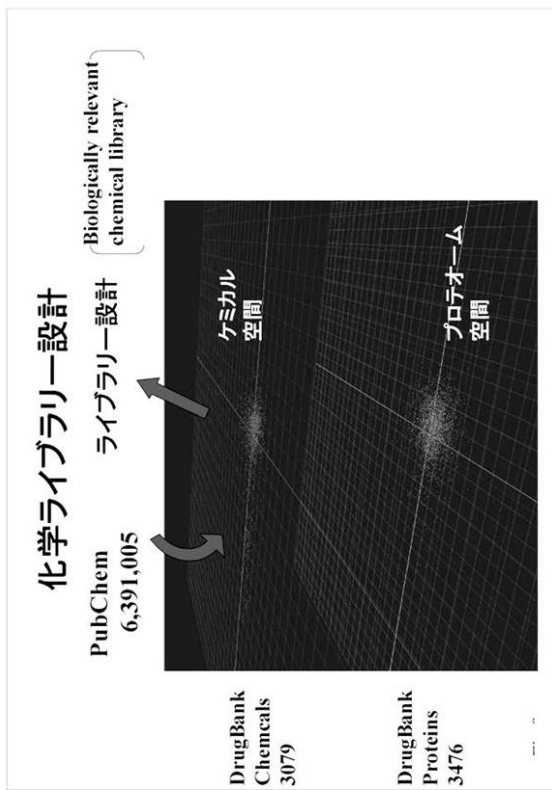
【 図 8 】



【 図 9 】



【 図 10 】



【 ☒ 1 1 1 】

ライブラリー例: PubChem 化合物の生物活性アノテーション G.O. Function

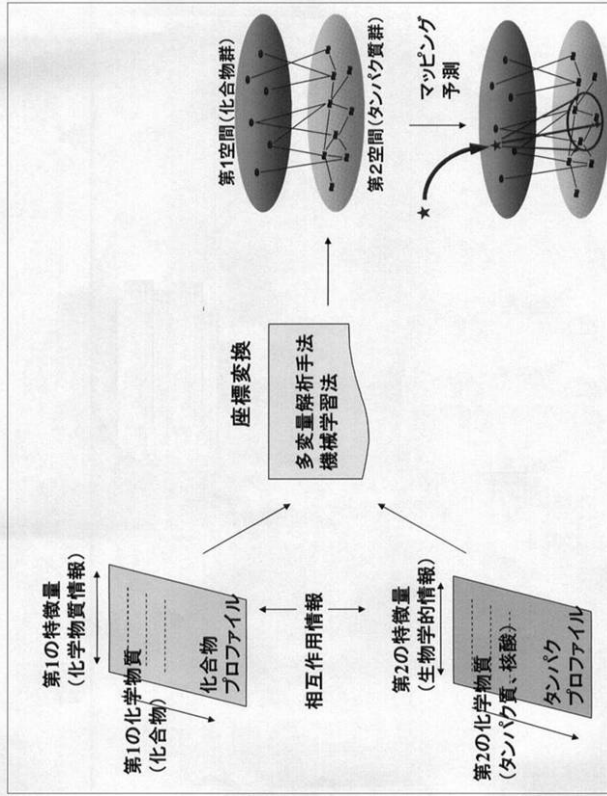
	score 27	score 26	score 25	score 24	score 23	score 22	score 21
motor activity	0	0	0	12659	56234	146521	30441
catalytic activity	0	132	7498	91972	809264	2220698	806438
helicase activity	198	94476	633290	1512369	2068061	2517613	44202
signal transducer activity	198	98111	675597	1623305	2387872	3033699	726701
receptor activity	0	119	3285	24141	233610	587076	147335
structural molecule activity	0	17186	424389	1075678	1577678	2119350	468872
transporter activity	0	12916	325688	804791	1032175	1336641	223948
carrier activity	77	52722	480419	1599411	3075508	3886455	2047957
binding	0	128	310	18706	348653	904792	220589
electron transporter activity	0	28167	505062	1217307	1779455	2617458	752057
molecular_functionunknown	0	0	0	1	25707	74223	275611
protein binding	0	0	0	123	7033	39455	6887
ion channel activity	0	0	0	0	333	79719	77269
ion transporter activity	0	119	2865	11430	45211	153815	42182
ion channel or pore class transporter activity	0	15869	330974	778901	957776	1167142	180847
channel or pore class transporter activity	0	813	24053	500237	1376542	2203510	393771
antioxidant activity	0	0	0	6788	21683	96815	13285
oxidoreductase activity	0	2913	14793	324305	1419636	2695267	909689
transferase activity	0	672	10776	149433	951290	2245088	979300
hydrolase activity	0	12419	98644	359396	1276409	2848843	1203689
lyase activity	77	10685	17390	51148	568656	1351063	379692
isomerase activity	0	0	193	15063	234557	643863	219701
ligase activity	0	0	6	19484	316544	1100592	330559
enzyme regulator activity	0	1114	25007	56681	262072	576435	116670
transcription regulator activity	0	3822	36823	92885	248734	857861	19525
translation regulator activity	0	0	19	29088	44889	65483	13644

【 ☒ 1 2 1 】

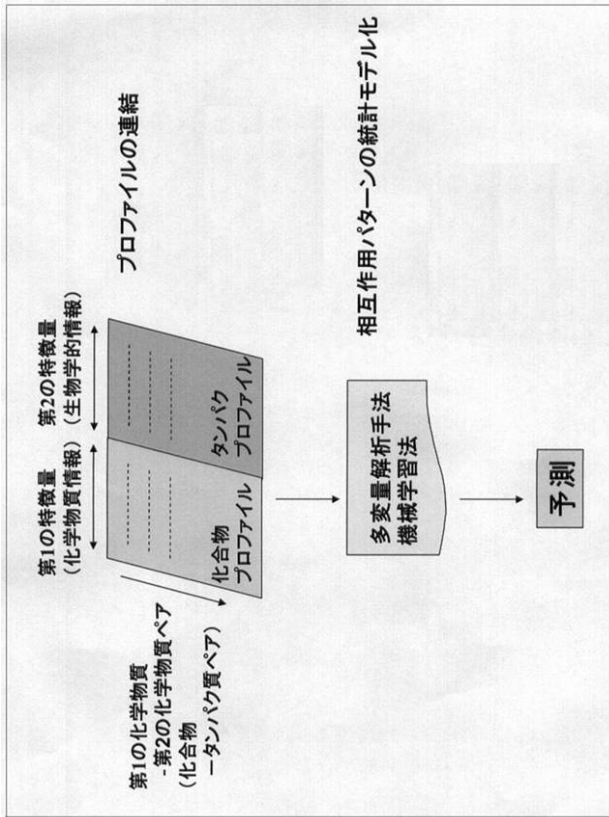
ライブラリー例: PubChem 化合物の生物活性アノテーション G.O. Process

	score 27	score 26	score 25	score 24	score 23	score 22	score 21
biological_process unknown	0	13996	40887	15286	16059	59686	6338
development	0	28342	253537	716385	1308320	1943767	364059
physiological process	0	1403	195556	916314	1875148	2873385	569565
behavior	198	34605	355924	723085	639465	741478	103319
metabolism	0	540	2211	214126	1235709	2301389	665608
catabolism	0	2180	10054	11156	105197	316669	68674
biosynthesis	0	0	49	15710	315536	1014521	300801
pathogenesis	0	2180	11384	67318	176416	588308	132143
cellular process	198	97030	741665	1731806	2718887	3446377	1077203
cell differentiation	0	12	1595	122701	788274	1410128	215008
macromolecule metabolism	0	78	4444	66108	372248	1116652	439101
secretion	0	1	11480	70182	56843	35007	4822
regulation of biological process	0	11402	189000	711117	1315999	1675585	267610
cellular physiological process	77	61946	775127	1953107	3400815	3953376	2351247
response to stimulus	0	2311	211647	986589	2113675	2968214	903314

【 図 1 3 】



【 図 1 4 】



---

フロントページの続き

(72)発明者 辻本 豪三

京都府京都市左京区吉田下阿達町46-29 京都大学大学院薬学研究科内

審査官 宮久保 博幸

(56)参考文献 国際公開第2005/069188(WO, A1)

特開2006-127248(JP, A)

特開2003-159095(JP, A)

特表2002-530727(JP, A)

国際公開第01/069440(WO, A1)

山下慶子, 能動学習法による創薬スクリーニング - 類縁蛋白質のリガンド情報を用いたGPCR  
リガンド探索 -, 第33回構造活性相関シンポジウム講演要旨集, 2005年11月25日,  
p. 63 - 66

Okuno, Y., GLIDA: GPCR-ligand database for chemical genomic drug discovery, Nucleic Ac  
ids Research, 2006年 1月, Vol.34, D673-D677

Yabuki, Y., GRIFFIN: a system for predicting GPCR?G-protein coupling selectivity using  
a support vector machine and a hidden Markov model, Nucleic Acids Research, 2005  
年 7月, Vol.33, W148-W153

(58)調査した分野(Int.Cl., DB名)

G06F 19/10

JSTPlus/JMEDPlus/JST7580(JDreamIII)

PubMed