◇本日の講義スケジュール



序章 13:00~13:20

- ・自己紹介
- ・時代の流れと、研究環境および研究内容や研究スタイルの変化

第一章 13:20~13:40 (10分休憩)

化学データサイエンスおよび人工知能の適用イメージ

第二章 13:50~15:00 (10分休憩)

化学データサイエンスおよび人工知能の適用事例:様々なアプローチが可能

第三章 15:10~15:50 (10分休憩)

化学データサイエンスおよび人工知能の適用手順やパターン

第四章 16:00~16:50

化学データサイエンスおよび人工知能の適用上での留意事項

- ①化学分野の留意点
- ②データサイエンス実施上での留意点

まとめと提案 16:50~17:00

- ①化学データサイエンスおよび人工知能のまとめ
- ②今後の展開についての「オートノマス創薬」の提案
- ・KY法の展開(クラス分類と重回帰型)
- ③自由討論



◇「ケモメトリックス」とは

ケモメトリックスとは:

計量化学(けいりょうかがく、chemometrics)とは、数理科学、統計学、 機械学習、パターン認識、データマイニングなどの手法により、(広義の) 化学分野における諸問題を解決しようとする分野である。

ウィキペディアより:https://ja.wikipedia.org/wiki/計量化学

*ちなみに、chemometricsなので、「化学計量学」 by Yuta

*化学の分野でデータサイエンスを適用する時の基本的な学問分野となる。

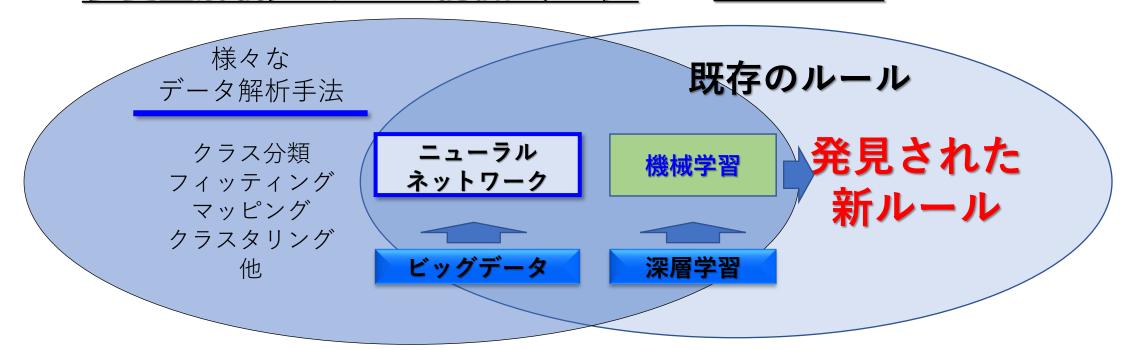


口今後のデータサイエンス手法(DS)と人工知能の展開

現在における多変量解析/パターン認識(DS)と人工知能

多変量解析/パターン認識と人工知能は 機械学習により繋がっている

多変量解析/パターン認識(DS) 人工知能





□本日の解説内容と順番

1. AI (人工知能) に関する概要

- 2. データサイエンスで利用される化学パラメータ
- 3. データサイエンス手法の概要
- 4. KY法の説明



□AI(人工知能)の種類と適用

AIは学習し、内容を<u>まとめる</u>ことは出来る。

AIは学習して、<mark>理解</mark>することは困難である。

この場合のAIは「ニューラルネットワーク型」である。



ロパーセプトロンからニューラルネットワーク

- *ニューラルネットワークは単純なネットワーク構造を利用したパーセプトロンを 基本として発展
- *パーセプトロンは簡単な二分類問題も解決できないということで、衰退
- *パーセプトロンの限界を打破したアプローチとしてニューラルネットワークが提案

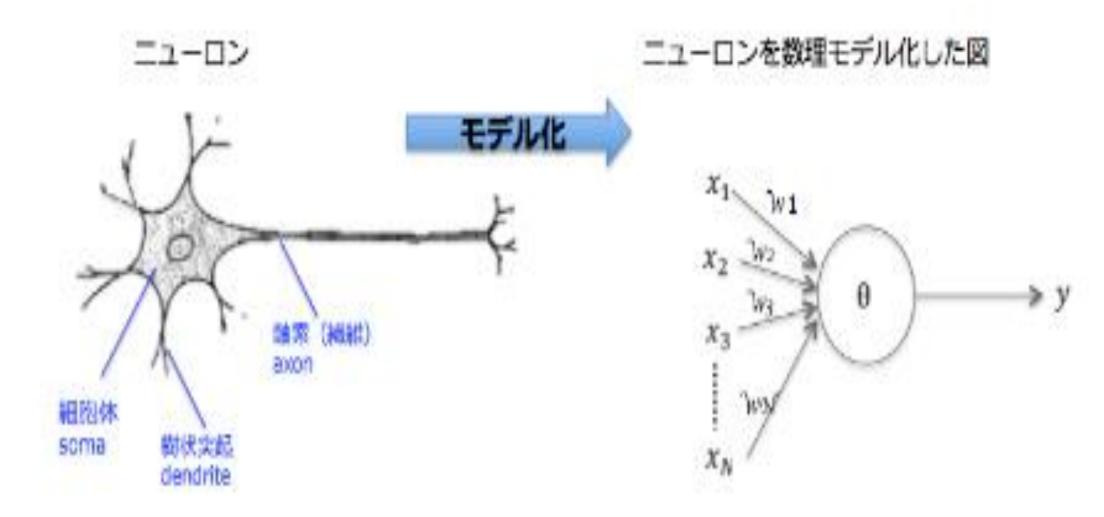
パーセプトロン:線形分類機

ニューラルネットワーク:非線形分類機

- *パーセプトロンは脳の機能を模したアプローチとして開発され、最終目標は AI(人工知能) への展開であった
- *ニューラルネットワークはその改良系で、ネットワーク構造が多層構造となった
- *最近の深層学習はニューラルネットワークのネットワーク構造をさらに複雑にした



ロパーセプトロン

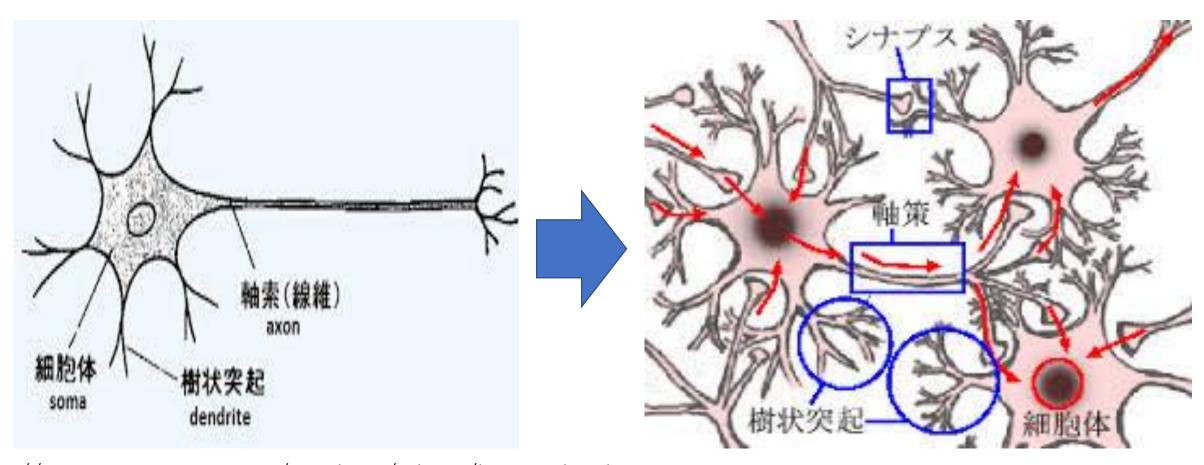


https://udemy.benesse.co.jp/ai/neural-network.html



単二ューロンモデル

ネットワークニューロンモデル



http://www.tamagawa.ac.jp/teachers/aihara/kouzou.html

http://www.sys.ci.ritsumei.ac.jp/project/theory/nn/nn.html

パーセプトロン

入力データ ネットワーク 出力データ 出力結果の判定 教師データ

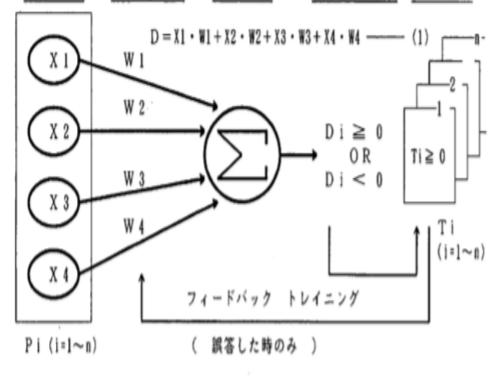
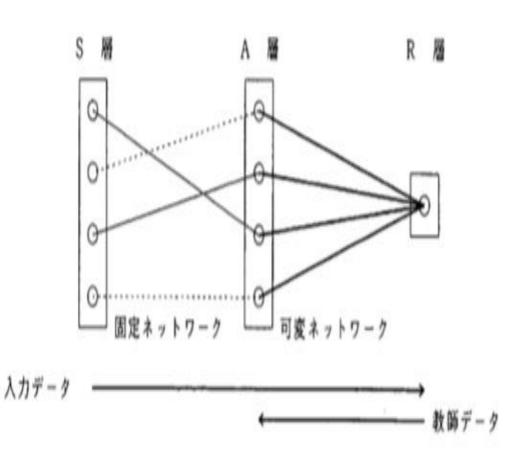


図 2 . パーセプトロンの"学習"の流れ

単二ューロンモデル



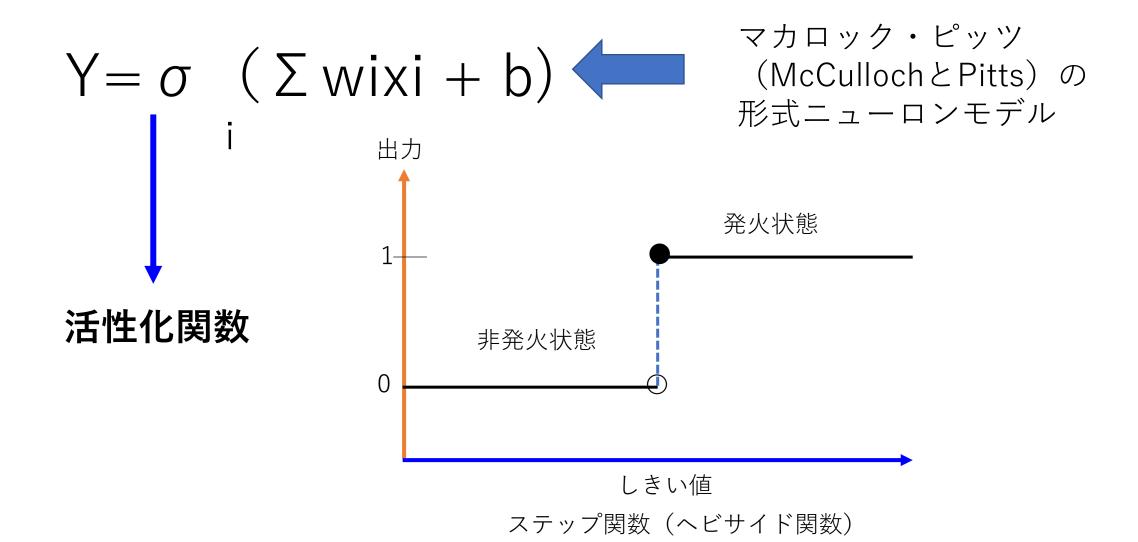
ニューラルネットワーク



ネットワークニューロンモデル

□形式ニューロンモデル

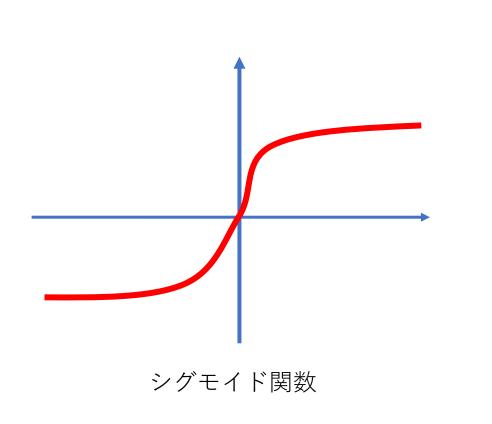


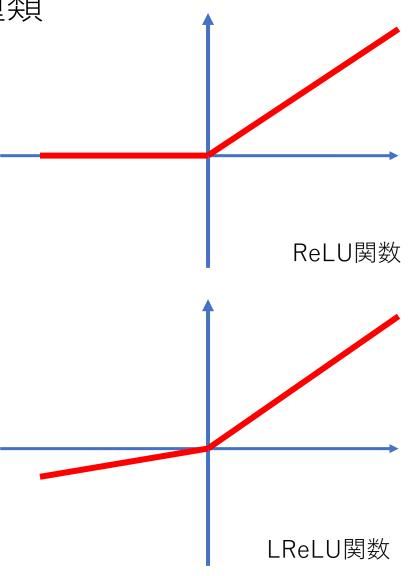


ニューラルネットワークの機械学習で利用されている



ステップ関数の種類





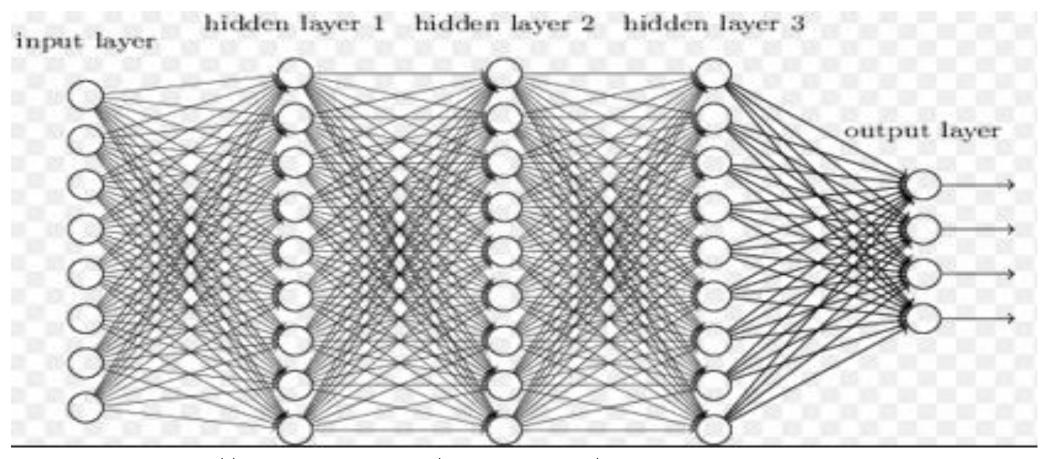
□機械学習型人工知能



- ・実行環境:機械学習手法の発展(深層学習:ディープラーニング)
- ◇ニューラルネットワークから深層学習へ

中間層を増やすことで分類性が向上することは明白であったが、効率的なバックプロパゲーションアルゴリズムが無かったので 展開が遅れていた。しかし、効率的に学習結果を前に戻すアルゴリズムが開発されたことで深層学習が開発された。 現在では、中間層が多層になった深層学習が次世代ニューラルネットワーク(人工知能)として脚光を浴びている。





https://nnadl-ja.github.io/nnadl_site_ja/chap6.html

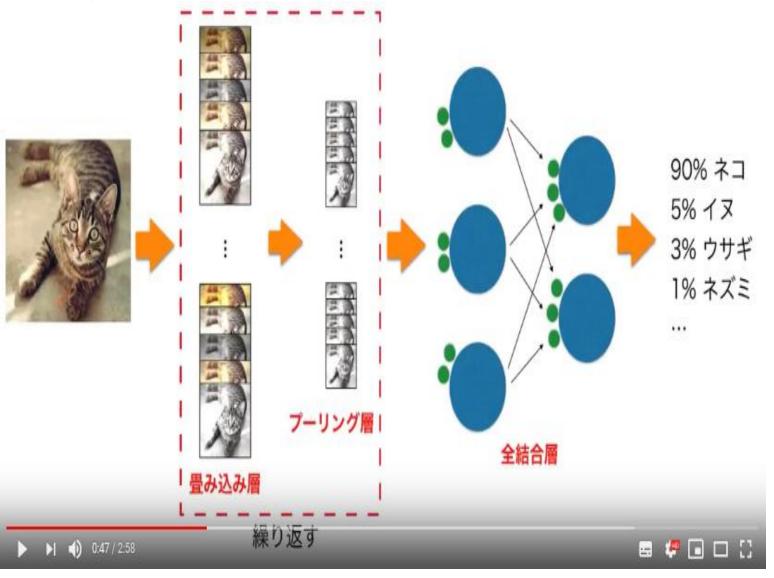
□機械学習型人工知能

In Silico Data
Miracles by the KY-methods

・畳み込みニューラルネットワーク

畳み込みニューラルネットワーク (<u>C</u>onvolutional <u>N</u>eural Network)

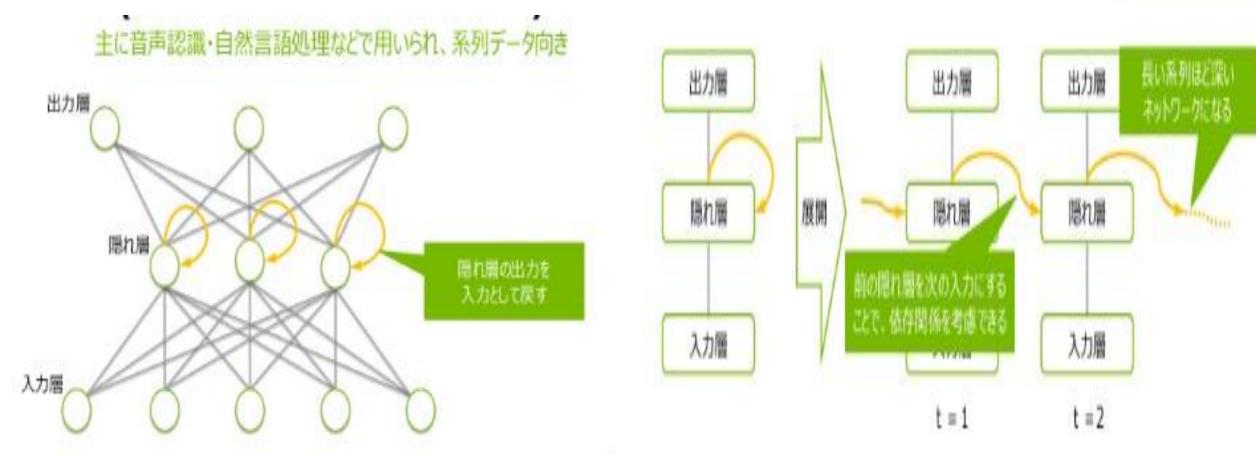
CNNは、全結合層だけでなく **畳み込み層(Convolution Layer)**と **プーリング層(Pooling Layer)**から 構成されるニューラルネットワーク。



〕機械学習型人工知能



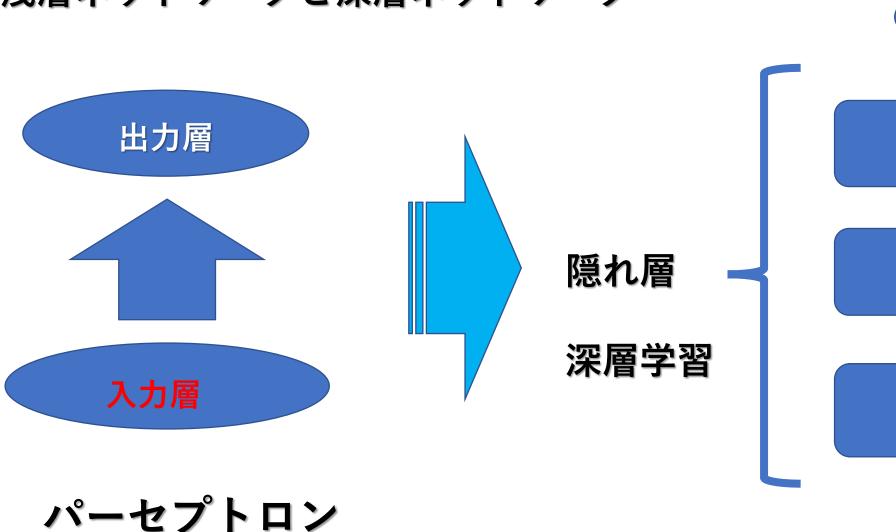
- ・開発歴:リカレント(再帰型)ニューラルネットワーク(RNN)
- ◇リカレント(再帰方)ニューラルネットワーク(Recurrent Neural Networks)
 - ・中間層の出力データを再び入力データとして利用することを特徴とする
 - ・時系列データの解析や言語処理等に利用される

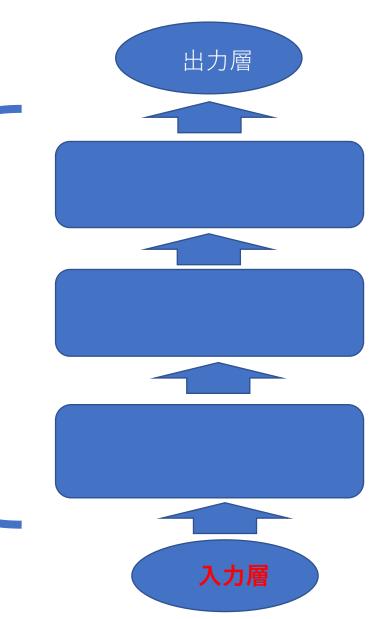


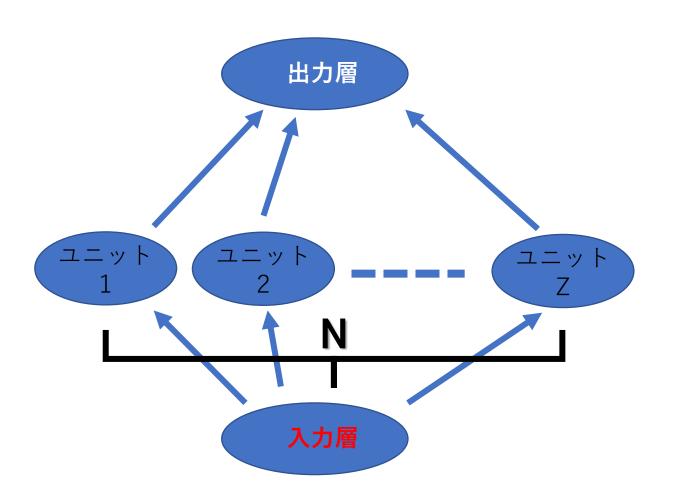
□ネットワークの深層化による効果

In Silico Data
Miracles by the KY-methods

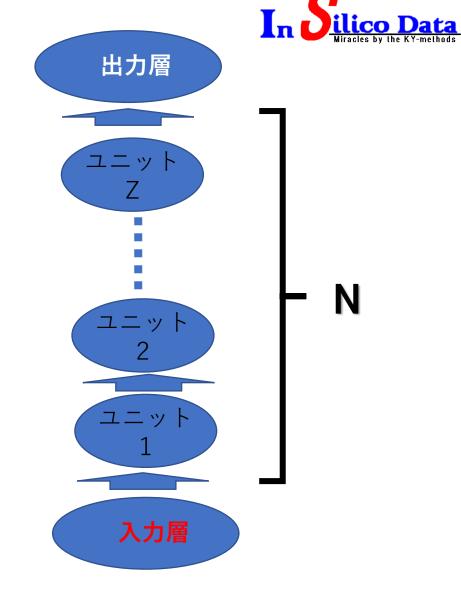
浅層ネットワークと深層ネットワーク







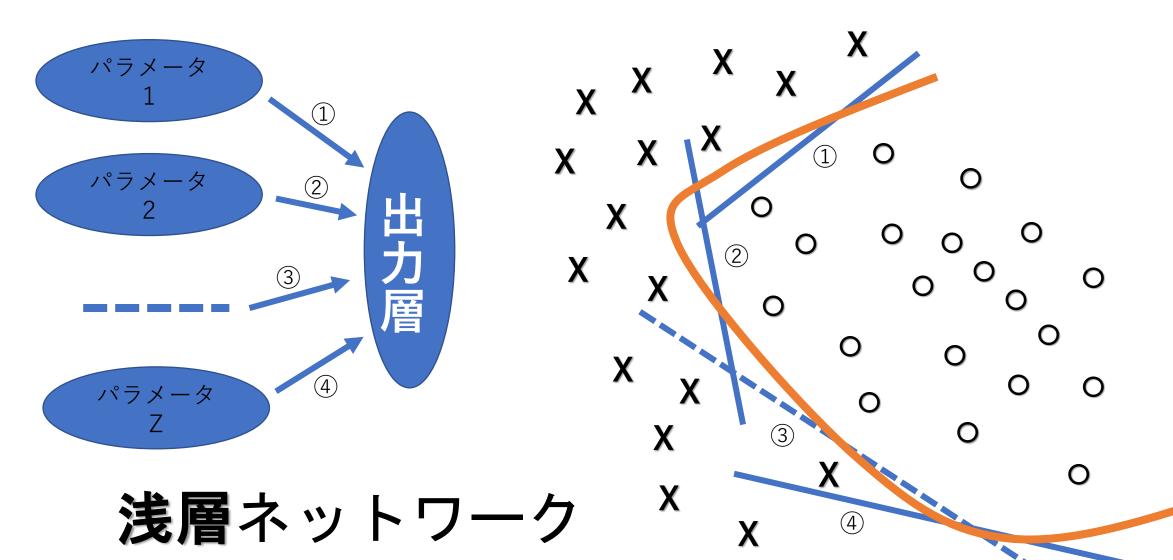
浅層ネットワーク 初期型ニューラルネットワーク



深層ネットワーク

□浅層ネットワークによる非線形化





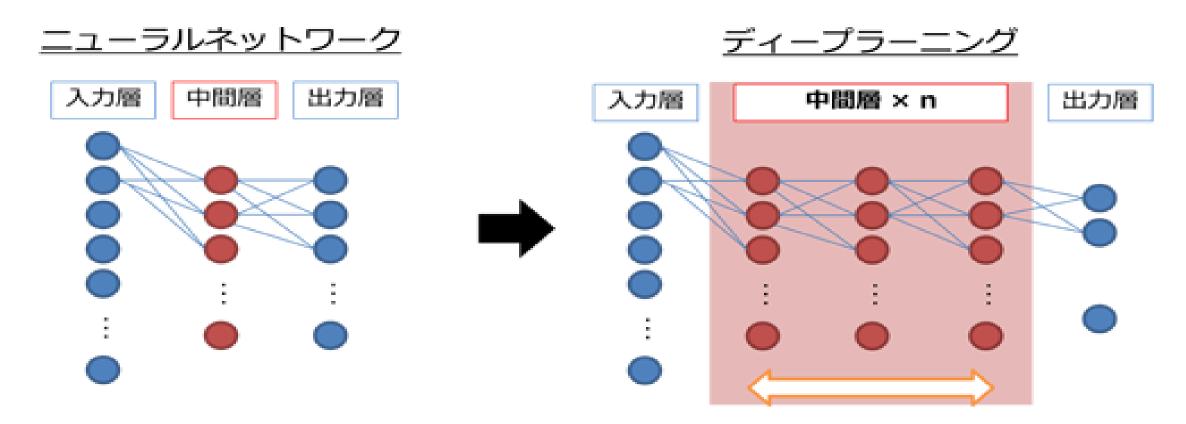
中間層のユニット数が大きければ、任意の関数を表現できる。非線形性が高い。

□深層ネットワークによる非線形化 入力層 ユニット1 X ユニット2 X ユニットZ 3 出力層 深層ネットワーク 4

浅層ネットワークよりも遥かに複雑な関数を表現できる。

口機械学習における問題





多層化による**勾配消失**と過剰適合の問題があったが、近年、アルゴリズムの改良とデータ量の増大、そして膨大なデータを処理できる計算装置(コンピュータ)の爆発的な性能向上によって問題が解消されてきた。



□本日の解説内容と順番

1. AI (人工知能) に関する概要

2. データサイエンスで利用される化学パラメータ

3. データサイエンス手法の概要

4. KY法の説明

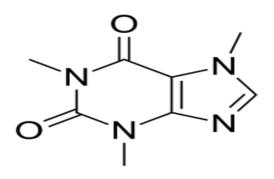
□化学データ(パラメーター)の種類



化合物に起因するデータを「化学データ(パラメーター)」とする 構造式の有無や内容により、以下の3種類に分類される

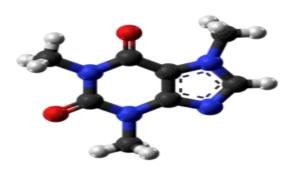
構造式不要

物性データ 機器スペクトルデータ 医療関連データ



2次元構造式関連データ

分子式関連パラメーター トポロジカルパラメーター 部分構造パラメーター フラグメントパラメーター



3次元構造式関連データ

トポグラフィカルパラメーター 3次元構造関連データ 力場関連パラメーター 量子化学関連パラメーター

□化学データ(パラメーター)の種類



◇化合物に由来する化学データ(パラメーター)

化合物関連データ解析を行うにはパラメーター、即ち数値データが必要となる

化合物構造式由来のパラメーター

■ トポロジカルデータ

分子構造インデックス:原子数 (原子種)、結合数 (結合種)、リング数、その他様々なインデックス値:HOSOYAインデックス、分子結合インデックスMC値

パス値インデックス、

■ トポグラフィカルデータ

化合物の3次元的形状に関するパラメータ

化合物全体構造 :ボックスパラメータ、対称パラメータ、

立体格子パラメータ、その他

化合物部分構造 :ステリモルパラメータ、

■ 物理化学データ

分子に関する様々な物性データ : 分子屈折率、分子量、LOGP、融点、沸点

分子容積、分子表面積、その他

分子軌道法より得られる様々なパラメータ:電子密度、HOMO、LUMO、他

分子力学計算から得られるパラメータ: 種々歪みエネルギー 種々スペクトルより得られるデータ: 種々スペクトルデータ

■ その他のデータ

部分構造パラメータ: 部分構造の有無、部分構造数、

部分構造単位の様々なパラメータ値計算、

演算パラメータ1 : 記述子間の演算により得られるパラメータ (+-x + Log)

演算パラメータ2 : 他の解析手法より算出されたパラメータ

ダミーパラメータ : 有るパターン存在の有無(1/0)に関するパラメータ

機器スペクトル由来 のパラメーター

Mass
IR
H-NMR
C-NMR
UV
GC
HPLC

Raman X線分析 その他 化合物に起因する 解析目的パラメーター

> 薬理活性 毒性 ADME 物性 環境毒性

□化学データ解析実施上での必要データ

In Silico Data
Miracles by the KY-methods

◇汎用的なデータ解析の流れ図:入力データ関連

化合物関連データ



- *化合物構造式(1/2次元)
- *物性データ
- *機器スペクトルデータ
- *バイオ関連データ
- * 医療関連データ (画像等)

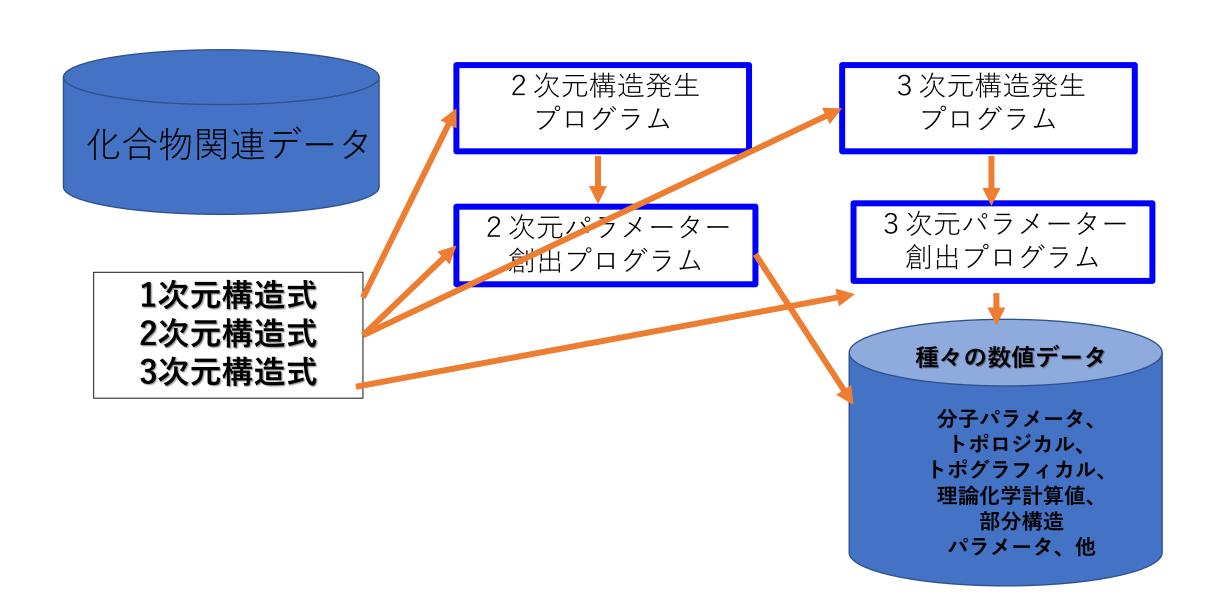
解析目的関連データ



- *薬理活性
- *毒性関連データ
- *機器スペクトルデータ
- *バイオ関連データ
- *患者関連データ



口化合物構造式を用いたパラメーター発生





◇汎用的なデータ解析の流れ図:データ前処理関連 解析目的と無関係な情報を有するパラメーター除去(特徴抽出)

パラメーター選択(特徴抽出)の実施

種々の数値データ

分子パラメータ、 トポロジカル、 トポグラフィカル、 理論化学計算値、 部分構造パラメータ、 演算パラメーター その他

*欠陥データ処理:

欠損データ、0値データ、同値データ、他

*統計的処理:

単相関、重相関、Fisher比、他

*データ解析手法による個別特徴抽出

主成分解析、パーセプトロン、遺伝的アルゴリズム、他

*パラメーター桁の調整

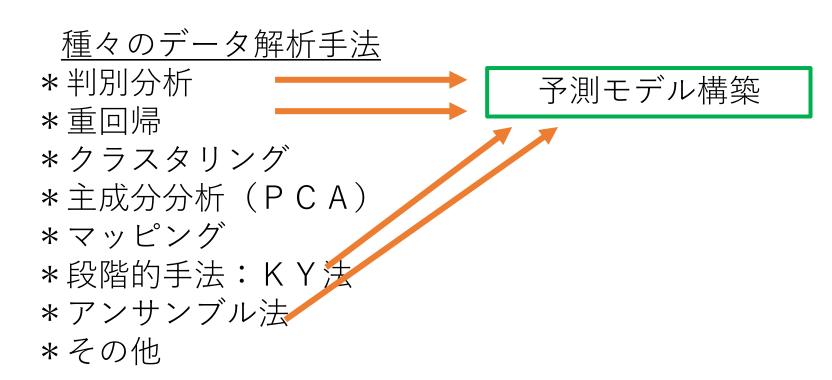
オートスケーリング



◇汎用的なデータ解析の流れ図:多変量解析/パターン認識関連

種々の数値データ

分子パラメータ、 トポロジカル、 トポグラフィカル、 理論化学計算値、 部分構造 パラメータ、他

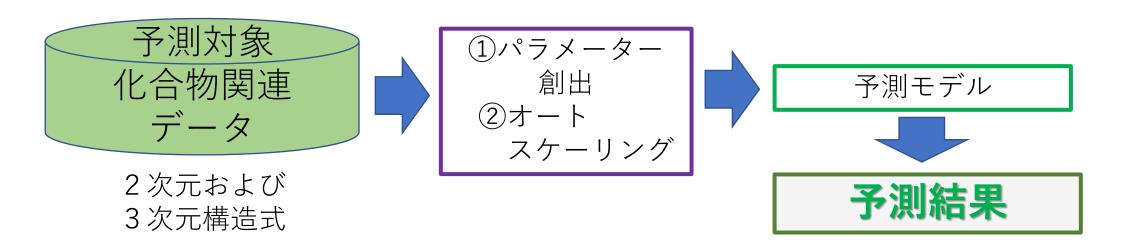


前処理後のパラメーター群

◇汎用的なデータ解析の流れ図:予測の実施



- *予測モデルを構成するパラメーターを化合物構造式から創出
- *構造式から発生できないパラメーターは外部パラメーターとして導入
- *オートスケーリングを実施



□データサイエンスに使う種々の化学パラメーター



◆化合物関連パラメーター

■ トポロジカルデータ

分子構造インデックス:原子数 (原子種)、結合数 (結合種)、リング数、その他様々なインデックス値:HOSOYAインデックス、分子結合インデックスMC値

パス値インデックス、

■ トポグラフィカルデータ

化合物の3次元的形状に関するパラメータ

化合物全体構造 :ボックスパラメータ、対称パラメータ、

立体格子パラメータ、その他

化合物部分構造 : ステリモルパラメータ、

■ 物理化学データ

分子に関する様々な物性データ : 分子屈折率、分子量、LOGP、融点、沸点

分子容積、分子表面積、その他

分子軌道法より得られる様々なパラメータ:電子密度、HOMO、LUMO、他

分子力学計算から得られるパラメータ: 種々歪みエネルギー

種々スペクトルより得られるデータ: 種々スペクトルデータ

■ その他のデータ

部分構造パラメータ: 部分構造の有無、部分構造数、

部分構造単位の様々なパラメータ値計算、

演算パラメータ1 : 記述子間の演算により得られるパラメータ (+-x + Log)

演算パラメータ 2 : 他の解析手法より算出されたパラメータ

ダミーパラメータ : 有るパターン存在の有無(1/0)に関するパラメータ



□AIに使う種々の化学パラメーター

深層学習による人工知能での化学パラメータ

- 1. フィンガープリント
- 2. フィーチュアベクトル
- 3. グラフコンボリューション

深層学習が発表され、化合物を扱う研究分野への適用初期は 上記パラメータが利用されていた。

これらのパラメータ(特に 2, 3番)は、化合物をトポロジーと考える人工知能研究者が持ち出したパラメータであり、あまり化学的な意味はないといえる。

構造一活性相関やケモメトリックスの分野では、古くから様々な化学パラメータが開発されているので、そちらを用いる方が化学的な解析には適している。



◆化合物関連パラメーター:トポロジカルパラメーター

トポロジカルデータの特徴 トポロジカルデータは化合物を構成する原子と結合とを、それぞれノードとエッジとに

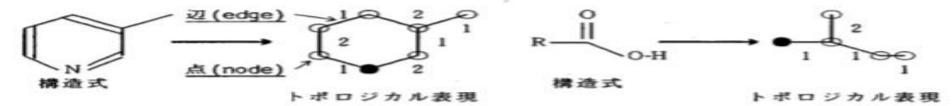
定義する。 トポロジー的な問題として化合物構造式を捕らえ、化合物の原子 (ノード) 間の相互的な関係(つながり状態)を数値データに変換したものである。

このトポロジカルデータの特徴を簡単にまとめると以下のようになる。

化合物の複雑な結合情報を数値データに変換できるので、通常の数値データ では説明出来ないような複雑な情報を扱う事が可能となる。 この結果、分

類能が飛躍的に向上することが期待される。

数値データの変換ルールと化合物構造式との関係が不明な時が多い。 ゴリズムが数値データへの変換の為のルールとなっている事が多く、最終目的である目的 変数に対する情報の説明や解釈が困難な事が多い。 即ち、分類の為のデータに陥り易く、 分類だけが目的の時は強力なパラメータとなりうるが、そのパラメータの持つ意味(情 報)を解釈する事が重要となる解析には不向きである。



このトポロジカルデータは現在様々なものが提唱されている。 特に有名なものとして 化合物の物性予測に用いられる事の多いHOSOYA INDEX と、構造活性相関分 野で利用実績の多い分子結合インデックス(M. C.) (Molecular Connectivity Inde x) 等が有名である。



◆化合物関連パラメーター:トポロジカルパラメーター

□ MC I 値の算出法 まず化合物を構成している個々の結合についてCκ値を求める。 続いて、このCκ値 を化合物中の総ての結合について総和した値が分子に対するMC I 値となる。

$$M C I = \sum_{k=1}^{m} C_k = \sum_{k=1}^{m} \frac{1}{[L_1 \cdot L_1]_{*}^{1/2}}$$

k: ある一つの結合の I D 番号

i: 結合kを形成する原子2個のうちの一つの原子に関するID j: 結合kを形成する原子2個のうちi以外の原子に関するID

上式中、L: は原子iの結合の多重度であり、L: は原子jの多重度を示している。 この多重度とは現在注目している原子から飛び出している結合の数を意味し、この時水素 原子とつながっている結合の数は無視して計算する。 例) C: 値の求め方

$$-\frac{1}{3} - \frac{k}{3} - \frac{1}{3} - \frac{1}{3} - \frac{1}{3} - \frac{1}{3}$$

$$-\frac{k}{3} - \frac{k}{3} - \frac{1}{3} - \frac{1}{3} - \frac{1}{3} - \frac{1}{3}$$

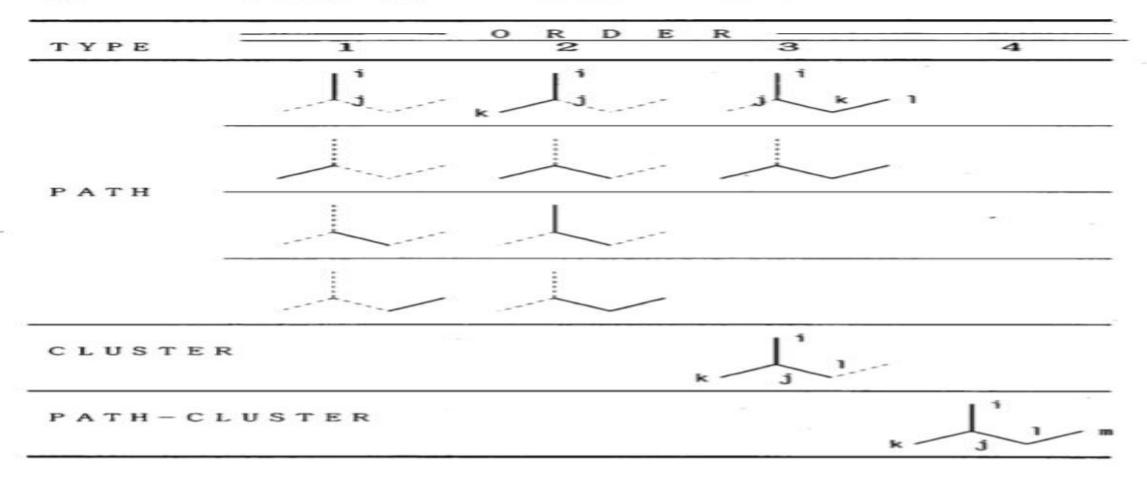
$$-\frac{k}{3} - \frac{k}{3} - \frac{1}{3} - \frac{1}{3} - \frac{1}{3} - \frac{1}{3}$$

$$-\frac{1}{3} - \frac{k}{3} - \frac{1}{3} - \frac{1}{3} - \frac{1}{3} - \frac{1}{3} - \frac{1}{3}$$



◆化合物関連パラメーター:トポロジカルパラメーター

例) MCIにおける次数と結合タイプの概念及びCx 計算式





- ◆化合物関連パラメーター:トポロジカルパラメーター
- □ MCIへの結合次数及び結合タイプの導入 Cκ値を求め、この値を基準としてMCIを求める時、化合物構造式の複雑さを情報 として取り入れるべく結合次数(BOND ORDER)という概念と結合タイプ(BO ND TYPE)という2つの概念を導入する。
- ・結合次数(BOND ORDER)は C_{κ} を求める時の対象となる結合と、その結合を 形成する原子の数を拡大してゆくものである。
- ・結合タイプ (BOND TYPE) とは、結合が複数集まって一つのCx を形成する時の集合形態に関する情報である。
- □ 結合次数 (BOND ORDER) について 結合次数は基本となる C に値を求める時に対象とする結合や原子数を規定するもの である。 次数が小さければMCIの値は大きく、次数が増大するにつれてMCIの値は 小さくなる。
- □ 結合タイプ (BOND TYPE) について TYPEはC、としてまとまった単位(特に次数が大きくなった時)の形を規制するものである。
 - ・PATHは最も単純な形をしており、結合が直線上に繋がっているものを意味する。 この時、次数が1のものは直線であり、PATHとみなす。
 - ・CLUSTERは分岐した形状を持つC。となる。 従って、次数が3以上で現れる
 - · PATH-CLUSTERはC,内部にPATH部分とCLUSTER部分を持つ。



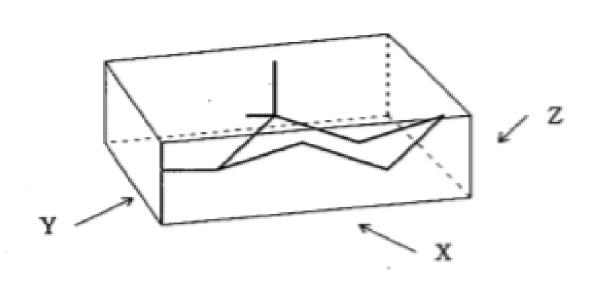
◆化合物関連パラメーター:トポロジカルパラメーター

```
例)
    「ン・: 結合次数 I、PATHタイプの C×値を基本として求めたMC I値
    ' × rc: 結合次数 4、PATH-CLUSTERタイプのCx 値を基本として
                    PATHタイプのCx値を基本として求めたMCI値に
                  、PATHタイプのC。値の計算にヘテロ原子を考慮して
        沙欠 娄女 1
                      s=1
        沙欠 类女 2
                      s=1
                       N.
                      \Sigma (\delta, \delta, \delta, \delta, \delta);
        沙欠 婆女 3
                      s=1
                       N-
```



◆化合物関連パラメーター:**トポグラフィカル**パラメーター

② 分子全体の形状に関する幾何学的情報(ボックスパラメータ)



化合物の3次元構造式をそのまま長方形の ボックスにいれる。 このボックスの各軸 の長さとその比とをパラメータとする。

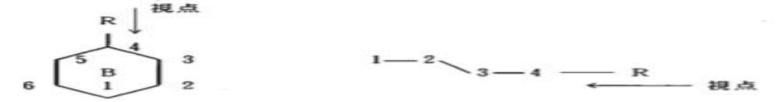
> パラメータ 1 = Xパラメータ 2 = Yパラメータ 3 = Zパラメータ 4 = X/Yパラメータ 5 = X/Zパラメータ 6 = Y/Z

このパラメータにより、分子全体の立体的な形状についての情報がえられる。 例えば、分子が平面に近い、細長い、立方体に近い等の情報である。

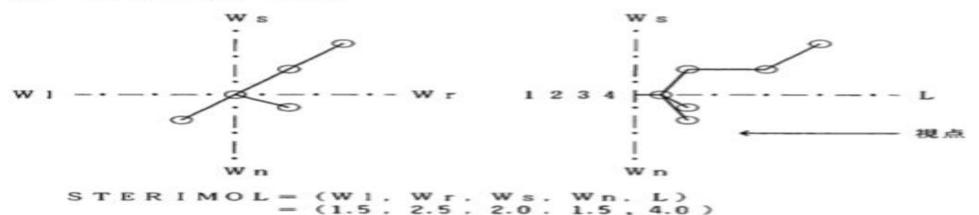


◆化合物関連パラメーター:**トポグラフィカル**パラメーター

①STERIMOL PARAMETER
このパラメータは化合物の置換基Rの3次元的立体情報を記述するのに用いられる。
特に、重回帰手法によるHANSCH/FUJITA法等に用いられて数多くの実績を有する構造活性相関には重要なパラメータである。



パラメークは化合物の基本構造部分(図中1~6で示されるB部分)と置換基R部分とに 分けた時、基本構造部分と置換基R部分とが直結している結合をRの方からBに向かって 見た時の置換基Rの占める空間上の領域をそれぞれの軸方向について分割した時の値を要 素データとするものである。





◆化合物関連パラメーター:物理化学パラメーター

◇種々の物性:

分子量、 融点、 沸点、分子屈折率、 LogP、 $Hammet \sigma$ 、 その他

◇分子軌道法関連パラメーター:

電子密度、 HOMO、 LUMO、 分極率、 双極子モーメント、その他

◇分子力学関連パラメーター:

結合エネルギー、トーションエネルギー、 水素結合エネルギー、 その他



◆化合物関連パラメーター:物理化学パラメーター (LogP)

HANSCH-REOによるフラグメント付加方式によるLOGP値推算。

$$LOGP = LOG - [C] lipid$$
[C] aqueous

[C] lipid : 平衡状態における油層中の濃度[C] aqueous : 平衡状態における水層中の濃度



◆化合物関連パラメーター:物理化学パラメーター (LogP)

```
① フラグメント付加方式によるLOGP推算式
           i番目のフラグメントの出現回数
           i番目のフラグメントに対するフラグメント定数値
           j 番目の修正因子の修正定数値
LOGP値計算例)
 フラグメント定数
  フラグメント・
  N - N = 0
```



◇ 機器スペクトルパラメーター

■有機化合物のスペクトルデータベース **SDBS**

SDBS No: 1898 CAS Registry No.: 58-08-2

DOI:

Molecular Formula: C₈H₁₀N₄O₂ Molecular Weight: 194.2

SDBS-NO= 1898

CAFFEINE

Compound Name:

caffeine

1,3,7-trimethyl-3,7-dihydro-1H-purine-2,6-dione

1,3,7-trimethyl-3,7-dihydro-purin-2,6-dion, kaffein

1,3,7-trimethyl-3,7-dihydro-purine-2,6-dione

1,3,7-trimethylxanthine

1H-purine-2,6-dione, 3,7-dihydro-1,3,7-trimethyl-

3,7-dihydro-1,3,7-trimethyl-1H-purine-2,6-dione

theine

InChI:

InChI=1S/C8H10N4O2/c1-10-4-9-6-5(10)7(13)12(3)8(14)11(6)2/h4H,1-3H3

InChIKey:

RYYVLZVUVIJVGH-UHFFFAOYSA-N

Publisher:

National Institute of Advanced Industrial Science and Technology (AIST)

-



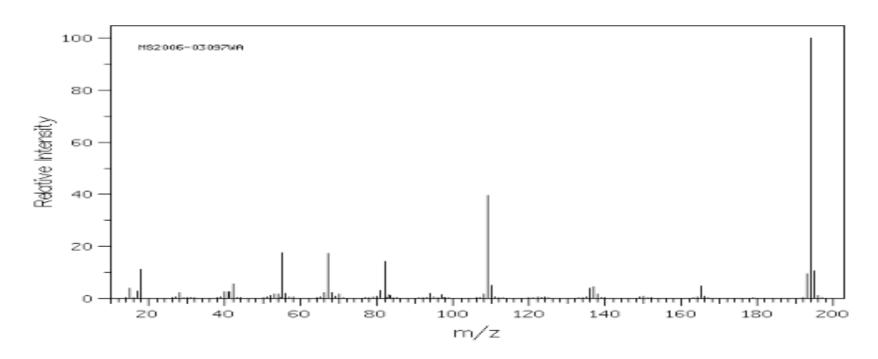
◇ 機器スペクトルパラメーター

SDBS-Mass

MS2006-03097WA caffeine C8H10N4O2 SDBS NO. 1898

(Mass of molecular ion: 194)

Mass

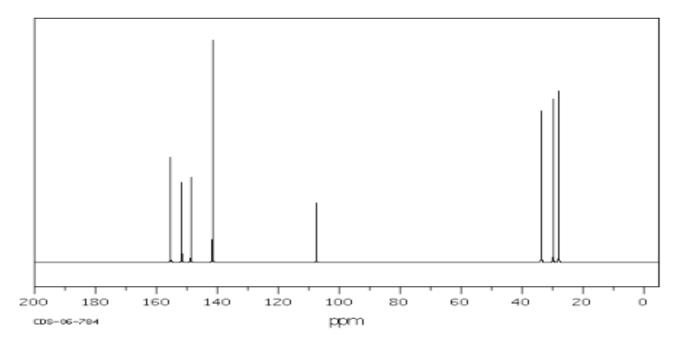




◇ 機器スペクトルパラメーター:

SDBS-¹³C NMRSDBS No. 1898CDS-06-784 C₈ H₁₀ N₄ O₂ caffeine

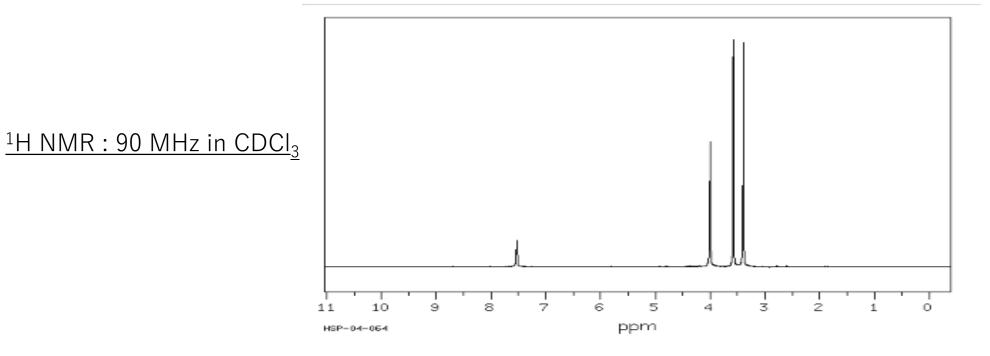
¹³C NMR: in CDCl₃





◇ 機器スペクトルパラメーター

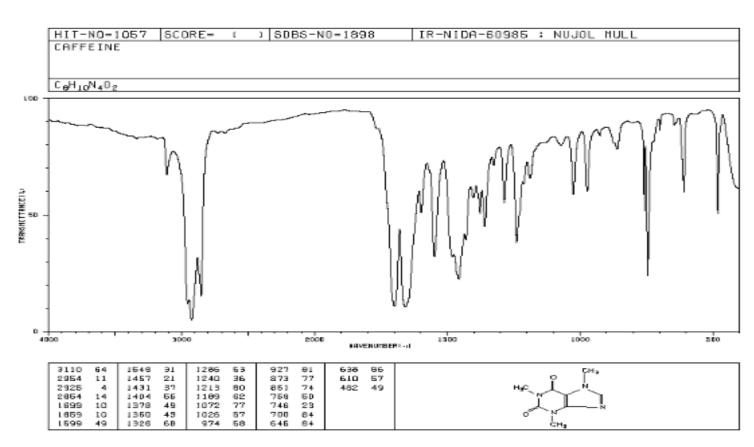
SDBS-¹H NMRSDBS No. 1898HSP-04-064 C₈ H₁₀ N₄ O₂ **caffeine**





◇ 機器スペクトルパラメーター

IR: nujol null





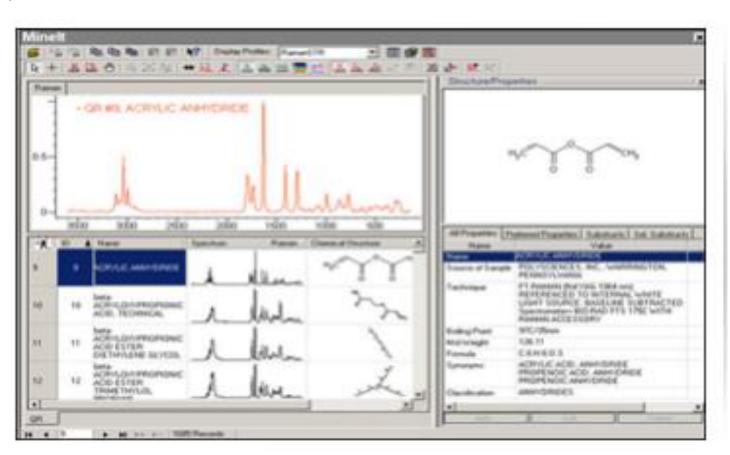
 $\langle \rangle$

機器スペクトルパラメーター

Raman データ

BIO-RADの スペクトルデータベースより

http://www.bio-rad.com/



http://www.bio-rad.com/ja-jp/product/raman-spectral-databases?ID=N0ZXPS4VY

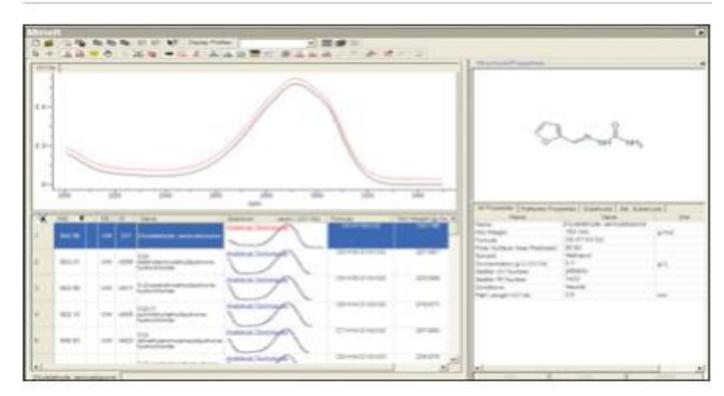


◇ 機器スペクトルパラメーター

紫外可視データベース

UV-Visスペクトル

BIO-RADの スペクトルデータベースより http://www.bio-rad.com/

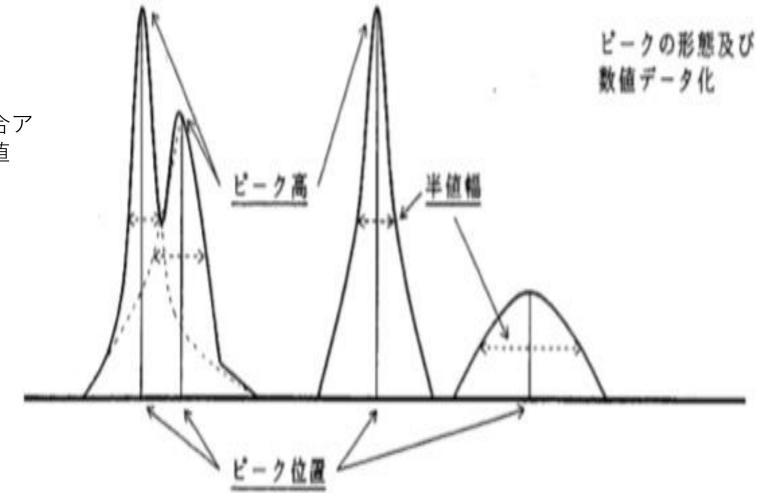


http://www.bio-rad.com/ja-jp/product/uv-vis-spectral-databases?ID=NH262L4VY



◆ 機器スペクトルパラメーター

*機器スペクトルデータは殆どの場合アナログデータなので、デジタルの数値 データに変換することが必要である。





◆ 機器スペクトルパラメーター

特徴:

- ①機器があれば数値データとして簡単に蓄積できる
- ②スペクトルチャートは様々な実験過程で種々蓄積される

留意点:

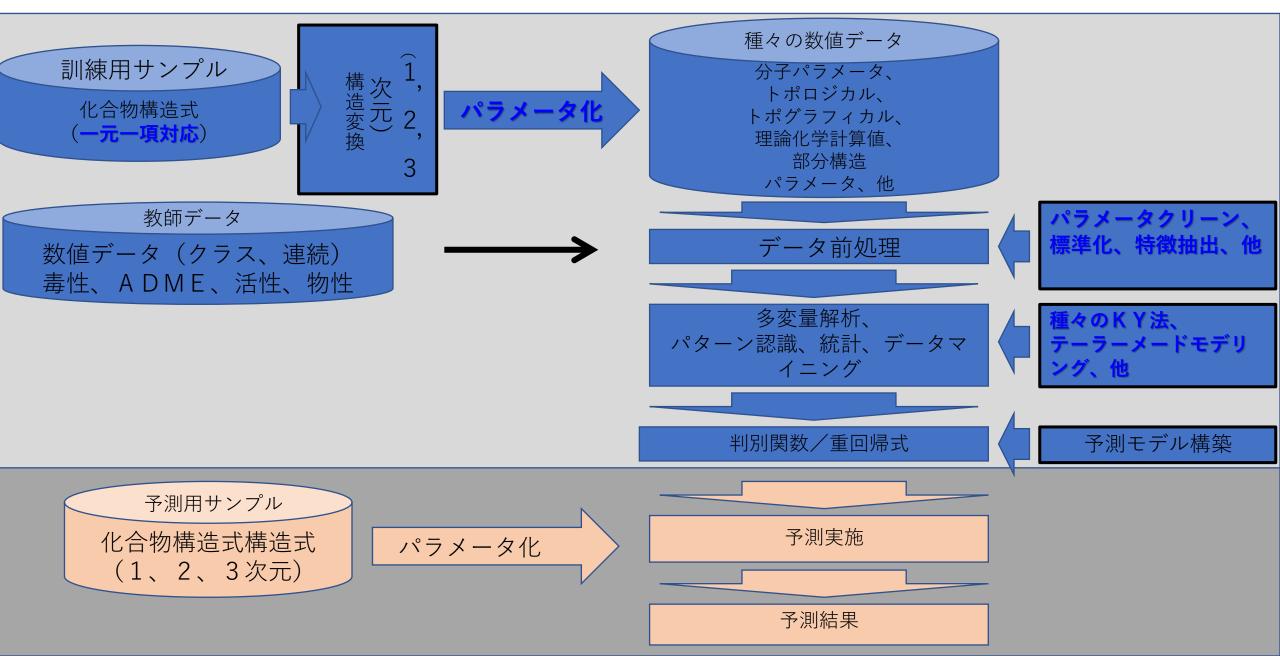
- ①一般的にパラメーター数が極めて大きくなりやすい 結果として、データ解析手法が限定、適用不可となる可能性が高くなる
- ②スペクトルデータは多重共線性が極めて高い ①と②の特徴により、データ解析実施においては次元圧縮・統合等が必要で、 解析手法もPLSやPCAと制限されることが多い
- ③スペクトルチャートの測定条件等統一が必要 出来れば測定機器メーカーや機種も統一
- 例:60MのH-NMRデータと90MのH-NMRデータを混在してのデータ解析は無意味
- ・測定条件等が統一されないとデータ解析の精度が保証されにくくなる



- □本日の解説内容と順番
- 1. AI (人工知能) に関する概要
- 2. データサイエンスで利用される化学パラメータ
- 3. データサイエンス手法の概要
- 4. KY法の説明

◇汎用的なデータ解析の流れ図:解析全体の流れ





□ データ解析手法の機能的分類



統計関連情報の扱い方によりデータ解析手法は大きく二種類に分類される

1. パラメトリック手法

解析に用いるパラメーターの母集団分布情報を用いて解析する手法 分布(正規分布等)等の情報、平均値、分散、サンプル数は大きいものが良い 統計の検定等を実施するときに利用される

2. ノンパラメトリック手法

解析に用いるパラメーターの母集団分布情報が不明な環境で解析する手法 分布型は問わない、サンプル数も小さくても実施可能、他 制限事項が少ないので、広範囲にわたって適用できる 基本的に多変量解析/パターン認識として展開される手法はノンパラメトリック



◇ データ解析手法の機能的分類

解析アプローチの差異に伴い大きく以下の解析手法に分類される 個々の手法での代表的な手法をリストする

- 1. 判別分析: 二クラス分類、多クラス分類
- 2. 重回帰(フィッテイング):単回帰、重回帰
- 3. クラスター分析:階層型クラスタリング、非階層型クラスタリング
- 4. 種々マッピング:主成分分析(PCA)、非線形写像(NLM)
- 5. チャート表示:手法⇒レーダーチャート、顔チャート、ラインチャート、
- 6. 決定木:手法⇒C5.0、CART
- 7. アンサンブル学習法:手法⇒AdaBoost、ランダムフォレスト
- 8. 段階的手法:手法⇒KY法
- 9. 最適化手法:最小二乗法、シンプレクス法、遺伝的アルゴリズム

□ニクラス分類の基本概念



判別関数とドラグデザイン的情報解析

 $Y = a_1 x_1 \pm a_2 x_2 \pm \cdot \cdot \cdot \pm a_n x_n \pm const.$

Y: 目的薬理活性/毒性/他

Y ≧ 0

・活性あり

・毒性あり

Y < 0

・活性無し

・毒性無し

□構造−活性/毒性/他との相関解析

係数 **a**_i ≥ **0** の時 パラメータ **X i** の持つ情報は ・活性向上、・毒性強化要因 ・活性低下、毒性低下要因

構造-活性相関、構造-毒性相関

□ニクラス分類の基本概念



判別関数、重回帰式からの情報取り出し

□パラメータとウェイトベクトルの利用

1 要因の抽出および明確化:パラメータ利用

最終判別関数や重回帰式中で用いられているパラメータが持つ情報内容の評価

2. 寄与の方向性:ウエイトベクトル利用

②係数 $a_i < 0$ の時 ____、 パラメータXi の持つ情報は

・活性低下、毒性低下要因

3 寄与の相対的な貢献度:ウエイトベクトル利用

係数の絶対値比較による貢献度の評価

□データサイエンス手法の基本事項

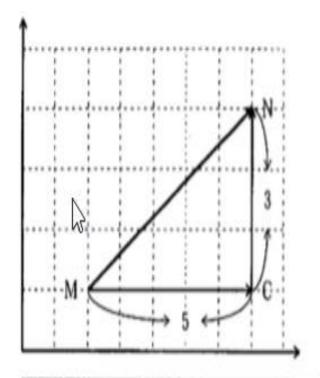


◇サンプル間の**距離**

N次元空間上でのサンプル間の距離の 算出は様々な多変量解析/パターン認識手 法の基本となる。

サンプル間の距離算出に様々なメトリック手法が展開されている。

右は、ユークリッド距離とシテイブロック 距離による距離計算事例である



2次元空間中でA点とB点との距離を計る時、ユークリッド距離では距離Dxxを算出し、シテイプロック距離では距離Dxcと距離Dcxとの距離を合わせたものをA点とB点との距離とする。

$$M = (2.1)$$

 $N = (7.4)$

ユークリッド距離 =
$$D_{MN}$$
 = $(5^2 + 3^2)^{1/2} = 5.83$
シテイプロック距離 = D_{MN} + D_{CN} = 5 + 3 = 8

図1. ユークリッド距離とシテイブロック距離

◇サンプル間の距離:**距離基準**



(1) ミンコフスキー (MINKOWSKI)距離

$$D = \left[\sum_{i=1}^{d} (X_{w_i} - X_{w_i})^k \right]^{1/k}$$

(2) ユークリッド(EUCLIDEAN) 距離

ユークリッド距離はミンコフスキー距離の式において、kが2の時にあたる。

$$D = \left[\sum_{i=1}^{d} (X_{Mi} - X_{Mi})^{2} \right]^{1/2}$$

(3) シテイブロック (CITY BLOCK)距離

シテイブロック距離は2パターン間の最短距離をとるのではなく、直交する2線の距離の総和をとるものである。

$$D = \sum_{i=1}^{d} \left[X_{Ni} - X_{Ni} \right]$$

(4) キャンベラ (CANBERA)距離

$$D = \frac{\sum_{i=1}^{d} \left[X_{Ni} - X_{Ni} \right]}{\sum_{i=1}^{d} \left[X_{Ni} + X_{Ni} \right]}$$

(5) ハミング(HAMMING) 距離

ハミング距離は1/0のバイナリデータで利用される事が多いが、ORとAND の識別を効率良く行う事が出来る。

$$D = \sum_{i=1}^{d} \left[X_{Mi} + X_{Ni} - 2 X_{Mi} X_{Ni} \right]$$

尚、1/0のバイナリーデータを用いた時、ハミング距離とシテイブロック距離 は同じものとなる。

6) 谷本距離

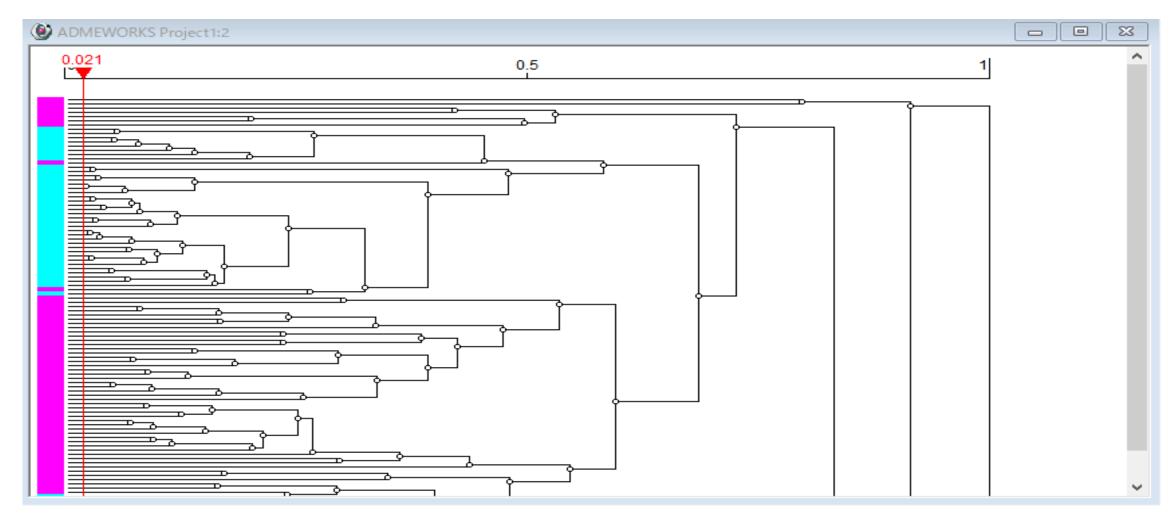
谷本距離はハミング距離がその性格上、1が少ないデータは不利に評価されるという欠点を改良したものである。

$$\Sigma = \frac{\sum_{i=1}^{d} \left[X_{Mi} + X_{Mi} - 2 X_{Mi} X_{Mi} \right]}{\sum_{i=1}^{d} \left[X_{Mi} + X_{Mi} - X_{Mi} X_{Mi} \right]}$$

◇個別手法の特徴と適用内容:クラスター分析



ModelBuilderの画面より





◇個別手法の特徴と適用内容:クラスター分析 クロスター分析は以下の基準にて様々な手法に分類される

クラスター化 アプローチ

- ■階層的クラスタリング 解析結果はデンドログラム として出力される
- ■非階層的クラスタリング 解析結果は単に クラスターの数とメンバー

クラスター化 アルゴリズム

- ①Division Method (分割法)
- ② Aggregative Method (凝集法)

融合法の種類





◇個別手法の特徴と適用内容:クラスター分析 クラスタリングの融合手法

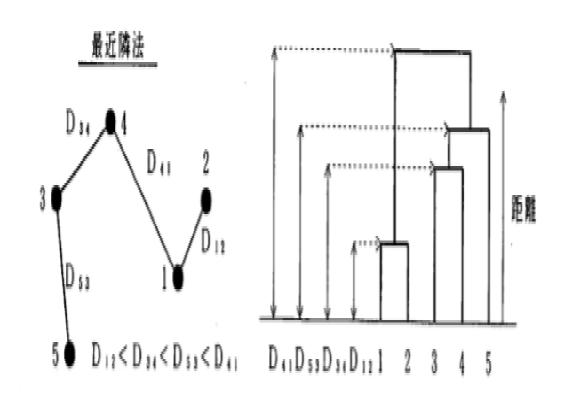


図2. 最近隣法によるクラスタリング手続きとデンドログラム

図3. 重心法によるクラスタリング手続きとデンドログラム

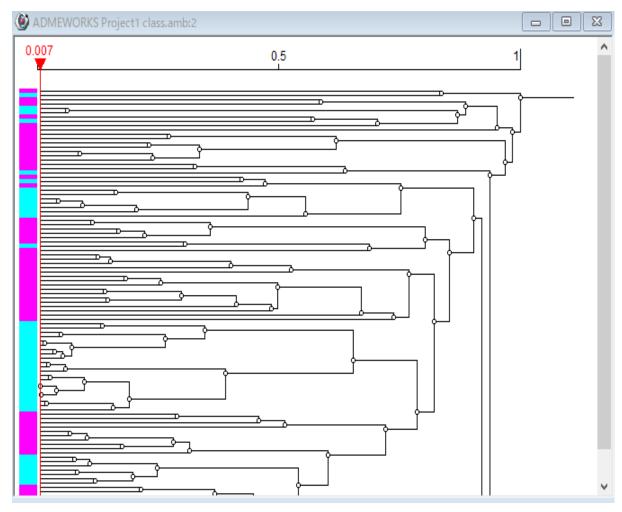
◇個別手法の特徴と適用内容:クラスター分析



- * クラスタリングはサンプル同士の相対的な位置関係や近隣関係を俯瞰して見ることが出来る
- * メトリックや融合手法の差異によりクラスタリング結果が大きく変化するので要因解析は注意

ModelBuilderの画面より









◇個別手法の特徴と適用内容:マッピング

■主成分分析(PCA)

- データ解析的にはリニアプロジェクション手法
- サンプル空間中のサンプルの位置関係を変えない
- オリジナルのサンプル空間を俯瞰する方向性を変えて空間を見直す手法
- •PCA適用前と後とで次元(パラメーター)数の変化はない

■非線形写像(NLM)

- ・データ解析的にはマッピング(写像)手法
- ・サンプルを人間が可視できる二次元/三次元上に分散
- ・サンプル空間を強制的に2/3次元に圧縮する
- ・非線形写像実施後は、サンプル空間の次元は2/3次元となる



◇個別手法の特徴と適用内容: Non Linear Mapping (NLM)

*NLMにより、最初のN次元空間が可視できる2次元空間に変換された。 *NLMでは、元のN次元空間における サンプル間の位置関係を保ちつつ 新しい2次元空間上に再配置される。 *2次元空間上の第6サンプルは、 その他のサンプル(1~5、7~9)との相 互位置関係(空間上の距離)をN次元空間上における関係と略同じとなる

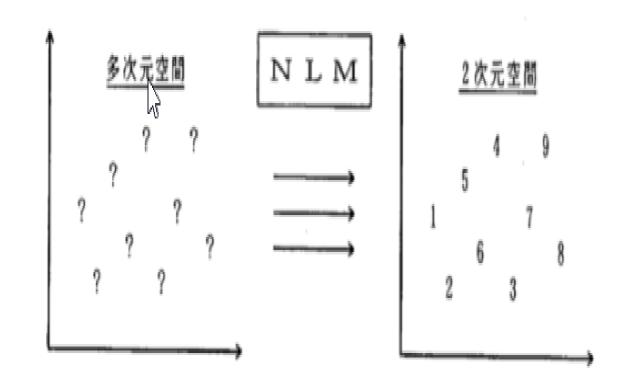


図5-6 ノンリニアマッピングによる多次元空間の2次元空間への写像



◇個別手法の特徴と適用内容: Non Linear Mapping (NLM)



◇個別手法の特徴と適用内容:

ニューラルネットワークによる次元圧縮

- * 入力層のN個のパラメーターから中間層のユニット数(図は3)に次元を変換する手法
- * NLMと異なり、サンプル間の相互 位置関係はキープされない
- * 新たな次元は、ニューラルネット ワークが解析目標とするクラスデータ 等を説明する情報を含んでいる

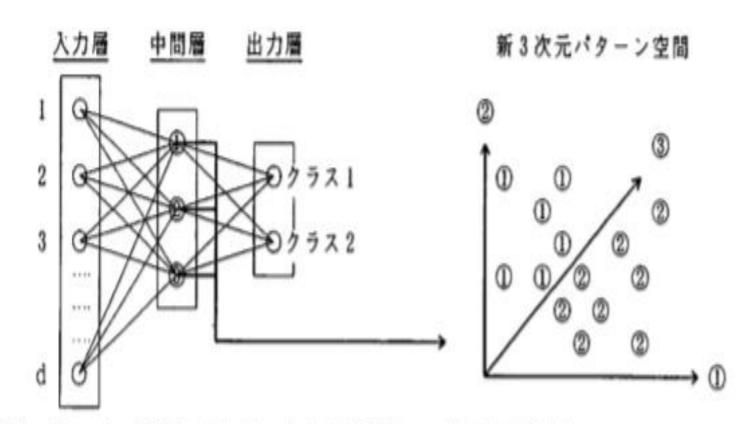


図 5 - 7. バックプロパゲーションによる写像。 d 次元 ⇒ 3 次元 (新たな次元は各パターンがクラス 1 と √ ラス 2 とに分類されやすいように 3 次元空間上に分布していることに注意)



□アンサンブル学習法

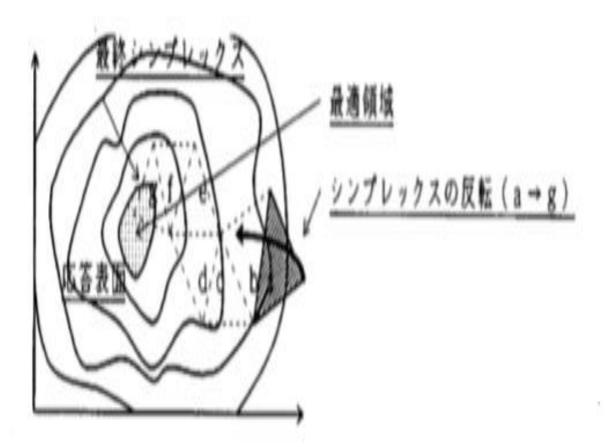
AdaBoostやランダムフォレスト等で採用されている解析手法 特徴:

弱分類機を条件を変えつつ多数用い、最終的な分類結果は個々の分類結果データを統合して最終結果とする。

従来の分類手法を組み合わせて使う「メタ解析手法」である。 *「**KY法**」もメタ解析手法となる



◇最適化手法:シンプレックス法



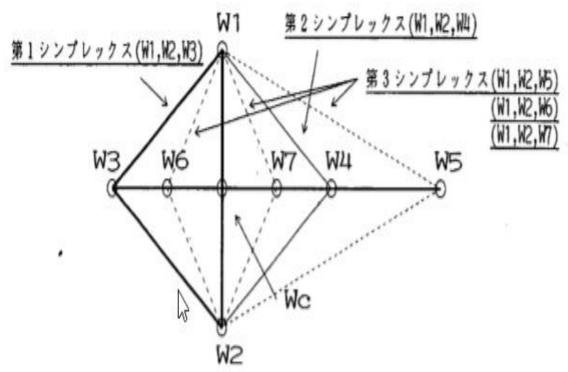
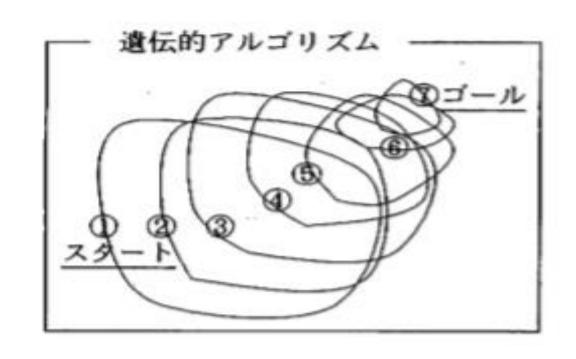


図6-2 シンプレックス反転に関するルール



◇最適化手法:遺伝的アルゴリズム

遺伝子の分裂/増殖/突然変異等の動きをシミュレーションするアプローチ



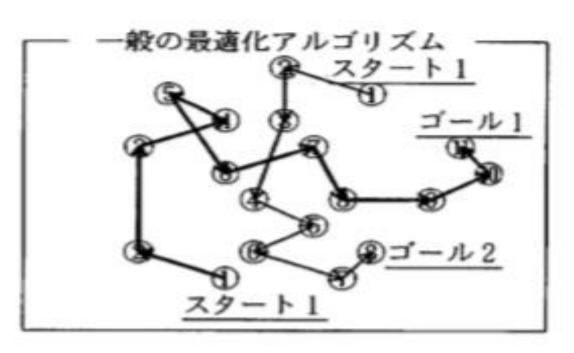


図 3. 遺伝的アルゴリズムと一般の最適化アルゴリズムによる最適領域の探索過程

- 1. 増殖 (MULTIPLICATION)
- 2. 交叉 (CROSSING-OVER)
- 次然変異(MUTATION)
- 4. 淘汰/選択 (SELECTION)



◇最適化手法:遺伝的アルゴリズム

- 1. 増殖 (MULTIPLICATION)
- 交叉(CROSSING-OVER) 情報の取替により、情報の変換を果たすものである。

10110100101110110010010111 遺伝子A

10011100110000110010010111 新遺伝子BA

10110100101111010101110100011 新遺伝子AB

10011100110000101110100011 遺伝子B

異なる遺伝子パターンをもつ2本の遺伝子がある時、これら2本の遺伝子を交叉させる。 この結果、新たに2本の遺伝子が誕生するがこれらの遺伝子は親となるAおよびBの遺伝子の形質を交叉させた点を中心としてそれぞれあわせ持つことがわかる。

3. 突然変異 (MUTATION) この突然変異では遺伝子の持つ情報がある部位で変化して対立遺伝子となる。

突然変異

10110100101難10110010010111 遺伝子A'

4. 淘汰/選択 (SELECTION)

様々な変換パターンにより構築された遺伝子が淘汰により悪い形質(問題解決に取って望ましくない)を持つものが取り除かれてゆくことである。 生物学的にはまわりの環境に適合した形質を持つ固体だけが生き残り(選択され)、次世代へと情報(形質)を伝えてゆくことを意味する。

								淘汰	
1	0	1	(********	1	1	1	遺伝子A		
0	1	0		1	1	0	遺伝子B		
0	0	1		0	0	1	遺伝子C	\longrightarrow \times	
*		***			***				
1	0	0		0	1	1	遺伝子乙	\longrightarrow \times	

◇最適化手法:ニューラルネットワーク



閉じたネットワーク構造を有するニューラルネットワーク ボルツマンマシンおよびホプフィールドネット

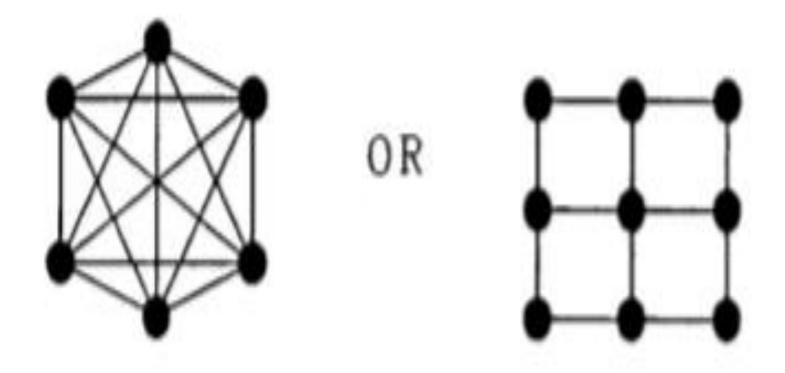


図1. ボルツマンマシン及びホップフィールドネットのネットワーク構造



- □本日の解説内容と順番
- 1. AI (人工知能) に関する概要
- 2. データサイエンスで利用される化学パラメータ
- 3. データサイエンス手法の概要
- 4. KY法の説明



□湯田が開発した新規データ開発手法:

KY法 (K-step Yard sampling methods)

KY法の基本技術:

多段階、リサンプリング、アンサンブル法をミックスして完成した最強力な**ニクラス分類手法**

- * 当初は二クラス分類手法のみであったが、現在はフィッテイング (重回帰)手法にも展開可能となっている。
- *判別関数や回帰式の作成様式の違いにより、二クラス分類および フィッテイング手法それぞれに3種類のKY法が創出された。
- ①2モデルKY法
- ②1モデルKY法
- ③モデルフリーKY法



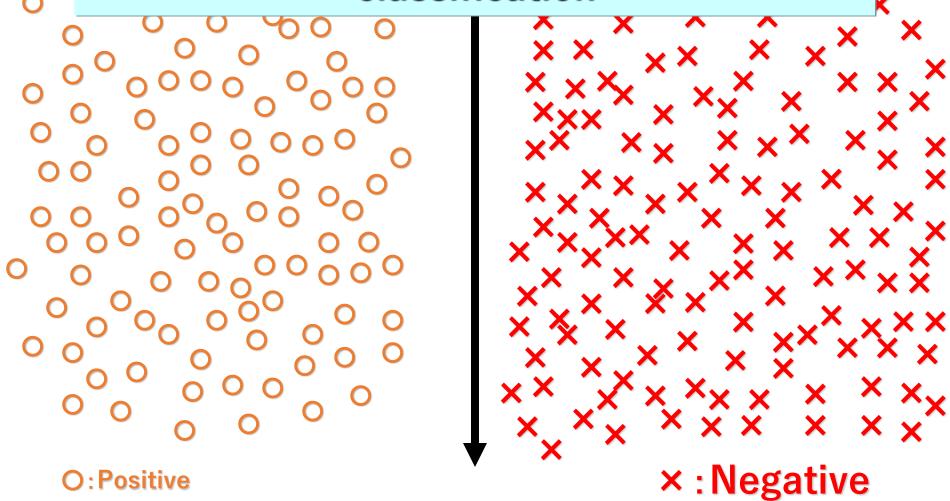
Development of "K-step Yard sampling method" and Apply to the ADME-T In Silico Screening

Kohtaro Yuta

Sample space: two cluster samples Int



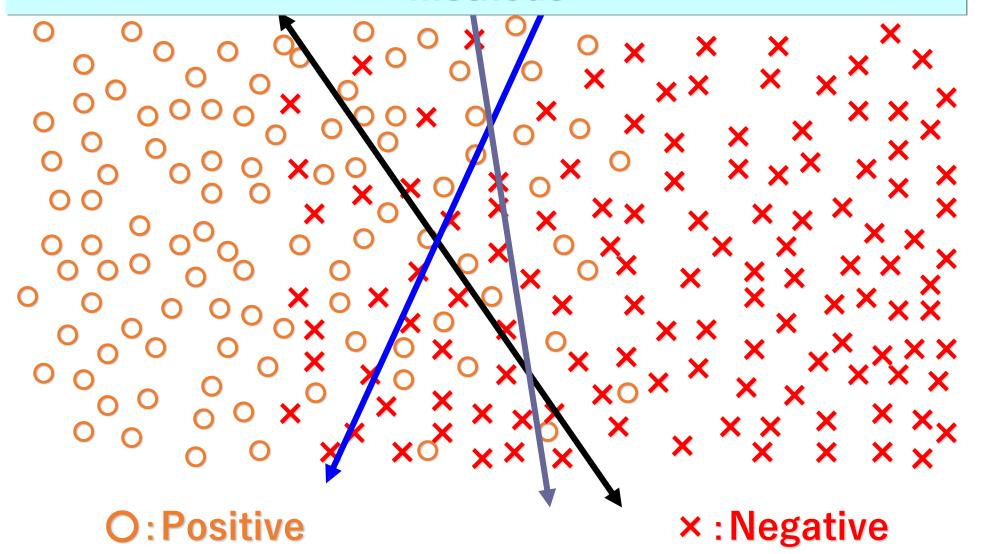
Discriminant function for perfect classification



Sample space : highly overlapped space



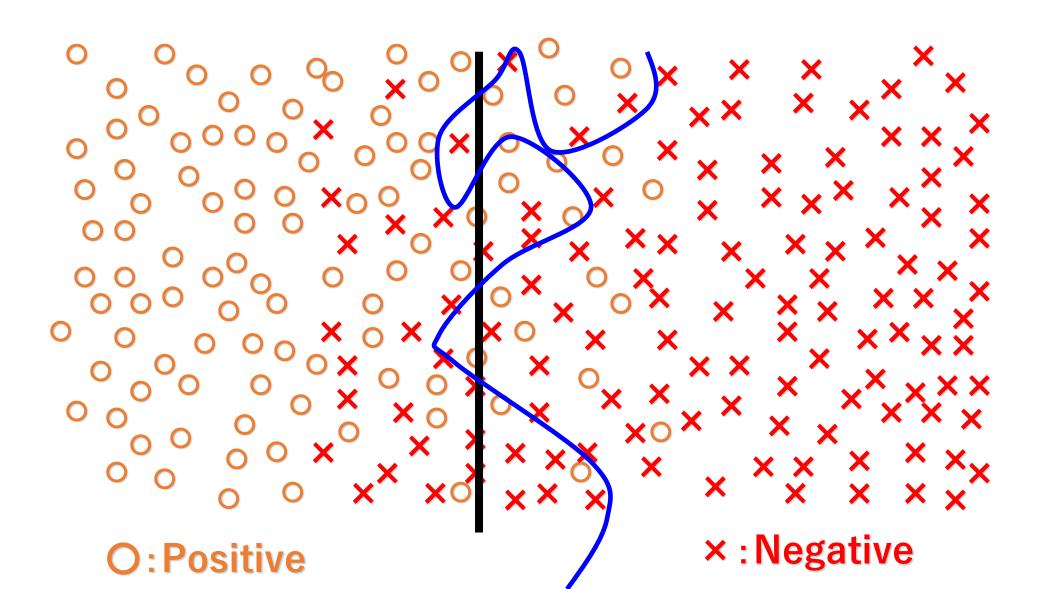
Discriminant function generated by various methods



Sample space : highly overlapped space



Discriminant function: Linear and non-linear



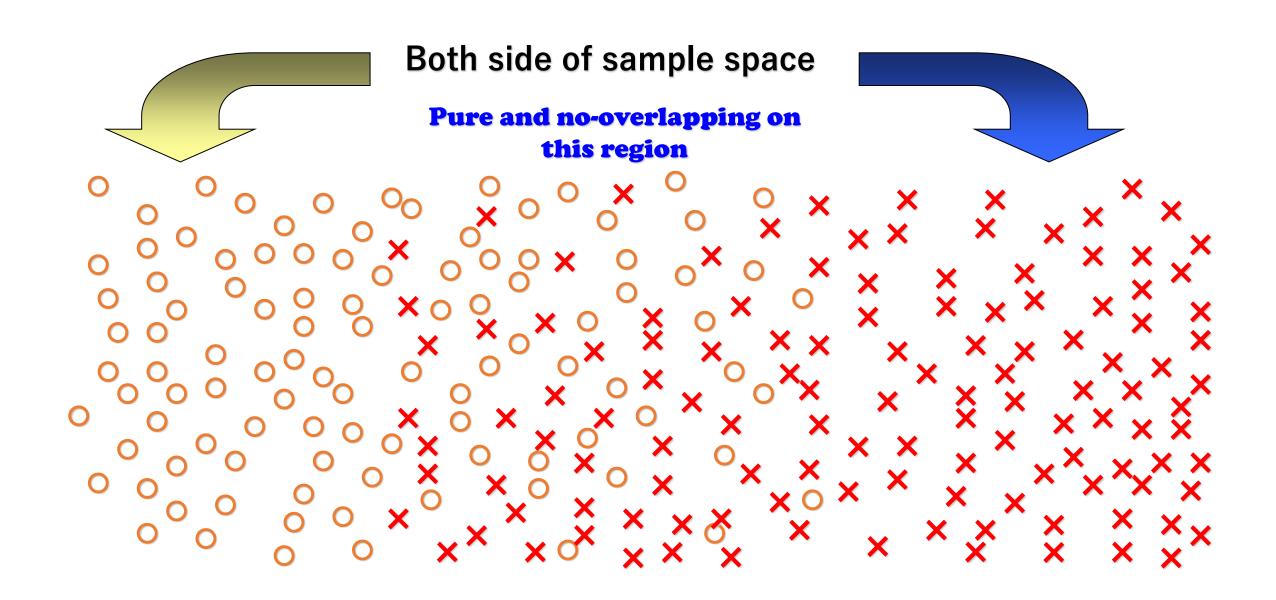


Building process to the features of "K-step Yard sampling method"

Step1: Yard sampling methods

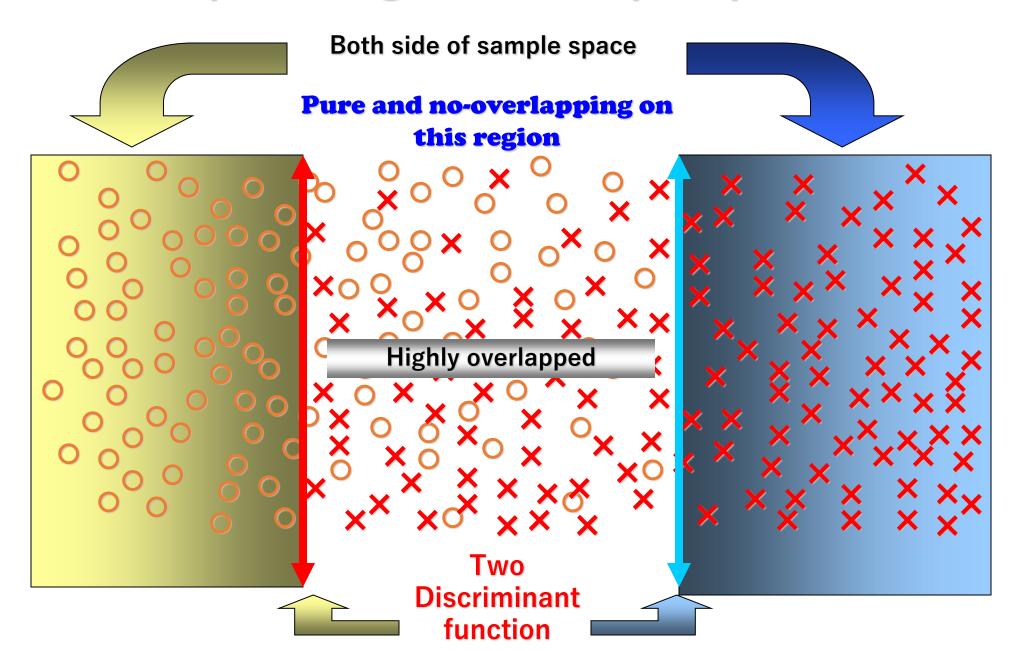
Spatial region on sample space





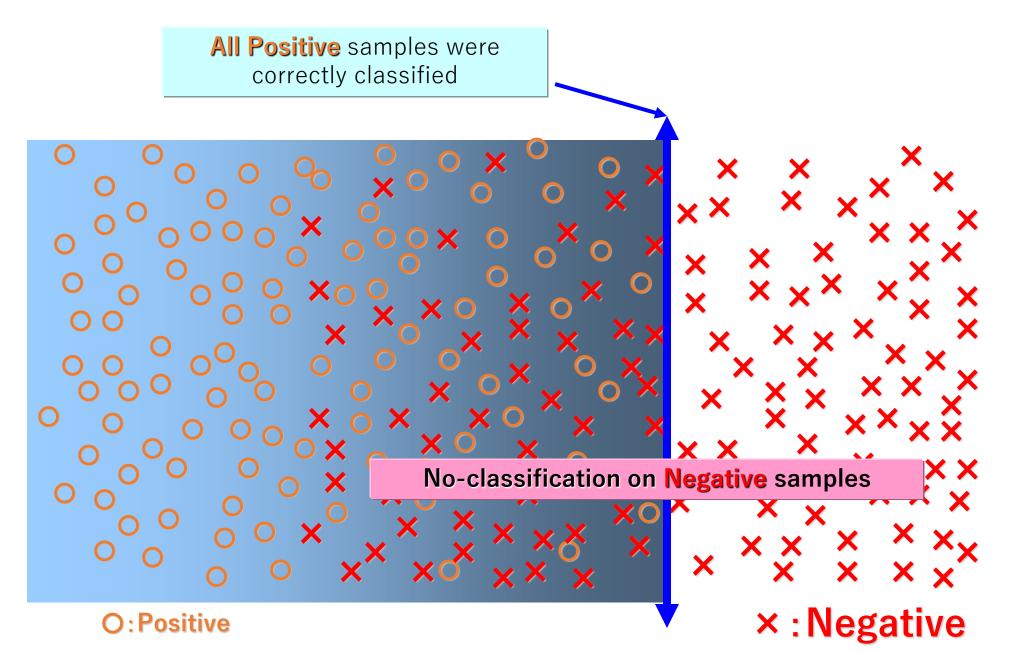
Spatial region on sample space





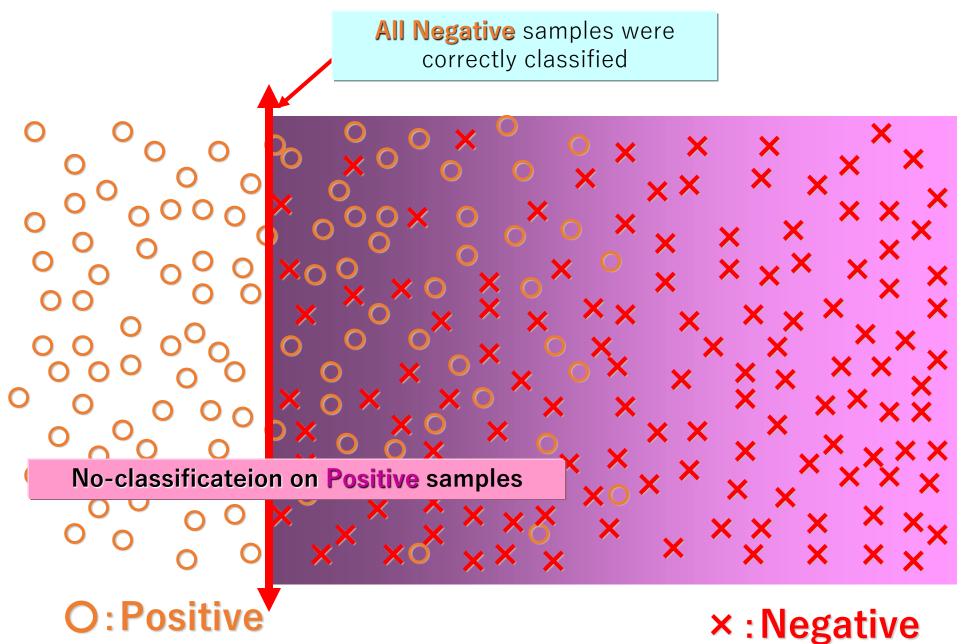
Property of AP (All Positive) model





Property of AN (All Negative) model





Combination of AN and AP models



Not to be classified High reliability High reliability **Positive** Negative Grey Zone **Zone** Zone

Linear and non-linear discriminant on AP and AN models



High reliability Not to be classified High reliability **Positive** Negative Grey Zone Zone Zone

Relations between Sample space & AN and AP models



Determination No conclusion **Determination** Positive AN model -Negative Positive — AP model Negative **Positive Negative** Grey Zone **Zone** Zone

Class determination by AN and AP models



☐ Sample Classification and prediction must be done by Combination of the results of AP and AN models.

AP model	AN model	Results
①AP; POSI,	AN; POSI	
②AP; POSI, ③AP; NEGA,	AN; NEGA AN; POSI	G R E Y G R E Y
4AP; NEGA,	AN; NEGA	N E G A



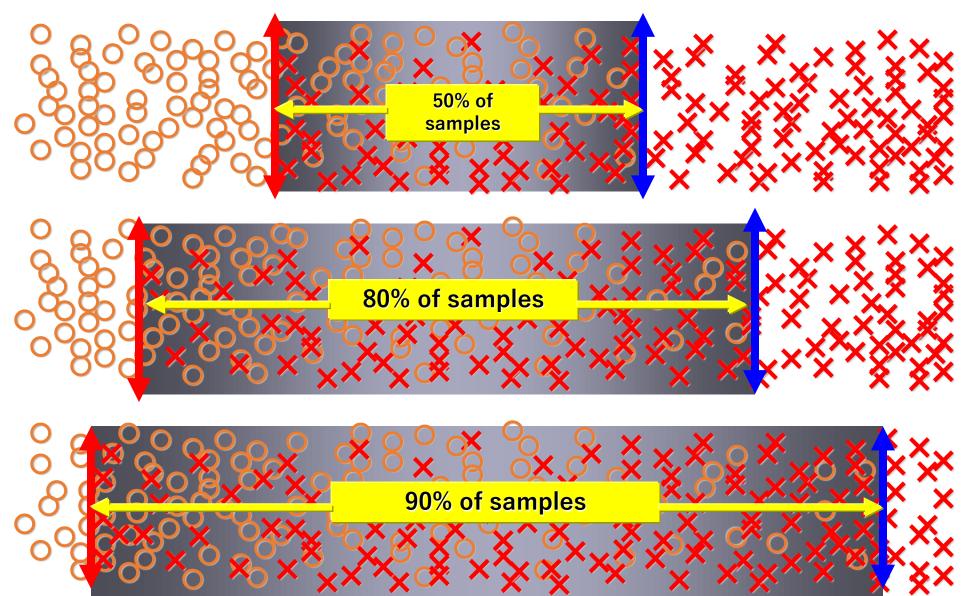
Building steps to the features of "K-step Yard sampling method"

Step2: K-step approach

Problems of Yard sampling methods

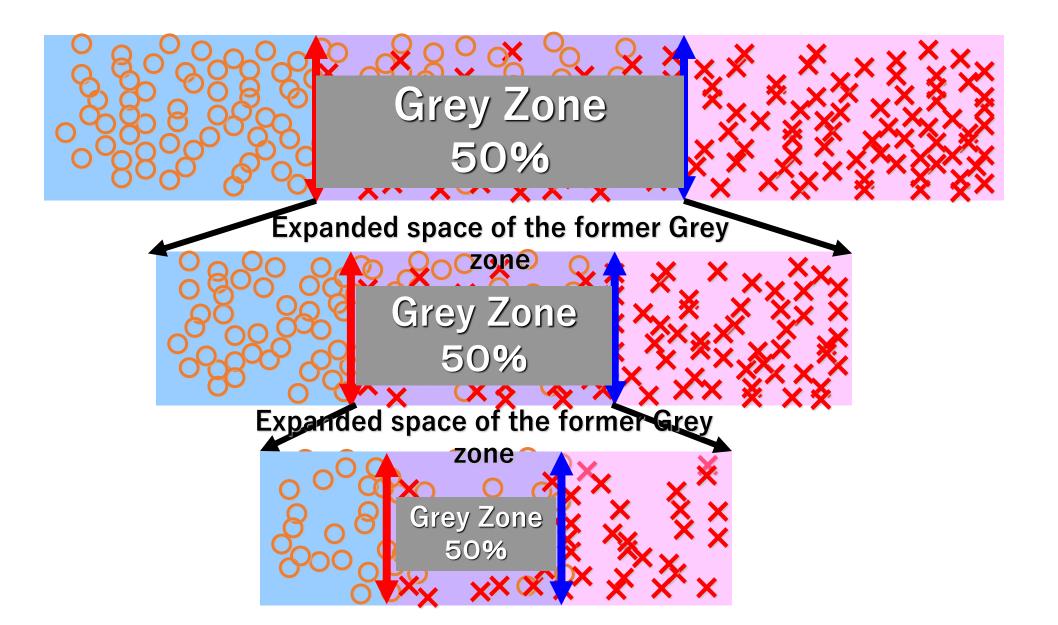


The ratio of Grey zone: Highly overlapped sample space





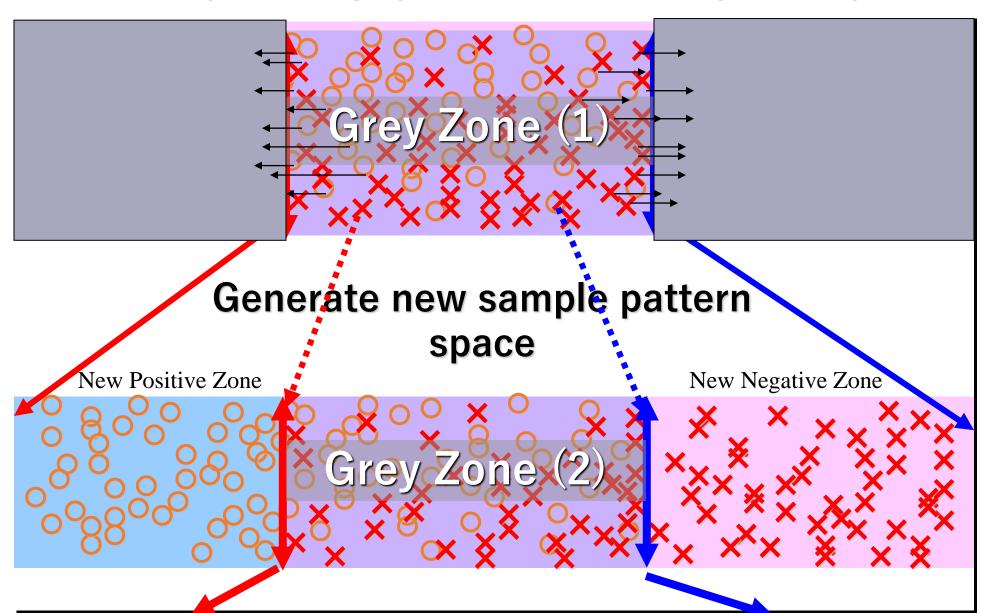
Steps to the K-step methods



"K-step Yard sampling (KY) Method"



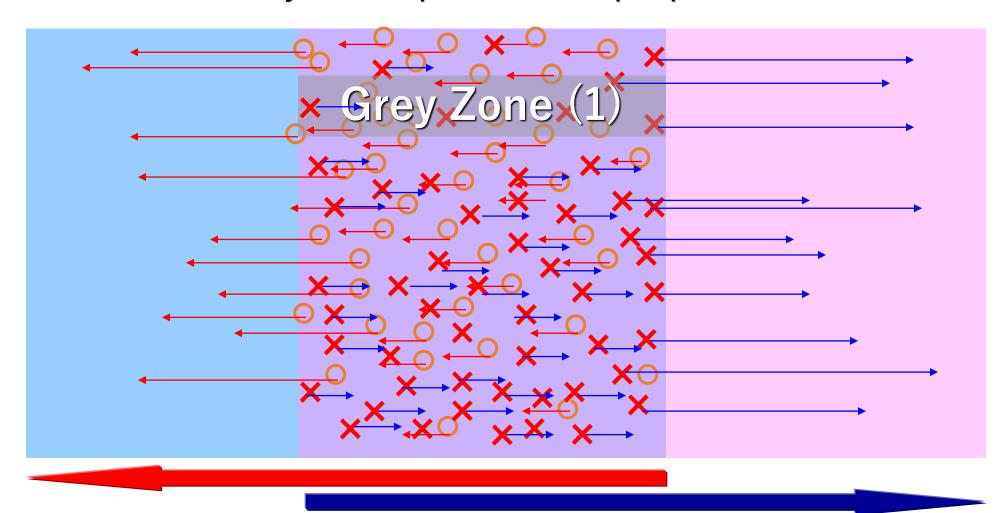
Improvement by repeated classification of Grey Zone samples



"K-step Yard sampling (KY) Method"



☐ Relocation of Grey Zone samples on new sample space





Applicability statement of "K-step Yard sampling method"

Classifying 7000 sample set of Ames test

Challenge for classification and prediction



K-step Yard sampling method KY-method



The most powerful and advanced data analysis method



The most difficult classification problem

6,965 sample of Ames test samples were,

Classified perfectly

Application test of "K-step Yard sampling"



Samples □ Samples □ Sam

- 1. Ames test data
- 2. Sample population

total:6,965

Mutagen; 2,932

Non-mutagen; 4,033

□ Result of KY-method

- 1. Number of steps : 23 steps ; 22 (2 models) + 1 (1 model)
- 2. Classification ratio: 100 %

□Used system

ADMEWORKS / ModelBuilder V 3.0.22

□ Used parameters (Initial condition)

Number of generated parameters: 838

Number of parameters for step 1:98

Confidency index (Samples (6965) / Parameters (98)): 71.1 > 4.0

Application test results by various D.A. methods

1. Linear discriminant analysis with linear least-squares method

```
Classification ratio : total; 73.50(6965), Mutagen;73.02(2932), Non mutagen;73.84(4033) 
Number of mis-classified : (1846), (791) (1055)
```

```
Prediction ratio (L100 out) 72.58% deviance(0.92%) (L500 out) 73.32% deviance(0.18%)
```

2. SVM (Support Vector Machine with Kernel)

```
Classification ratio : total; 90.87(6965), Mutagen;86.83(2932) Non mutagen; 93.80(4033) Number of mis-classified : (636), (386) (250) Prediction ratio (L500 out) 80.99% deviance(9.88%)
```

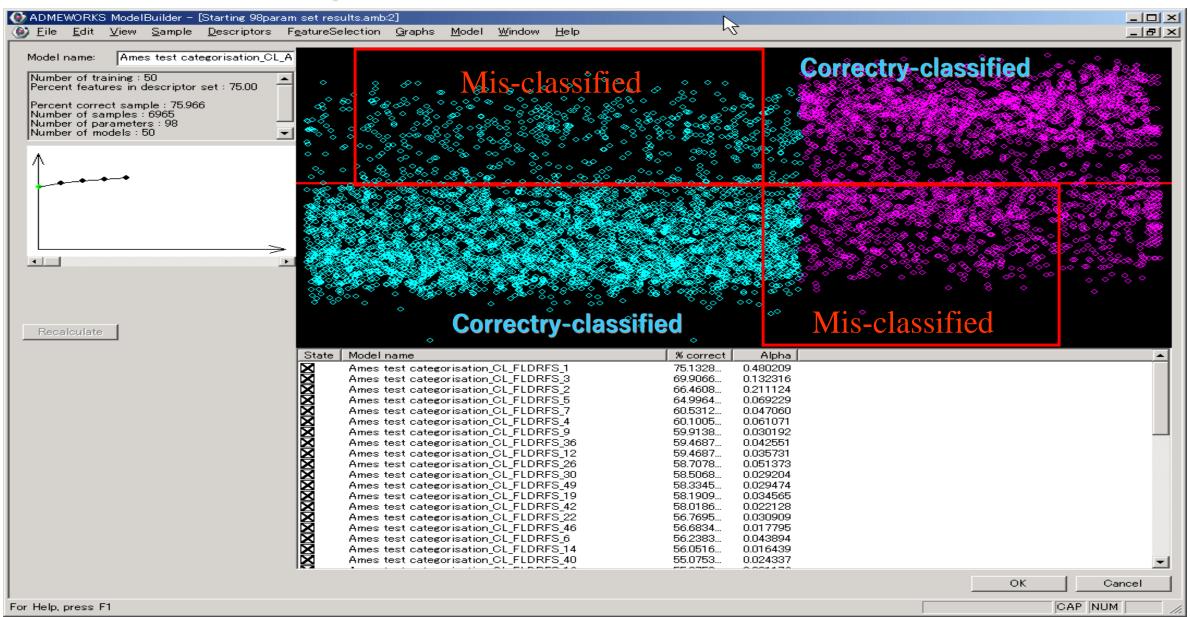
3. AdaBoost

```
Classification ratio : total; 77.24(6965), Mutagen;66.13(2932) Non-mutagen; 85.32(4033) Number of mis-classified : (1585) (993) (592) Prediction ratio (L500 out) 75.16% deviance(2.08%)
```

Classification results by AdaBoost



Sample distribution of 6,965 of 77.24%



"K-step Yard sampling (KY) Method"



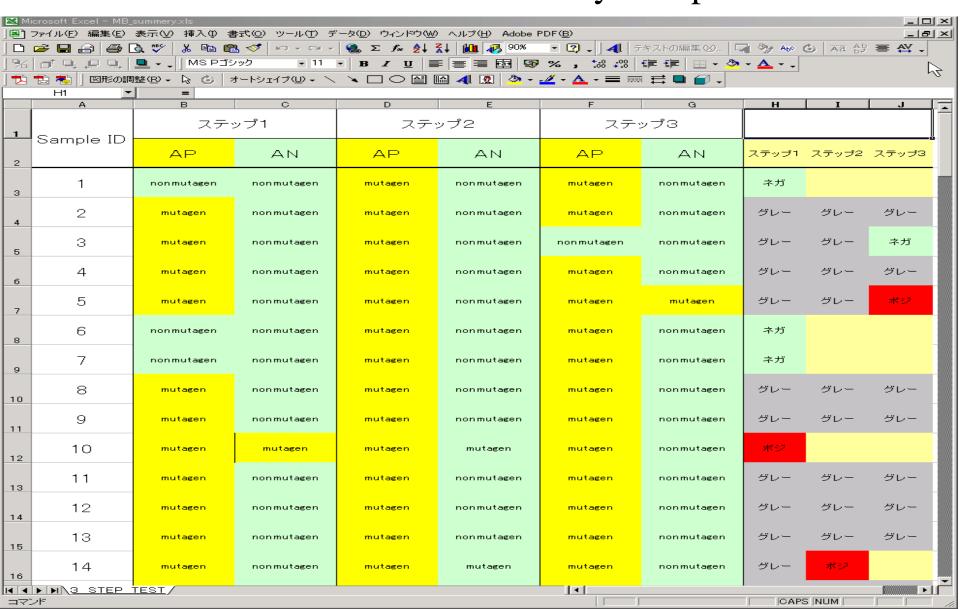
Total steps: 23 steps (2 models) + 1 step (1 model)

ステップID(KY法)	Starting samples(Total)	Mutagen (Initial)	Non-mutagen (Initial)	Grey sample (Initial)
	Final samples	Mutagen (Final)	Non-mutagen (Final)	Grey sample (Final)
	Determined samples(Total)	Determined samples(Mut.)	termined samples(Non-mu	Grey ratio(%) (Grey/Total)
1	6965	2932	4033	
	5864	2413		586
	1101	519	582	84.19
2	5864	2413		586
	5108	2142	2966	510
	756	271	485	87.1
3	5108	2142	2966	510
	4486	1919	2567	448
	622	223	399	87.83
4	4486	1919	2567	4480
	4133	1779	2354	413
	353	140	213	92.1:
5	4133	1779	2354	413:
	3794	1651	2143	379
	339	128	211	91.6
6	3794 3462	1651	2143 1977	379 ⁴
	332	1485 166	166	91.2
	3462	1485	1977	346
7	3090	1345	1745	340.
	372	140	232	89.29
8	3090	1345	1745	3090
	2826	1220		2820
	264	125	139	91.40
9	2826	1220	1606	282
	2592	1139	1453	259
	234	81	153	90.6
10	2592	1139	1453	259
	2384	1047	1337	238
	208	92	116	91.9

"K-step Yard sampling (KY) Method"



Classification results by 3 steps



Spatial features of "K-step Yard sampling"



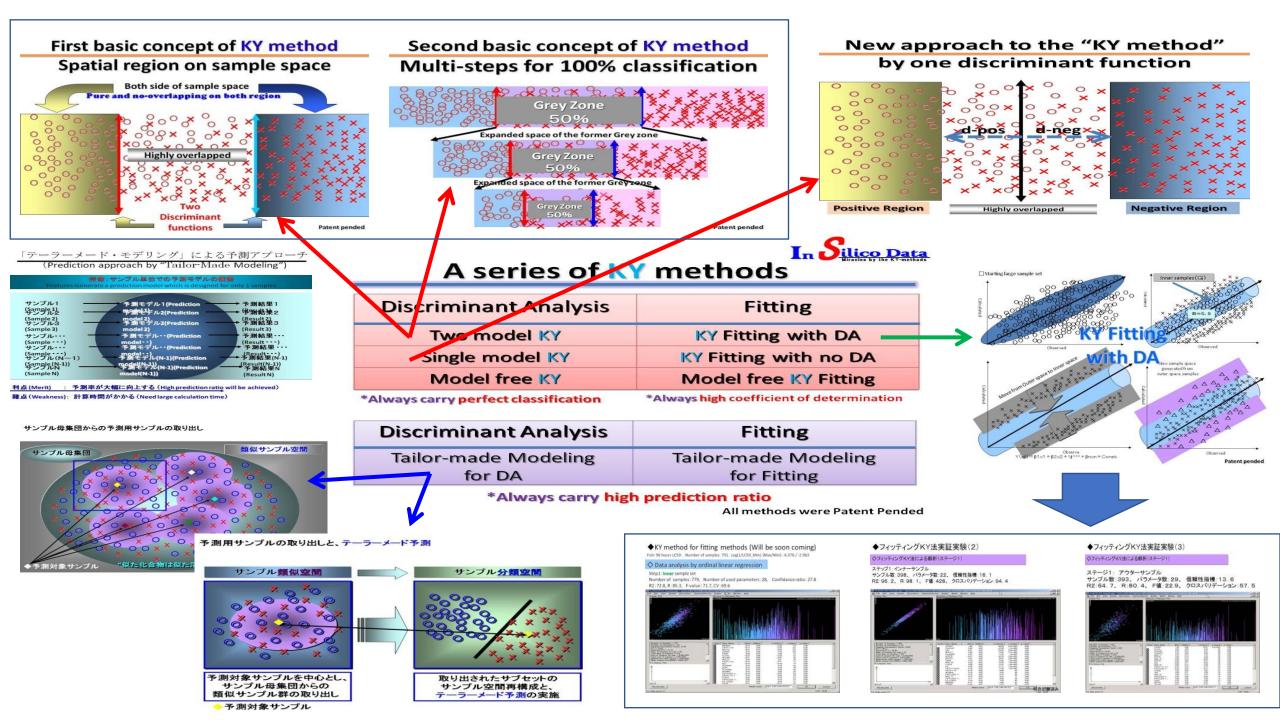
□Summary

Advantages

- 1. Sample number free approach
- 2. Sample distribution free approach
- 3. Perfect classification is achieved in any condition

■ Disadvantages

- 1. Relatively complex operation to generate discriminant functions
- 2. Need powerful computer power



TS-02 第四章 適用留意事項に続く