# Leading-edge
# Classification / Prediction Methods
# for applied on
# Toxicity Prediction Field

**Kohtaro Yuta**
**In Silico Data Ltd.**
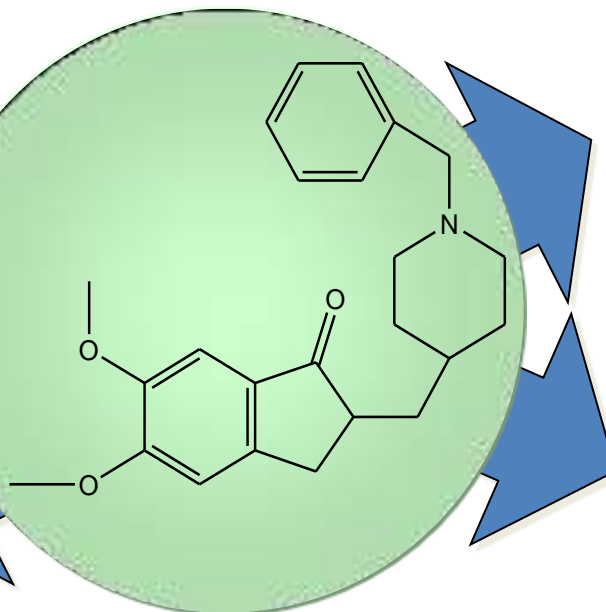**E-Mail : k-yuta@insilicodata.com**
**http://www.insilicodata.com**

# Drug properties and compound structure



**ADME properties**

**Pharmacological activity**

**Physicochemical properties**

**Toxicity**

All properties are fixed when the structure is determined.

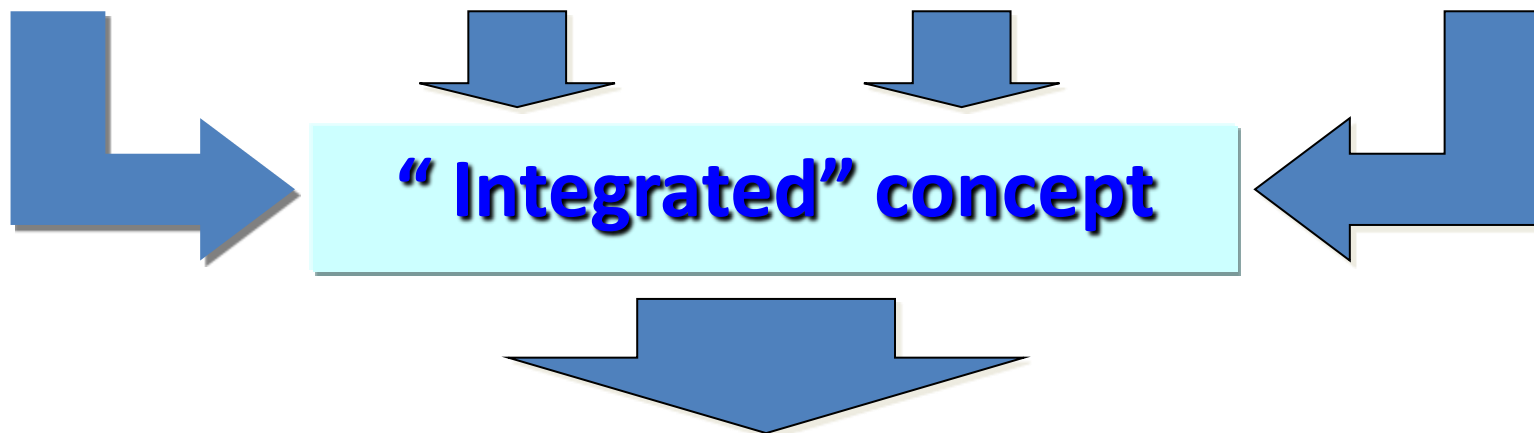There are no relations between any two properties.

All properties must be optimized for developing drugs.

# Why PR(Pattern Recognition) for toxicity screening

**Toxicity**

1. Unexplained & complex mechanisms
2. High structural diversity

CADD

Multi-variate Pattern recognition

Artificial intelligence

Factorial analysis

Black box

Know-how

Prediction

# Normal sample space : small overlapped space

# Toxicity sample space : large overlapped space

50% of samples
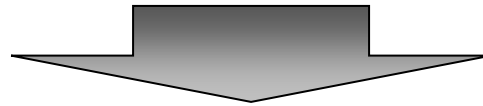
80% of samples

90% of samples

# Classification Result by AdaBoost
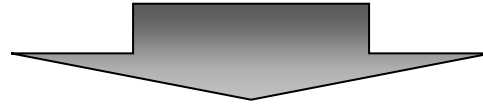
## 77.24% of Ames test 6,965 samples

# Perfect(100%) classification of

## Ames test 6965 pos/neg sample set

### K-step Yard sampling method
### KY-method

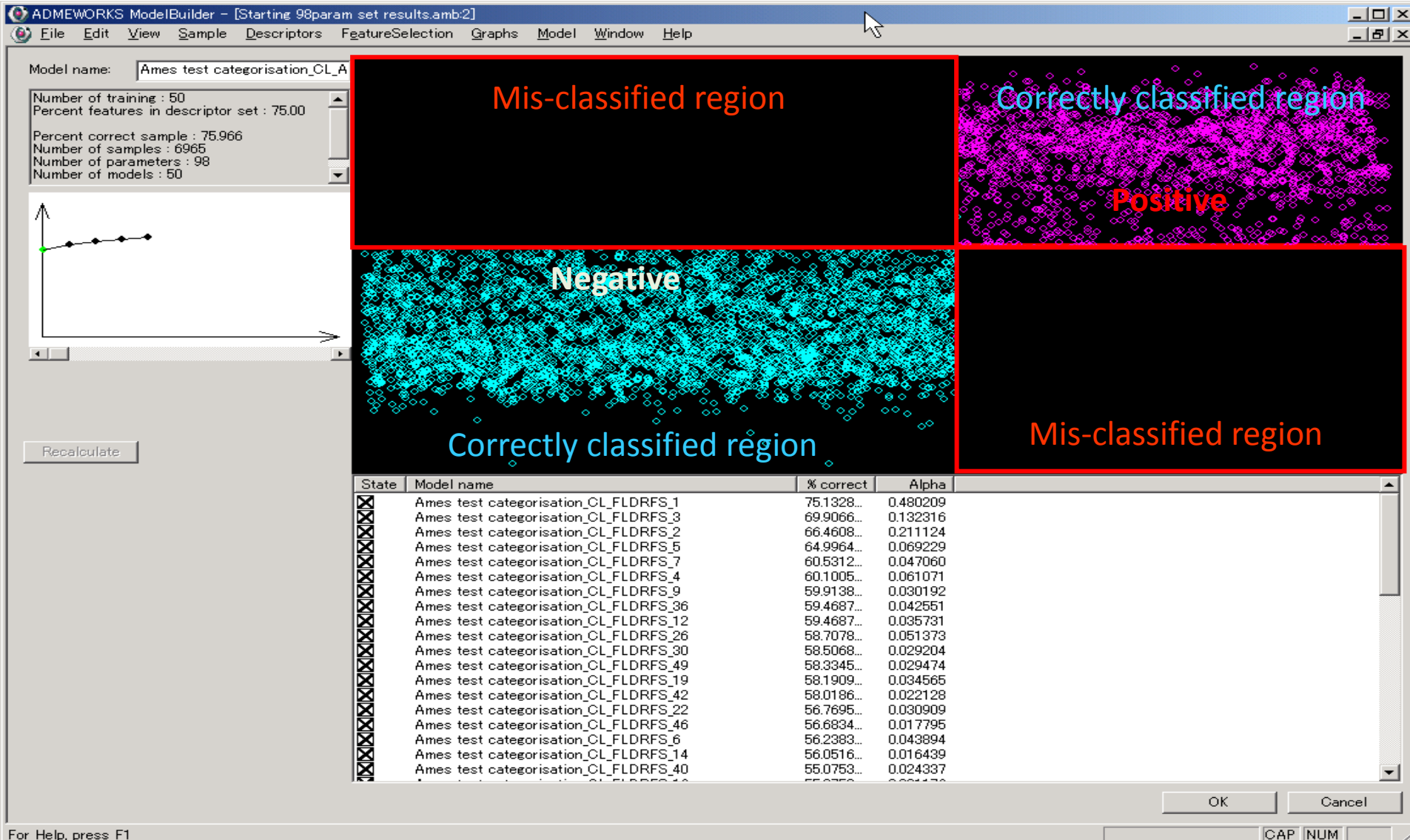**The most powerful and advanced data analysis method**

**The most difficult classification problem**
**6,965** sample of Ames test were,
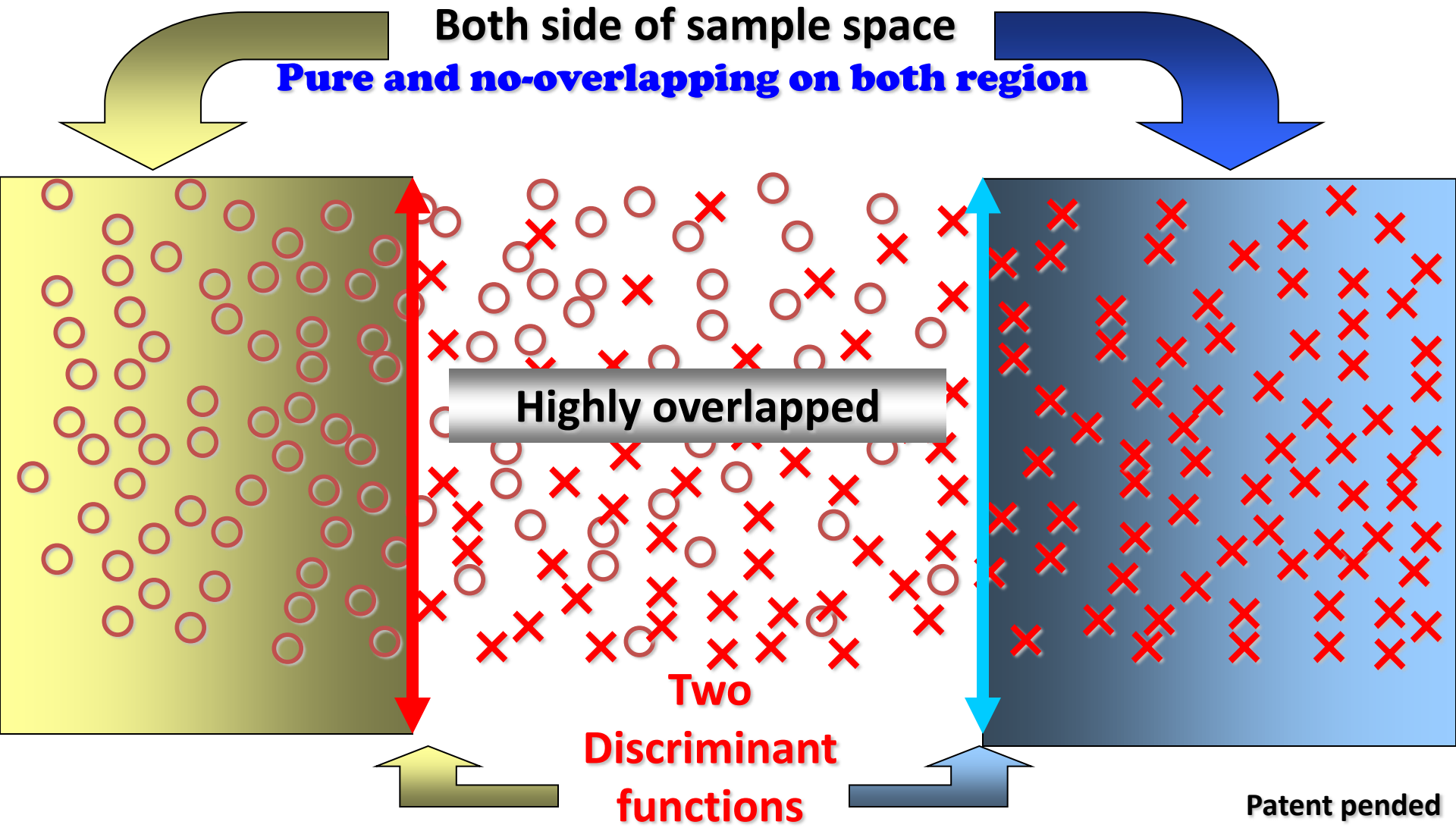**Classified perfectly**

# 100% correctly classified of Ames test 6,965 samples

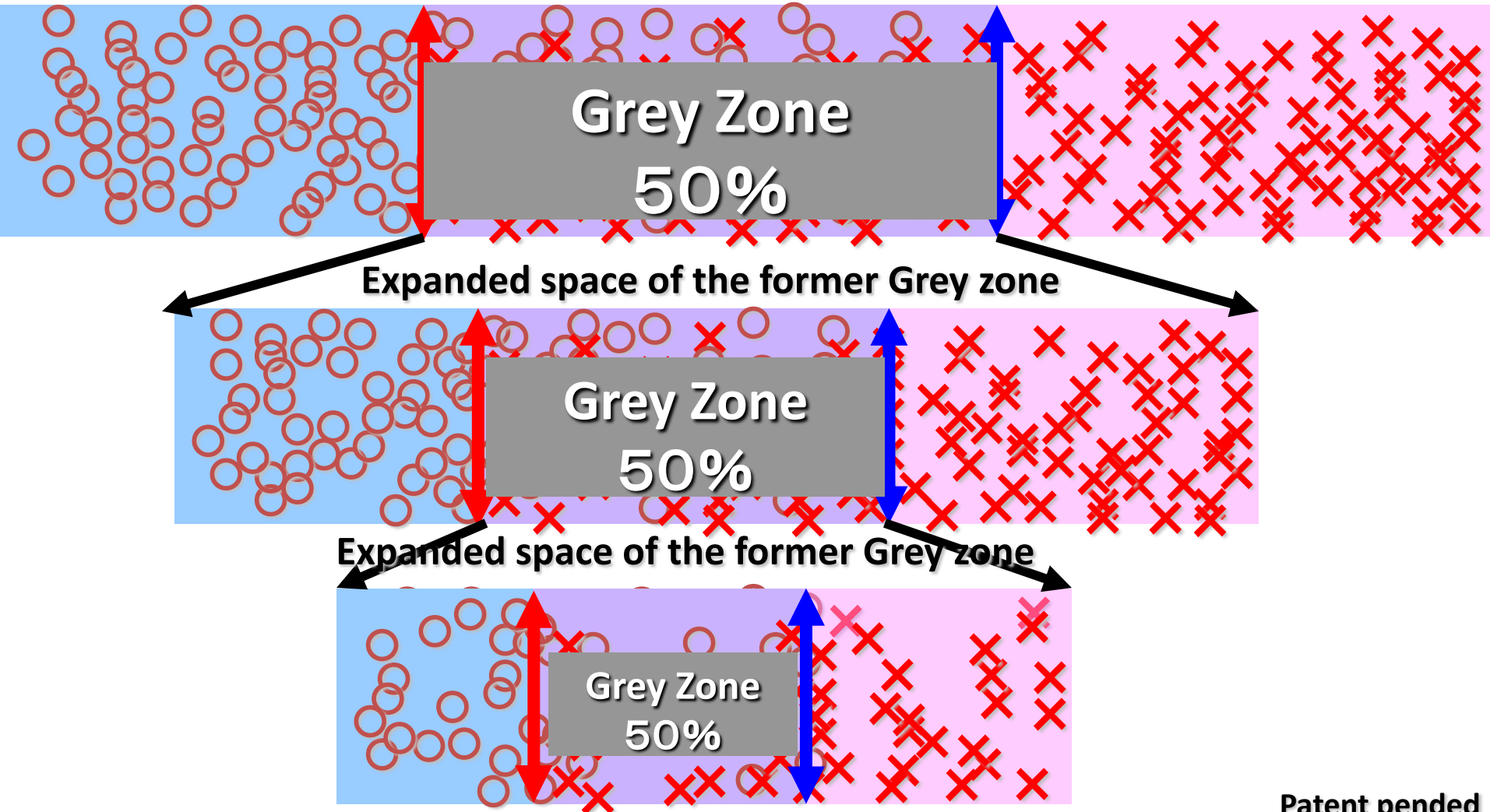## Artificial Display Image of Perfect Classification

# First basic concept of KY method
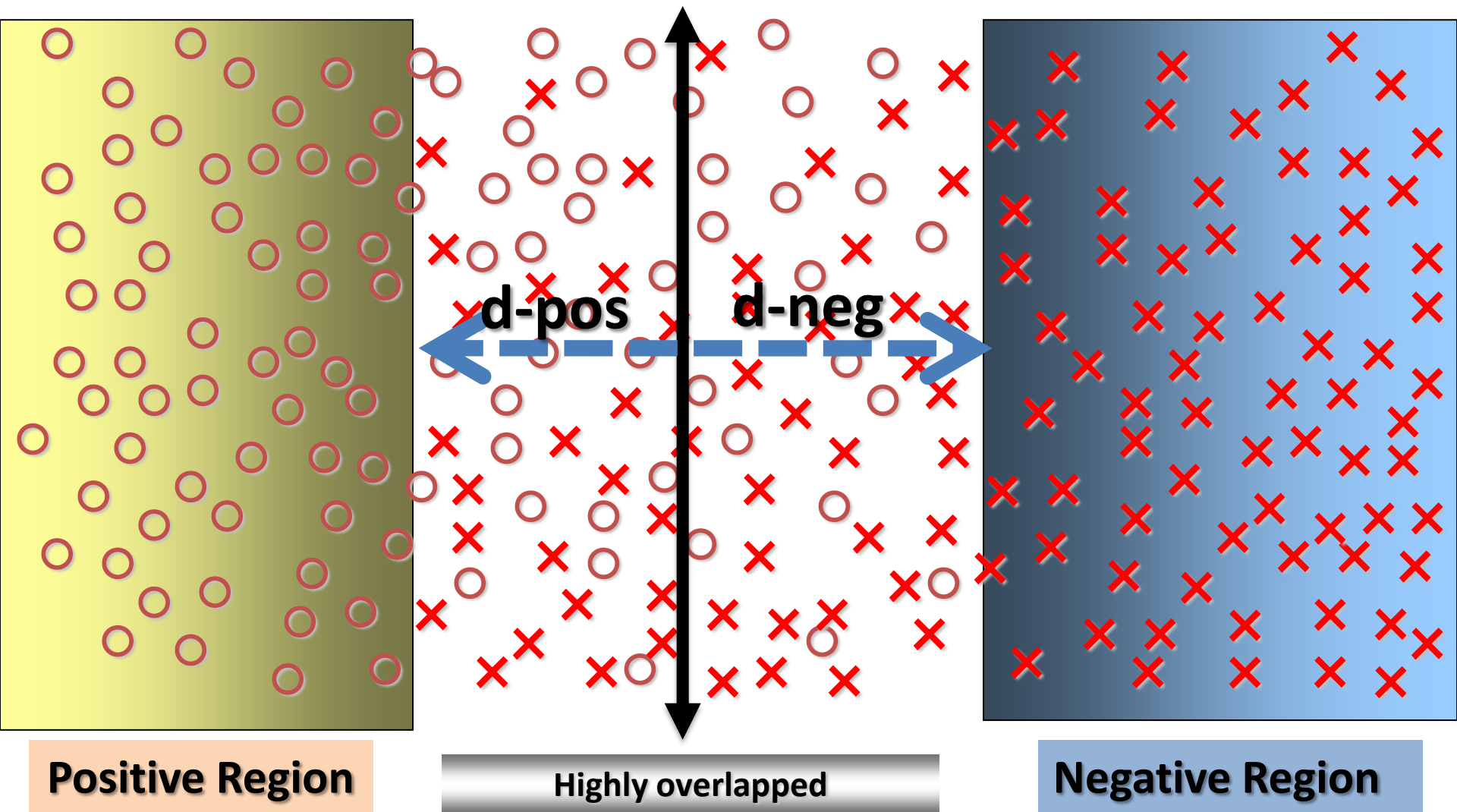
## Spatial region on sample space

Both side of sample space

Pure and no-overlapping on both region

Highly overlapped

Two Discriminant functions

Patent pended

# Second basic concept of KY method

# Multi-steps for 100% classification

Grey Zone 50%

**Expanded space of the former Grey zone**

Grey Zone 50%

**Expanded space of the former Grey zone**

Grey Zone 50%

# New approach to the "KY method" by one discriminant function



d-pos    d-neg

Positive Region

Highly overlapped

Negative Region

# A series of KY methods

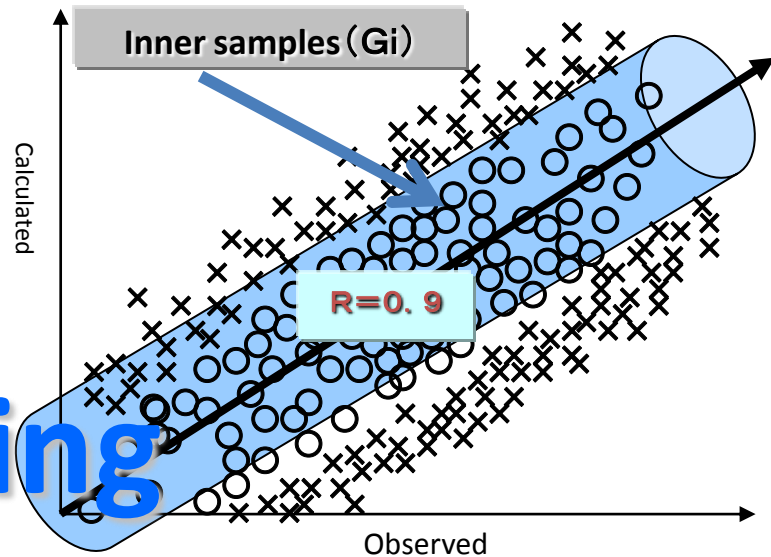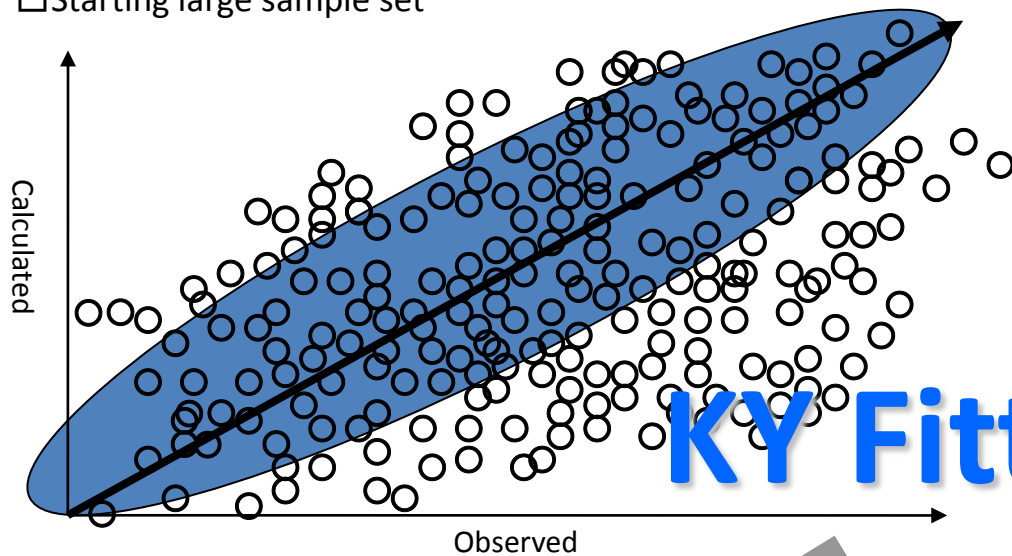| Discriminant Analysis | Fitting |
|---|---|
| Two model KY | KY Fitting with DA |
| Single model KY | KY Fitting with no DA |
| Model free KY | Model free KY Fitting |

*Always carry perfect classification

*Always high coefficient of determination

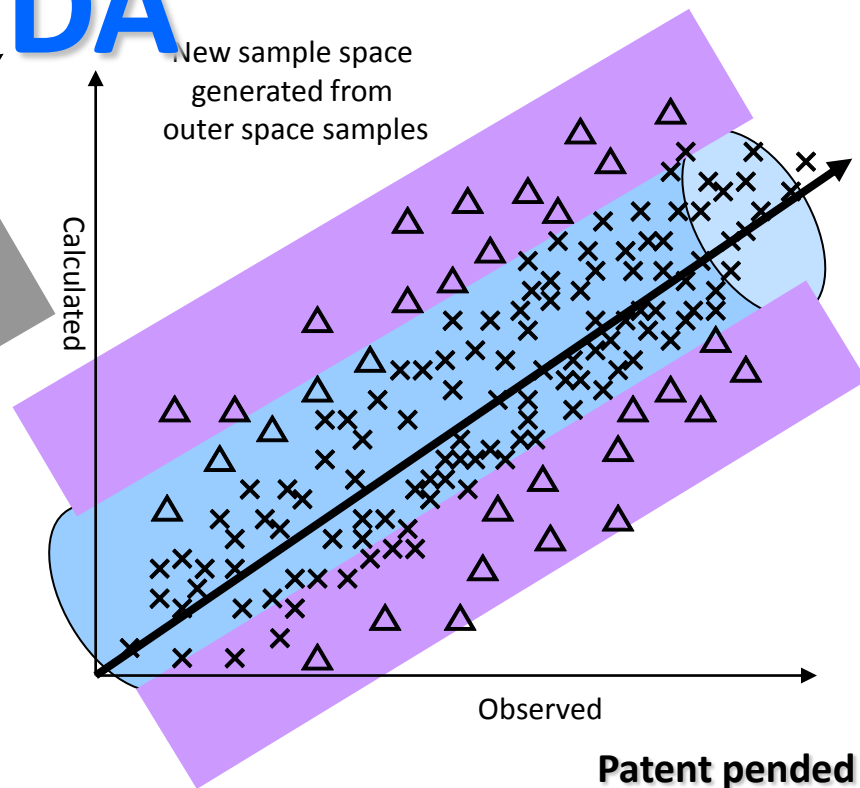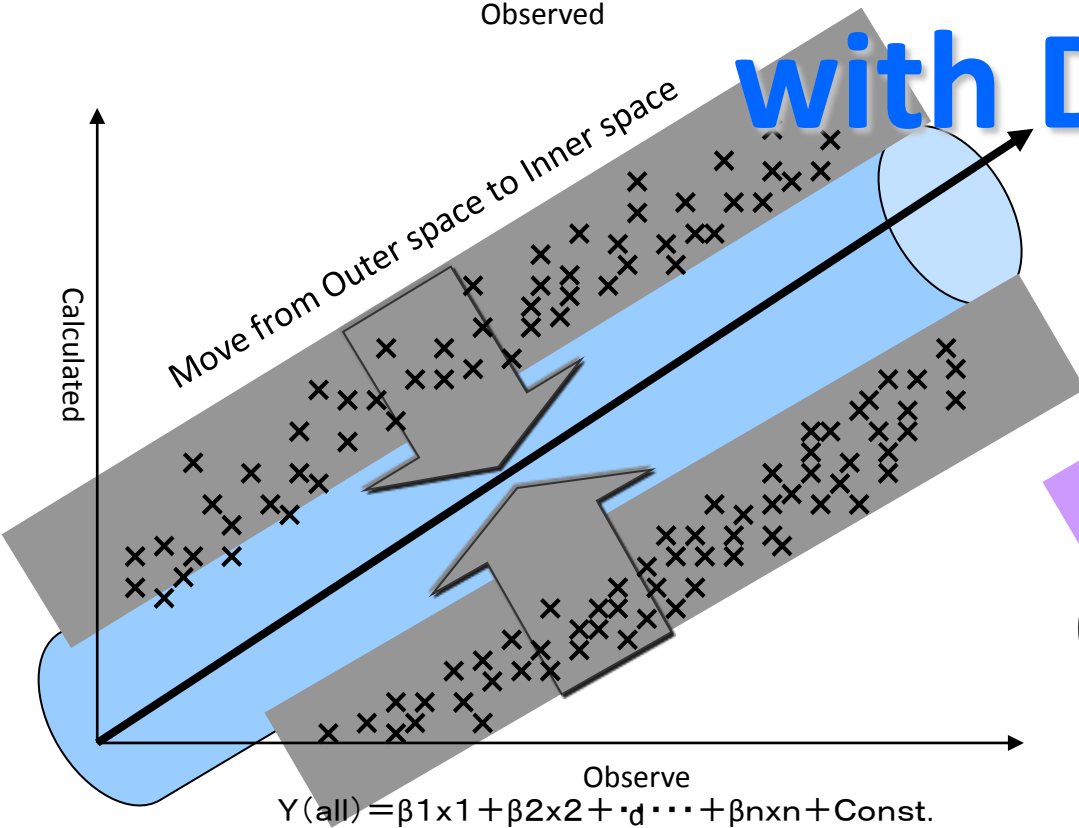| Discriminant Analysis | Fitting |
|---|---|
| Tailor-made Modeling for DA | Tailor-made Modeling for Fitting |

*Always carry high prediction ratio

All methods were Patent Pended

□Starting large sample set



**KY Fitting with DA**

Inner samples（Gi）

R＝0．9

Move from Outer space to Inner space

Observe

$$Y(all) = \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_n x_n + Const.$$

New sample space generated from outer space samples

Observed

Observed

Observed

Calculated

Calculated

Calculated

Calculated

**Patent pended**

# ◆KY method for fitting methods (Will be soon coming)

Fish: 96 hours LC50、Number of samples: 791、Log(1/LC50_Mm) (Max/Min) : 6.376 / -2.963

## ◇ Data analysis by ordinal linear regression

Step1：**Inner** sample set

Number of samples：779,　Number of used parameters：28,　　Confidance ratio：27.8

R2：72.8, R：85.3,　F-value：71.7, CV：69.6

# ◆フィッティングKY法実証実験（２）

ステップ１：インナーサンプル
サンプル数：398、　パラメータ数：22、　信頼性指標：18．1
R2：96．2、　R：98．1、　F値：428、　クロスバリデーション：94．4



特許出願済み

# ◆フィッティングKY法実証実験（３）

ステージ1： アウターサンプル
サンプル数：393、 パラメータ数：29、 信頼性指標：13．6
R2：64．7、 R：80．4、 F値：22.9、 クロスバリデーション：57．5

# 予測用サンプルの取り出しと、テーラーメード予測



サンプル**類似空間**

サンプル**分類空間**

予測対象サンプルを中心とし、
サンプル母集団からの
類似サンプル群の取り出し

取り出されたサブセットの
サンプル空間再構成と、
テーラーメード予測の実施

◆予測対象サンプル

サンプル母集団からの予測用サンプルの取り出し

# 従来手法による予測アプローチ
## (Prediction approach by traditional method)

**特徴：総てのサンプルを対象とした予測モデルの構築**
Features:Generate a prediction model which can handle all samples

サンプル1
(Sample 1)
サンプル2
(Sample 2)
サンプル3
(Sample 3)
サンプル・・・
(Sample ・・・)
サンプル・・・
(Sample ・・・)
サンプル（N－1）
(Sample (N-1))
サンプルN
(Sample N)

**予測
モデル
(Prediction
Model)**

予測結果 1
(Result 1)
予測結果2
(Result 2)
予測結果3
(Result 3)
予測結果・・・
(Result ・・・)
予測結果・・・
(Result・・・)
予測結果(N-1)
(Result(N-1))
予測結果 N (Result N)

**利点 (Merit)　　：　少ない数の予測モデル作成で済む（Small number of prediction models are generated ）**

**欠点(Weakness)：予測率の向上が困難である（Difficult to achieve high prediction ratio）**

# 「テーラーメード・モデリング」による予測アプローチ
## （Prediction approach by "Tailor-Made Modeling")

**特徴：サンプル単位での予測モデルの構築**
Features:Generate a prediction model which is designed for only 1 samples

サンプル1
(Sample 1)

サンプル2
(Sample 2)

サンプル3
(Sample 3)

サンプル・・・
(Sample ・・・)

サンプル・・・
(Sample ・・・)
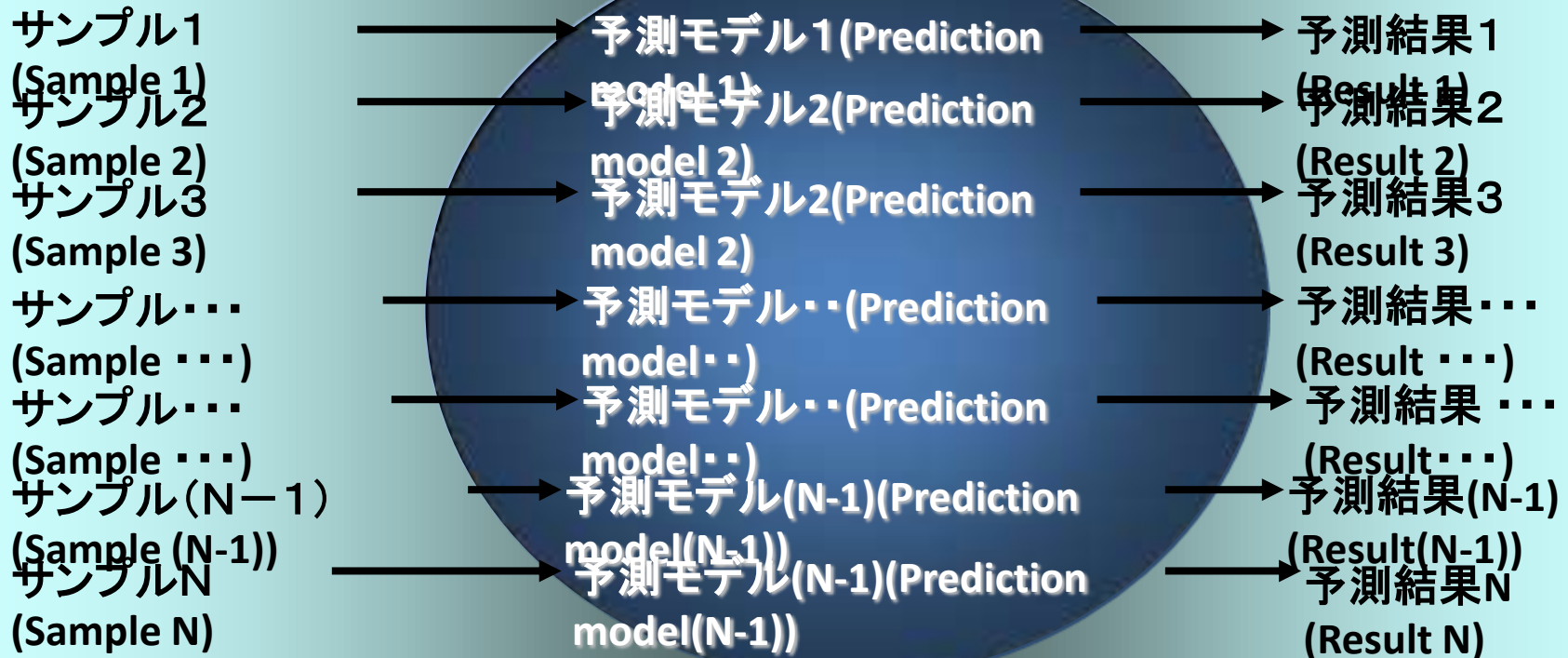
サンプル（N－1）
(Sample (N-1))

サンプルN
(Sample N)

予測モデル1(Prediction model 1)

予測モデル2(Prediction model 2)

予測モデル2(Prediction model 2)

予測モデル・・(Prediction model・・)

予測モデル・・(Prediction model・・)

予測モデル(N-1)(Prediction model(N-1))

予測モデル(N-1)(Prediction model(N-1))

予測結果1
(Result 1)

予測結果2
(Result 2)

予測結果3
(Result 3)

予測結果・・・
(Result ・・・)

予測結果 ・・・
(Result・・・)

予測結果(N-1)
(Result(N-1))

予測結果N
(Result N)

**利点 (Merit)** ： 予測率が大幅に向上する（**High prediction ratio will be achieved**）

**難点（Weakness）**： 計算時間がかかる（**Need large calculation time**）