

# データサイエンスの基礎

---

北陸先端科学技術大学院大学  
Hieu-Chi Dam

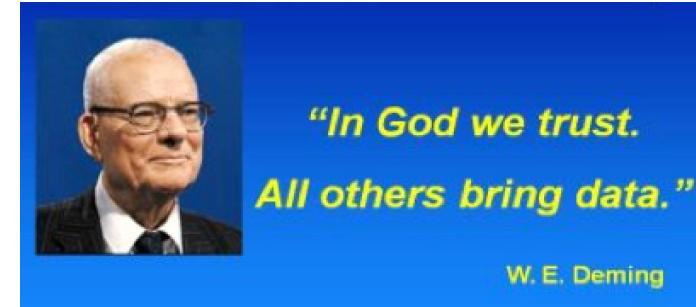
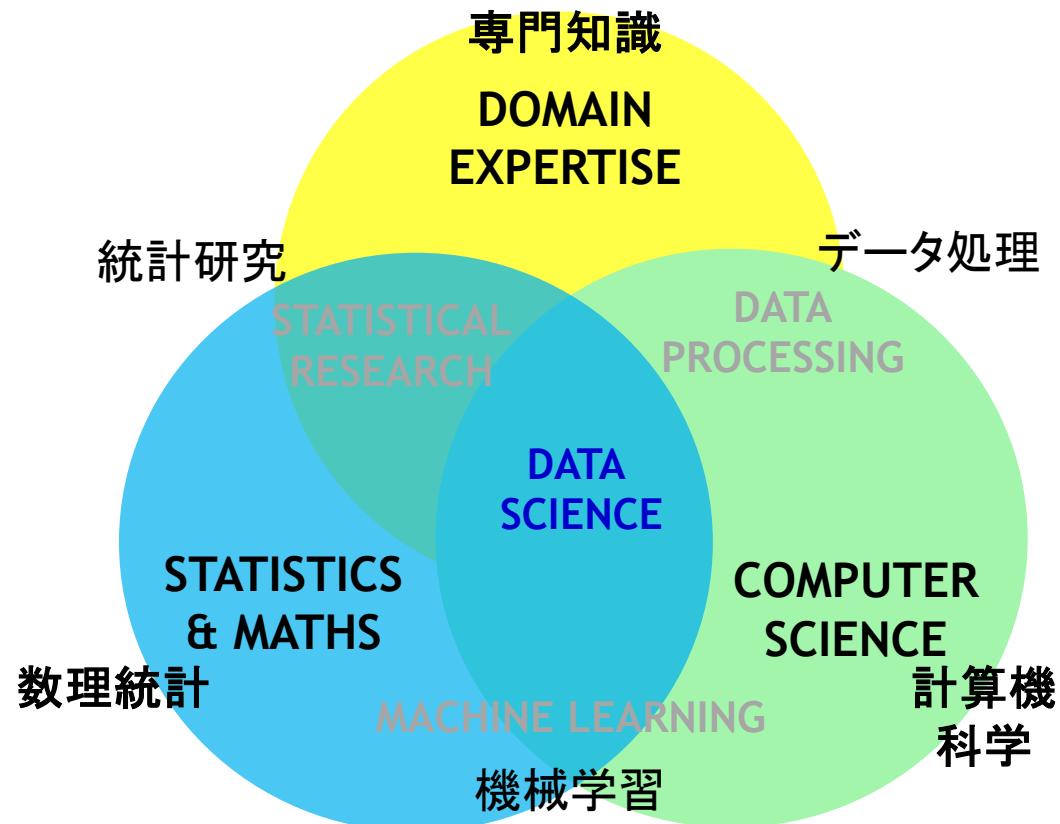
2024年7月8日（第1部 18:00 ~ 19:20）

# 目次

---

1. イントロ
  - i. データサイエンスの紹介
  - ii. 推論方式の紹介
  - iii. 数理統計と機械学習の紹介
2. データサイエンスの材料科学への応用事例紹介

# データサイエンスの定義



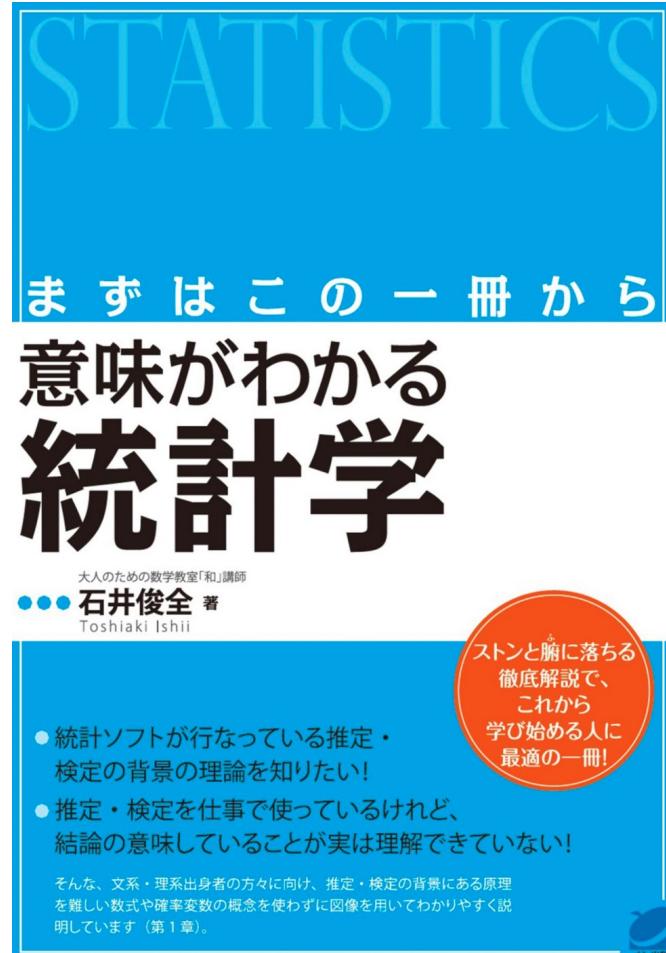
**Data Scientist: The Sexiest  
Job of the 21st Century**  
(Harvard Business Review, October  
2012)

# データの理解と活用

---

1. どのように人々がデータを集め、整理するのかを理解する
2. データの特徴(データ種類、データの特性)を理解する
3. 分析する時のタスクを知る
4. 与えられたタスクに適切なデータ分析手法について知っている
5. 分析ソフトウェアについて知っている
6. 分析から得られた結果について説明できる

# 参考書の紹介



イントロ  
その2：推論方式の紹介

# 演繹法, 帰納法, およびアブダクション

科学のパラダイム:

実験科学 - 第1パラダイム

- 実験, 観測などから結論を導く

理論科学 - 第2パラダイム

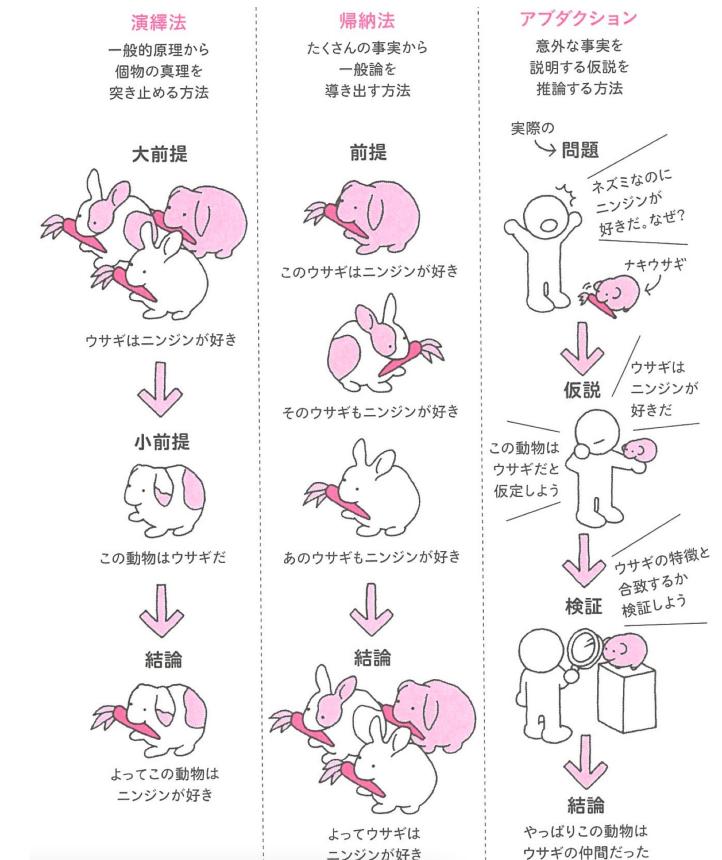
- ニュートン方程式など

計算科学 - 第3パラダイム

- 複雑現象の計算機シミュレーション, 計算機実験など

データ科学 - 第4パラダイム

- 理論, 実験, 計算機実験, データマイニング



# 演繹法, 帰納法, およびアブダクション

科学のパラダイム:

実験科学 - 第1パラダイム

- 実験, 観測などから結論を導く

理論科学 - 第2パラダイム

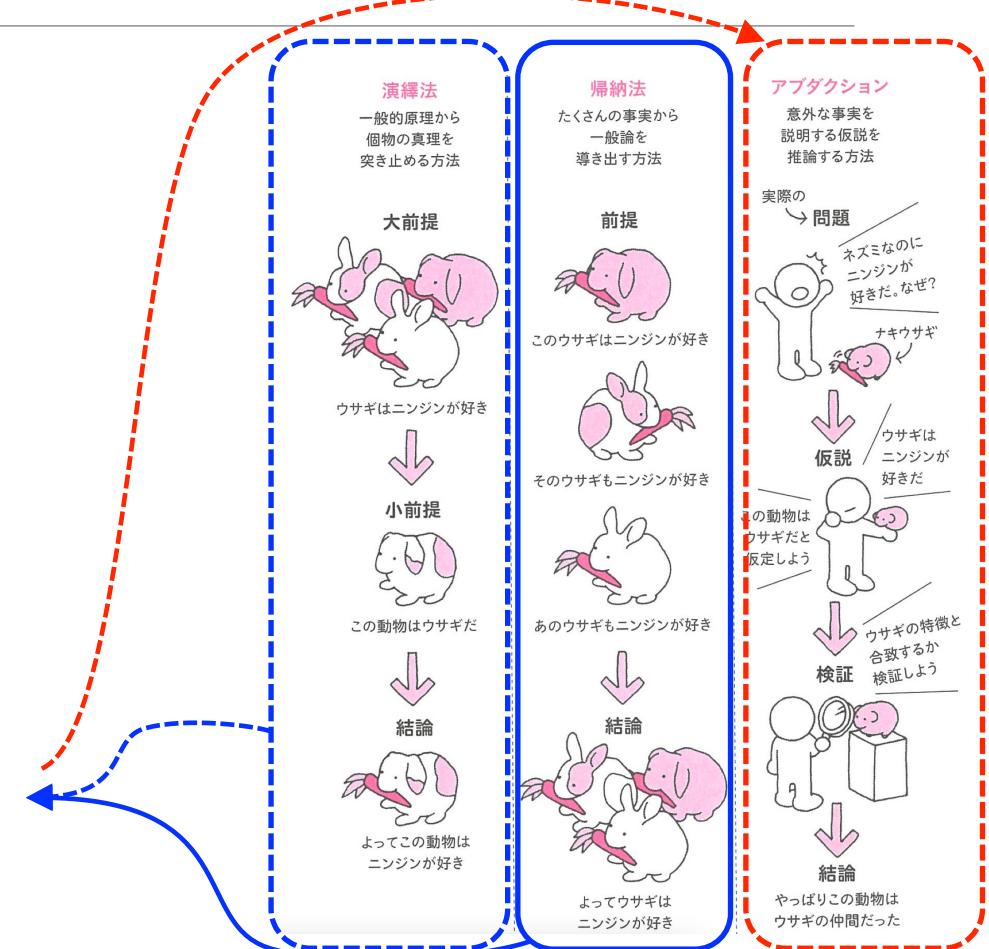
- ニュートン方程式など

計算科学 - 第3パラダイム

- 複雑現象の計算機シミュレーション, 計算機実験など

データ科学 - 第4パラダイム

- 理論, 実験, 計算機実験, データマイニング



哲学用語図鑑・田中正人(著)・プレジデント社

# 演繹法, 帰納法, およびアブダクション

## 演繹法

一般的原理から  
個別の真理を  
突き止める方法

## 大前提



ウサギはニンジンが好き

## 小前提



この動物はウサギだ

## 結論



よってこの動物は  
ニンジンが好き

大前提が正しければ結論が正しい(保証される)

妥当性の検証  
が必要

不確かさの定  
量評価が必要

諸前提が正しいであっても  
必ずしも結論が正しいと限  
らない(保証されない)

## 帰納法

たくさんの事実から  
一般論を  
導き出す方法

## 前提



このウサギはニンジンが好き



そのウサギもニンジンが好き



この動物は  
ウサギだと  
仮定しよう



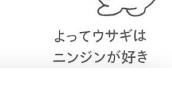
あのウサギもニンジンが好き

## 結論



ウサギの特徴と  
合致するか  
検証しよう

## 検証



やっぱりこの動物は  
ウサギの仲間だった

## アブダクション

意外な事実を  
説明する仮説を  
推論する方法

## 実際の問題



## 仮説



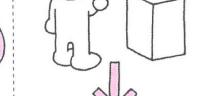
## 結論



## 検証



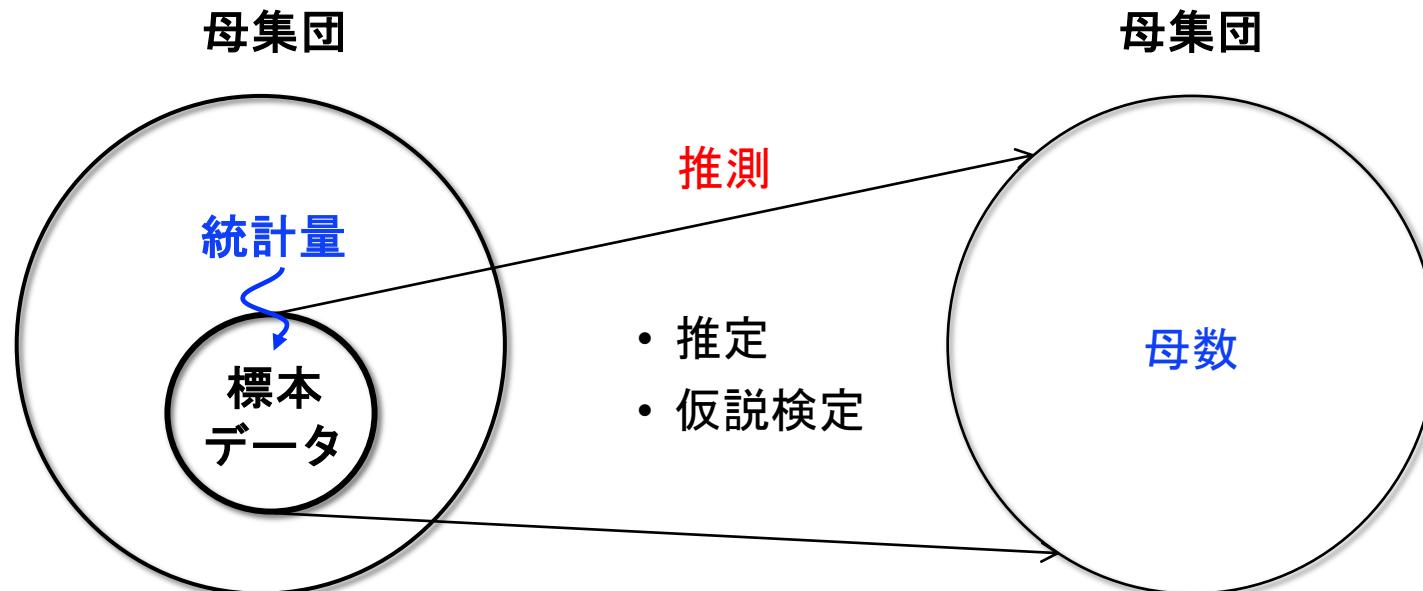
## 結論



# イントロ

## その3：数理統計と機械学習の紹介

# 統計学の基礎： 統計的推測



統計学的推測とは、標本データを解析して母数に関する結論を導くこと。

# 統計学の基礎： 統計的推測

---

- ✓ **母数**: 母集団全体の特性を表す値(例: 母平均、母分散など)。母集団全体のデータに基づいていて、理想的には真の値を示す。
- ✓ **統計量**: 母集団の一部である標本のデータに基づいて計算される値(例: 標本平均、標本分散)。統計量は、**母数を推定**するために使用される。
  - **推定**: すべてのデータをできないため、標本を使用して母数を推定する。
  - **不確実性**: 標本に基づく統計量は、標本の選び方や偶然の影響により、異なる値を有り得るため、統計量には不確実性が伴う。母数は母集団全体に基づいて、(理想的な)確定した真の値。

# 統計学の基礎： 数学的推論

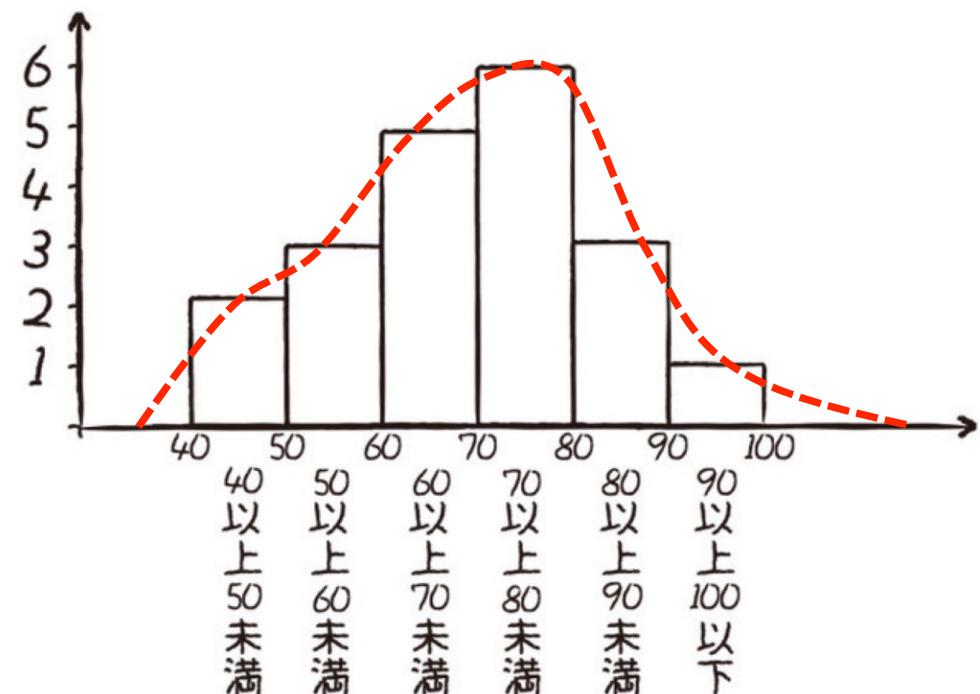
生データ

2人 43、47  
3人 52、52、54  
5人 61、67、67、68、69  
6人 70、71、71、73、76、78  
3人 82、84、84  
1人 91

整理された  
データ

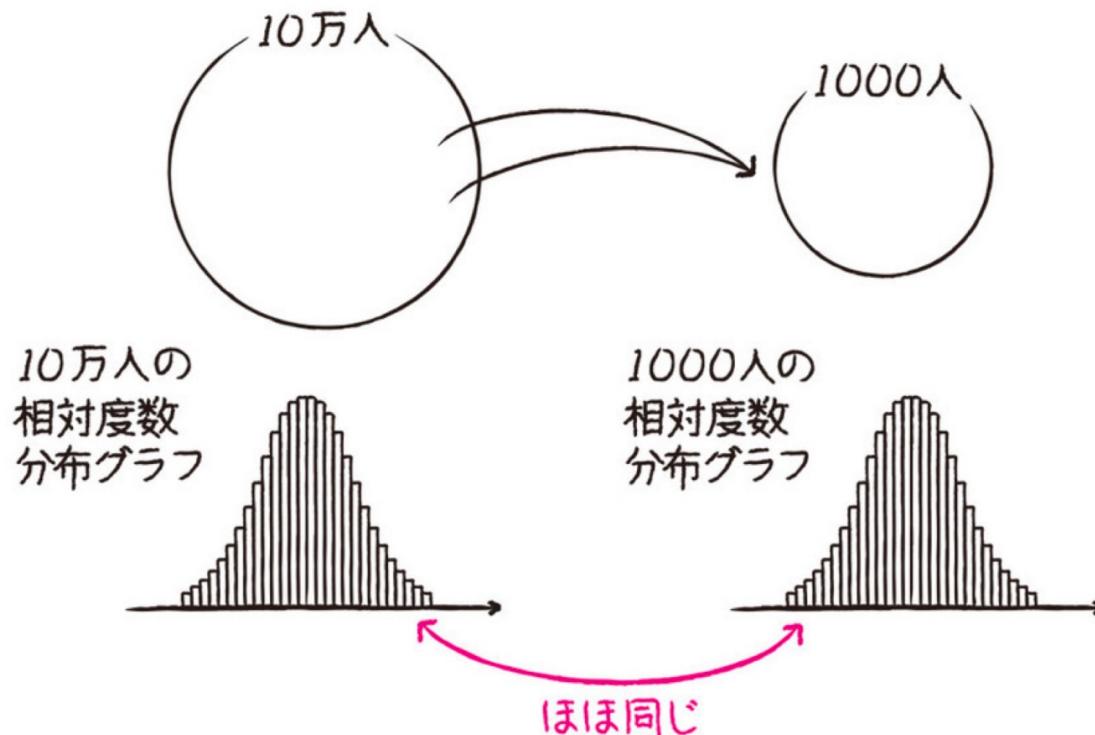
階級	階級値	度数
40 以上 50 未満	45	2
50 以上 60 未満	55	3
60 以上 70 未満	65	5
70 以上 80 未満	75	6
80 以上 90 未満	85	3
90 以上 100 以下	95	1

可視化されたデータ

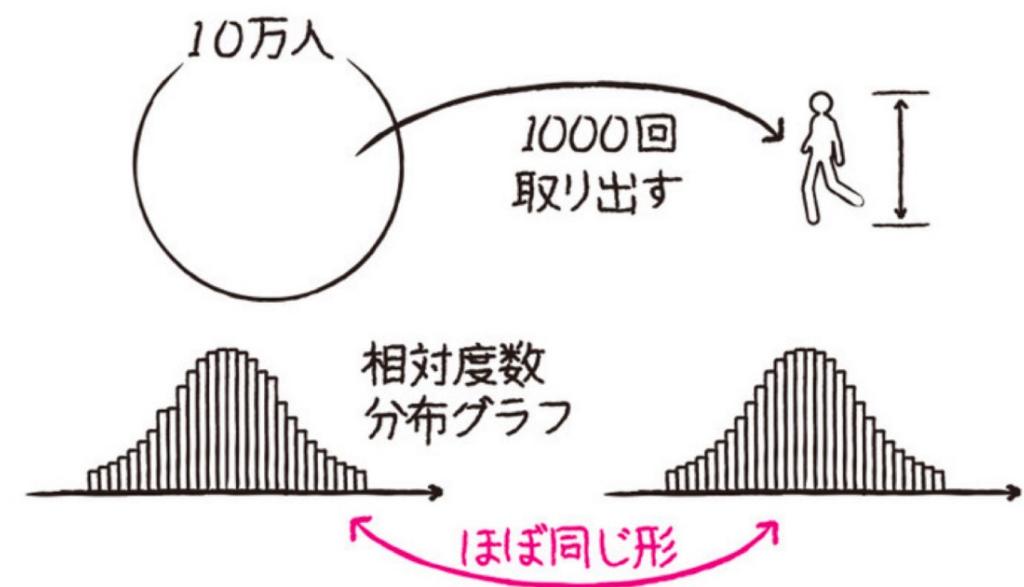


# 統計学の基礎：データサンプリング

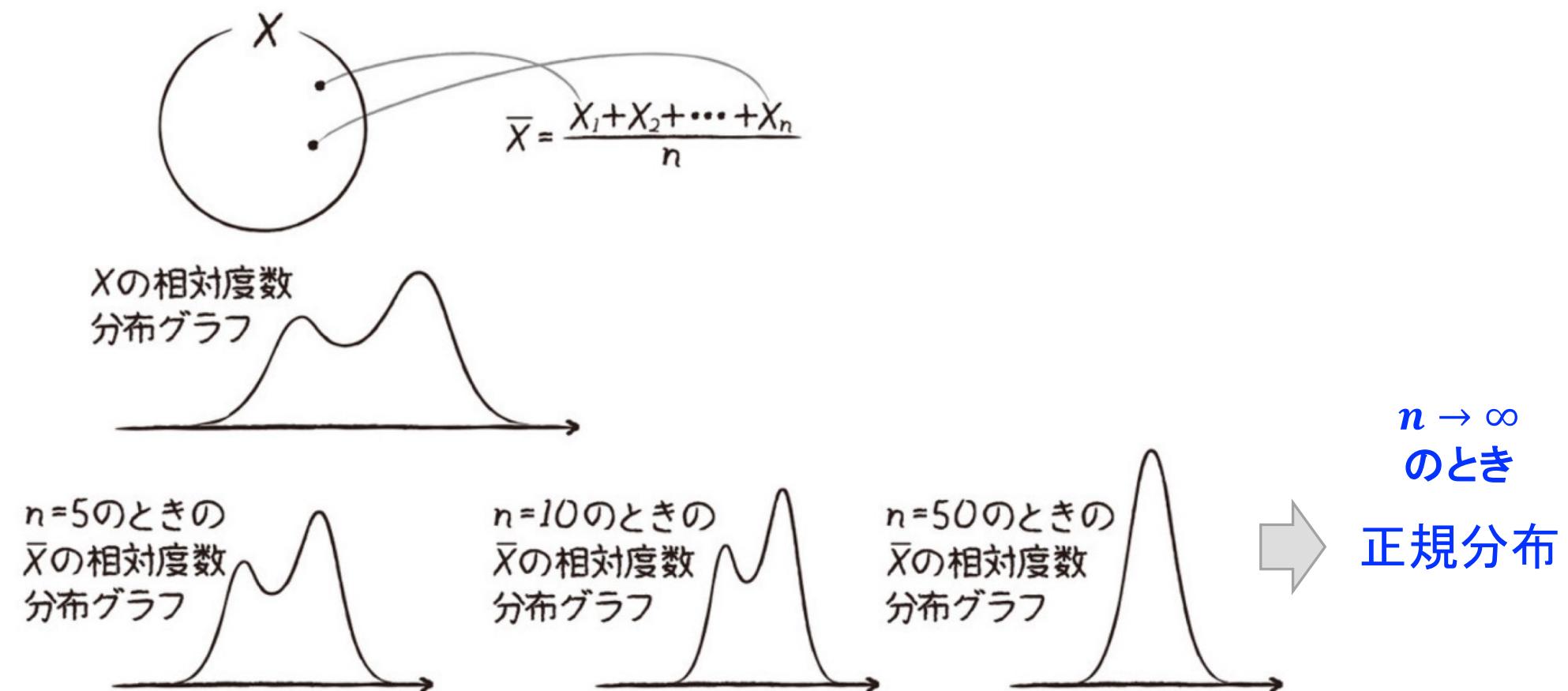
非復元抽出



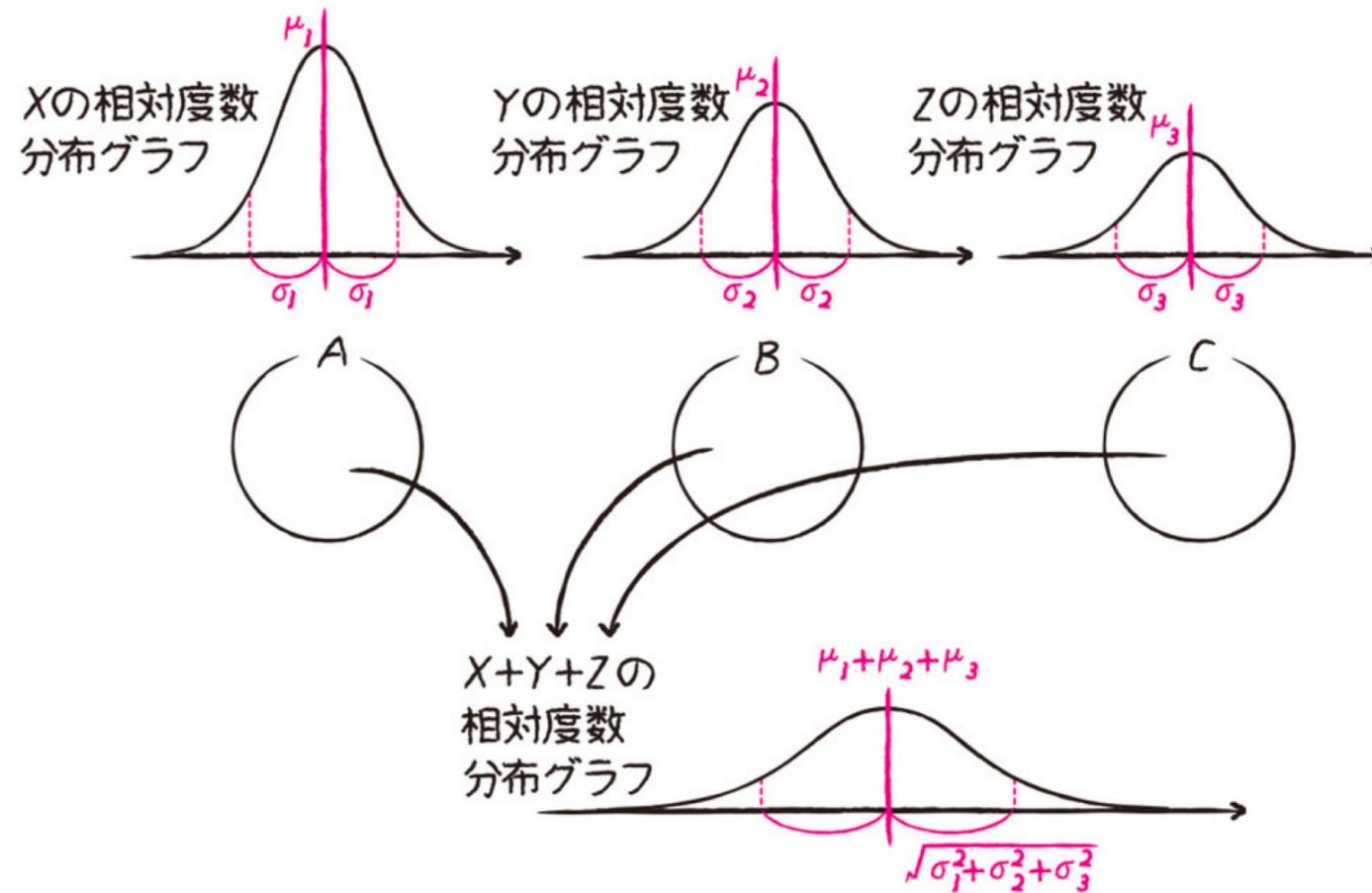
復元抽出



# 統計学の基礎： 中心極限定理



# 統計学の基礎：正規分布



# 統計学の基礎： 推定

---

不偏推定量による推定

最尤法による推定

# 統計学の基礎： 推定

---

与えられたデータは**現実に出現**している



考えられる数理モデルの集合を検討する



**最尤法による推定**

与えられた**データを再現する数理モデル**は  
**比較的高い確率**が割り振られている

# 統計学の基礎：コイン投げ

---



なぜそのように計算できるか？

- 本当なのか？ どれくらい信頼できるか？
- この計算はデータ解析なのか？
- 一種のデータ解析だとしたら、このデータ解析が何を基にしているか？
- 背景にはどのような仮説があるか？

100 回: 70 表 & 30 裏

$$p(\text{表}) = \frac{70}{100} = 0.7$$

# 統計学の基礎：コイン投げ

---



100回: 70表 & 30裏

仮定

$$p(\text{表}) = \theta$$

尤度

$$\mathcal{L}(\theta) = {}_{100}C_{70}\theta^{70}(1 - \theta)^{30}$$

考え方：

与えられたデータを再現する数理モデルは  
比較的高い確率が割り振られている

# 統計学の基礎：コイン投げ

---



100回: 70表 & 30裏

仮定

$$p(\text{表}) = \theta$$

尤度

$$\mathcal{L}(\theta) = {}_{100}C_{70}\theta^{70}(1 - \theta)^{30}$$

$$\theta \triangleq \arg \max_{\theta \in [0,1]} \log p(\mathcal{D}|\theta)$$

$$\triangleq \arg \max_{\theta \in [0,1]} \log \mathcal{L}(\theta)$$

$$\triangleq \arg \max_{\theta \in [0,1]} \log {}_{100}C_{70}\theta^{70}(1 - \theta)^{30}$$

# 統計学の基礎：コイン投げ

---



100回: 70表 & 30裏

仮定

$$p(\text{表}) = \theta$$

尤度

$$\mathcal{L}(\theta) = {}_{100}C_{70}\theta^{70}(1-\theta)^{30}$$

$$\frac{d}{d\theta} \log \mathcal{L}(\theta) = \frac{70}{\theta} - \frac{30}{1-\theta} = 0$$

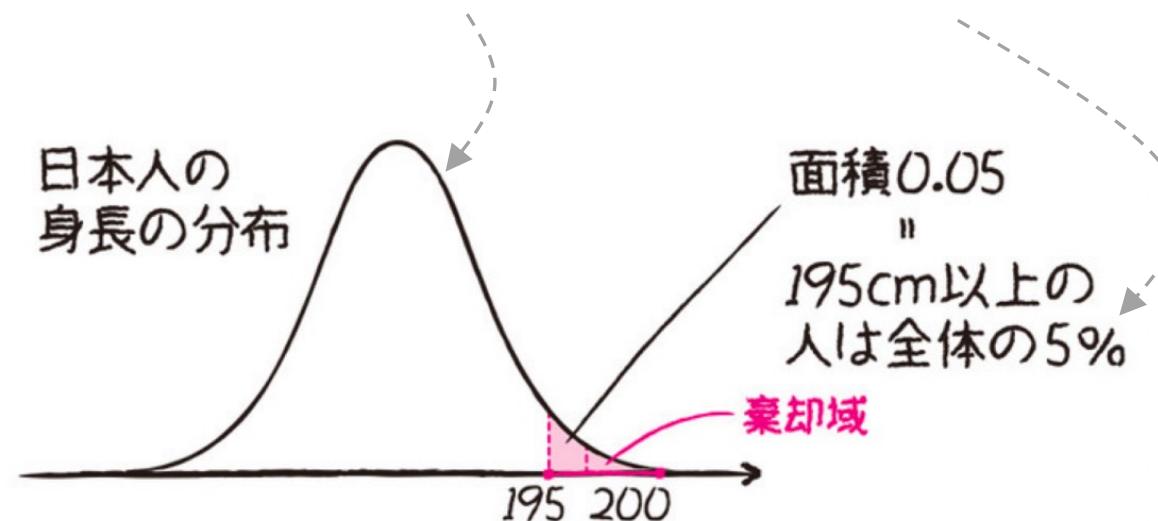
$$\theta = \frac{70}{100} = 0.7$$

# 統計学の基礎：仮説検定

状況：街（日本）で向こうを歩いている身長2mの人を見ました。

疑問：歩いてきた人は外国人ではないか？

理由：日本人だとすれば、大きすぎる。常識から考えて、あり得ないと思う。



# 統計学の基礎：仮説検定

「むこうを歩いている人が  $\leftarrow H_0$  を仮定 帰無仮説

日本人であるとすれば、

大きすぎるなあ。

ありえないよ。

日本人じゃないな。

外国人じゃないの」

$\leftarrow$  起こる確率が 5% 以下の  
できごとが起きた

$\leftarrow H_0$  棄却

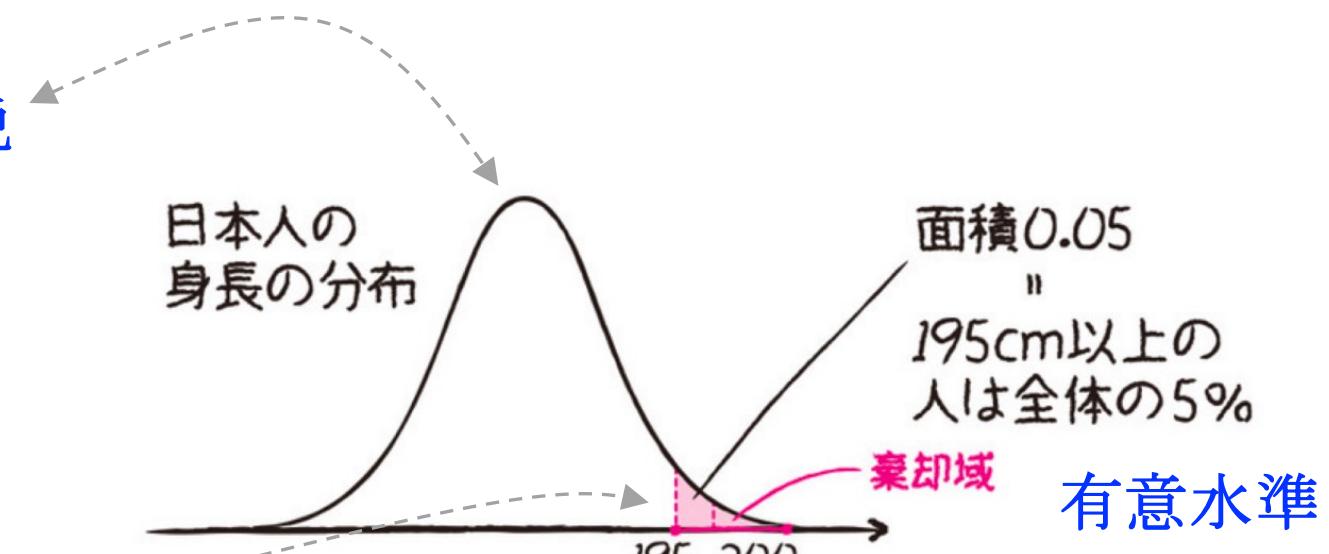
$\leftarrow H_1$  採択 対立仮説

日本人の  
身長の分布

面積 0.05  
" 195cm 以上の  
人は全体の 5%

有意水準

判断結果	帰無仮説 $H_0$ を採択	帰無仮説 $H_0$ を棄却
事実		
帰無仮説 $H_0$ が本当	正しい判断	第 1 種の誤り
帰無仮説 $H_0$ が嘘	第 2 種の誤り	正しい判断



# 統計学の基礎：仮説検定と背理法

状況：街（日本）で向こうを歩いている身長2mの人を見ました。

疑問：歩いてきた人は外国人ではないか？

理由：日本人だとすれば、大きすぎる。常識から考えて、あり得ないと思う。

背理法：

- ① 日本人だと仮定する。
- ② 日本人なら身長が“それほど”背が高くない。
- ③ 見えた人の身長が2m
- ④ ③が②と矛盾しているので①が正しくない  
(矛盾しない場合は①が正しくないことはまだ言えない)

対立仮説  $H_1$  ← 証明したいこと

帰無仮説  $H_0$  を仮定する ← 反対のことを仮定する

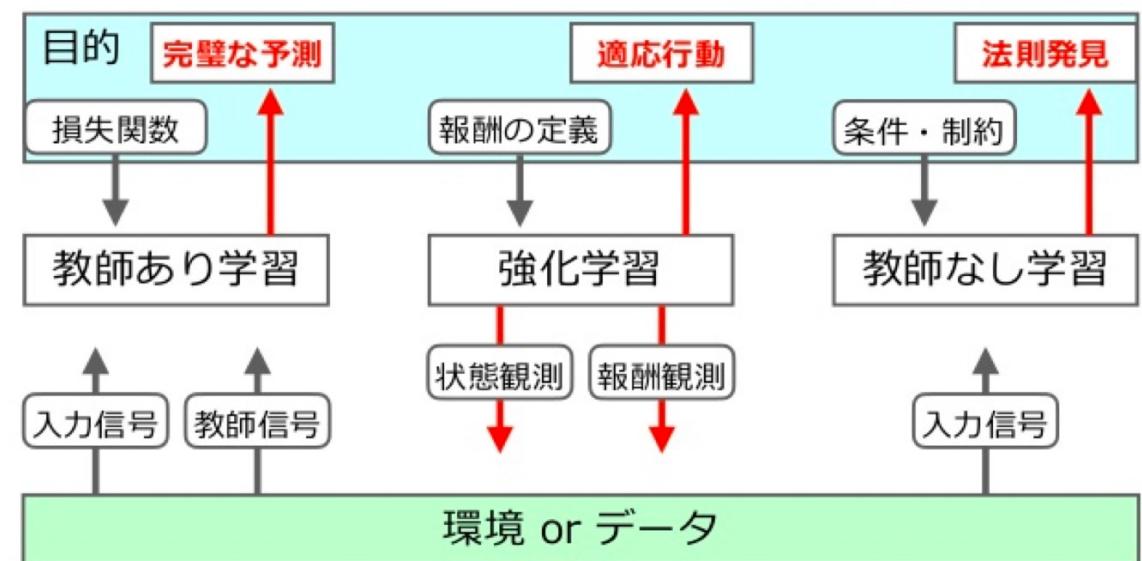
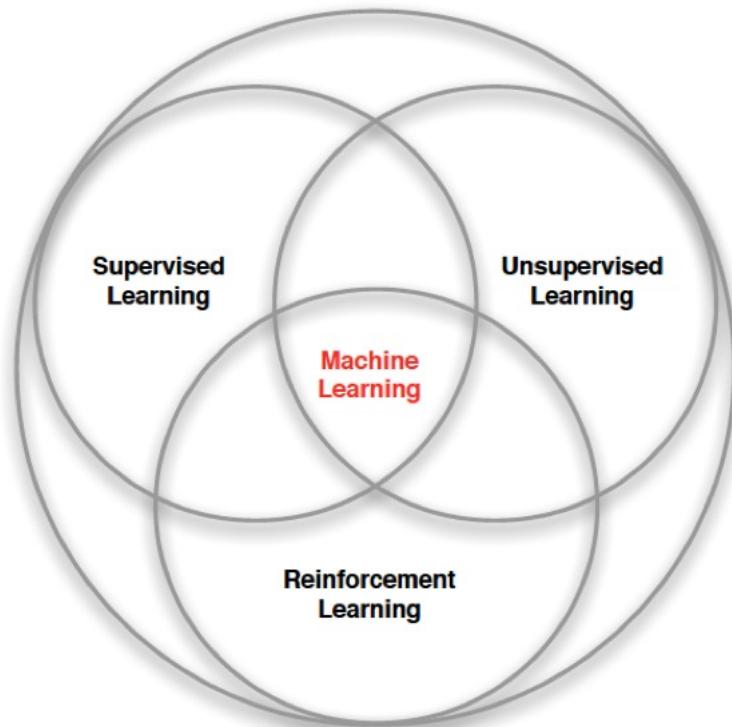
このもとで確率を計算

ありえそうもないことが起きている

対立仮説  $H_1$  を採択 ← 証明完了

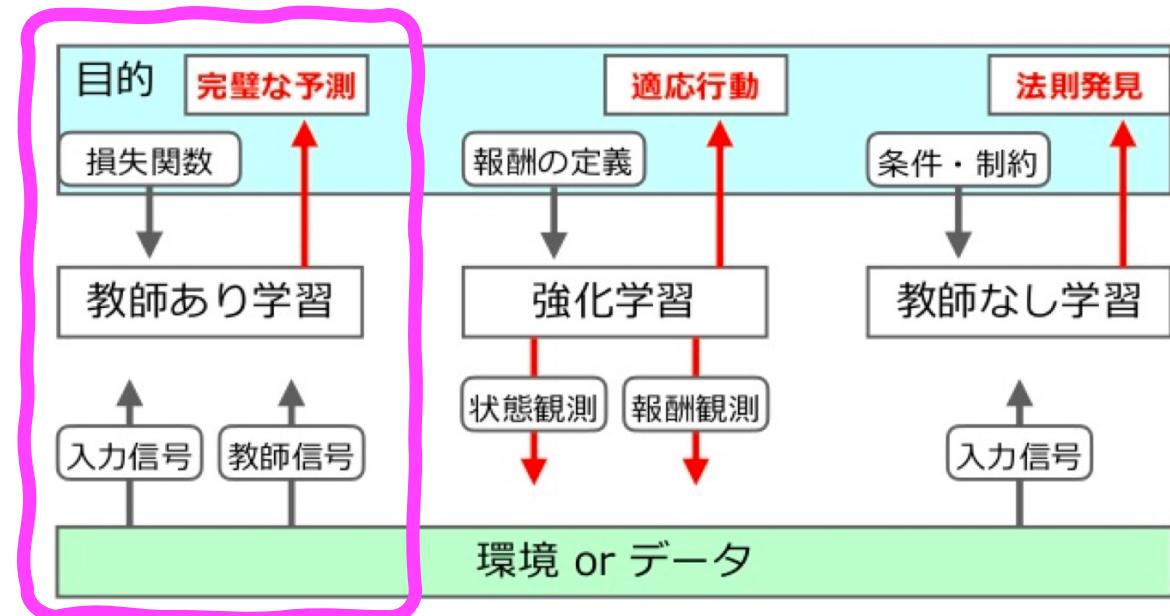
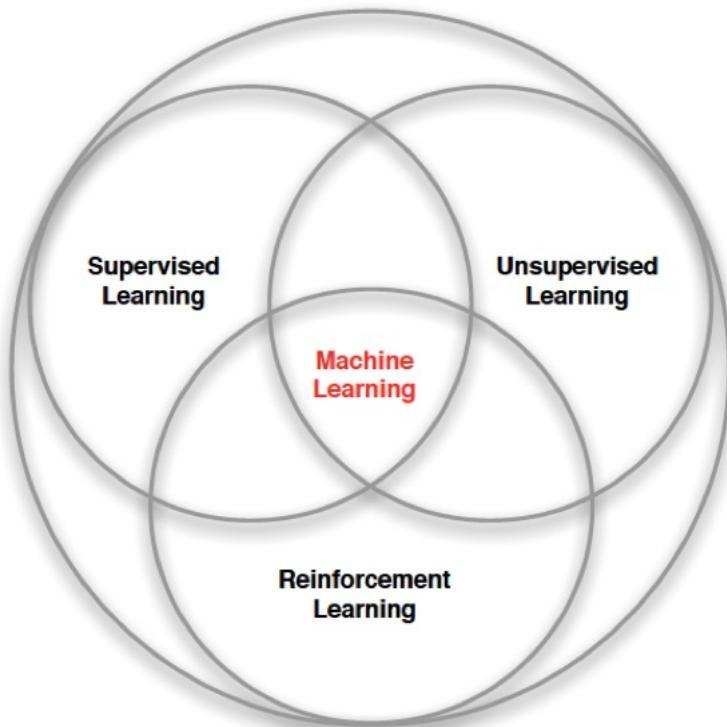
} 矛盾を導く

# 機械学習



データ所与である前提と異なり、強化学習は  
環境を探索して主体的にデータを獲得しつつ行動方策を最適化

# 機械学習

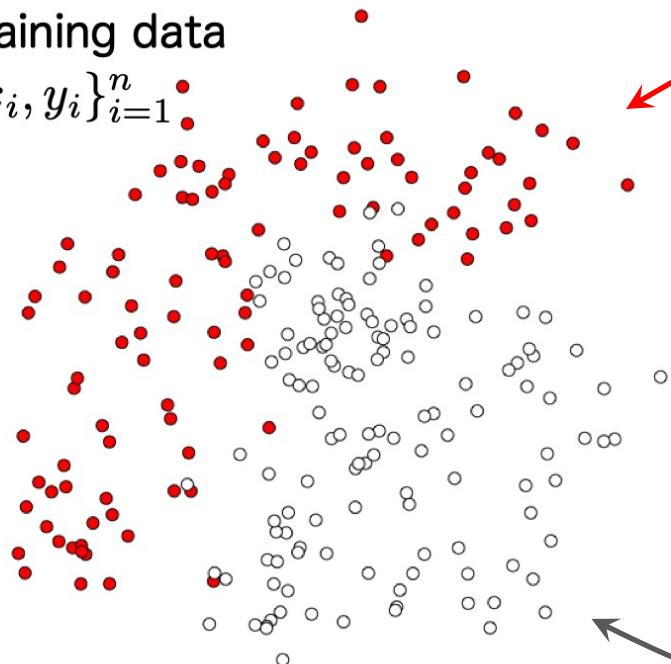


データ所与である前提と異なり、強化学習は  
環境を探索して主体的にデータを獲得しつつ行動方策を最適化

# 教師あり学習

Training data

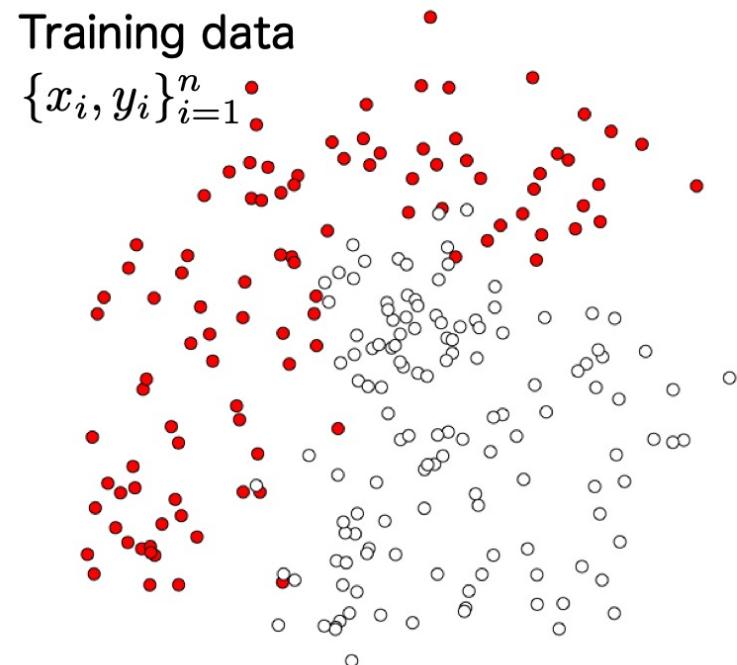
$$\{x_i, y_i\}_{i=1}^n$$



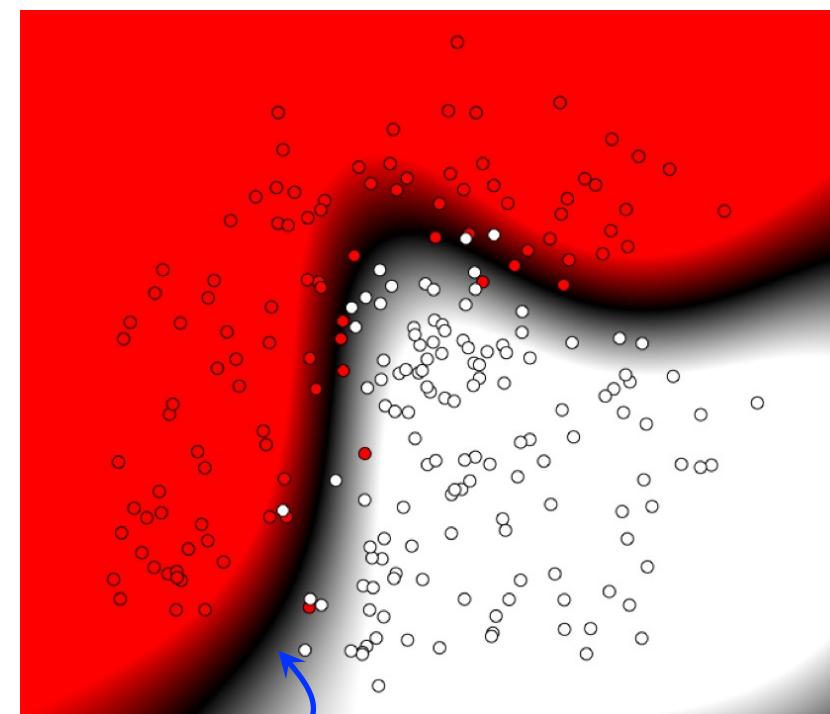
データ変換  
(属性, 特徴量など)



# 教師あり学習

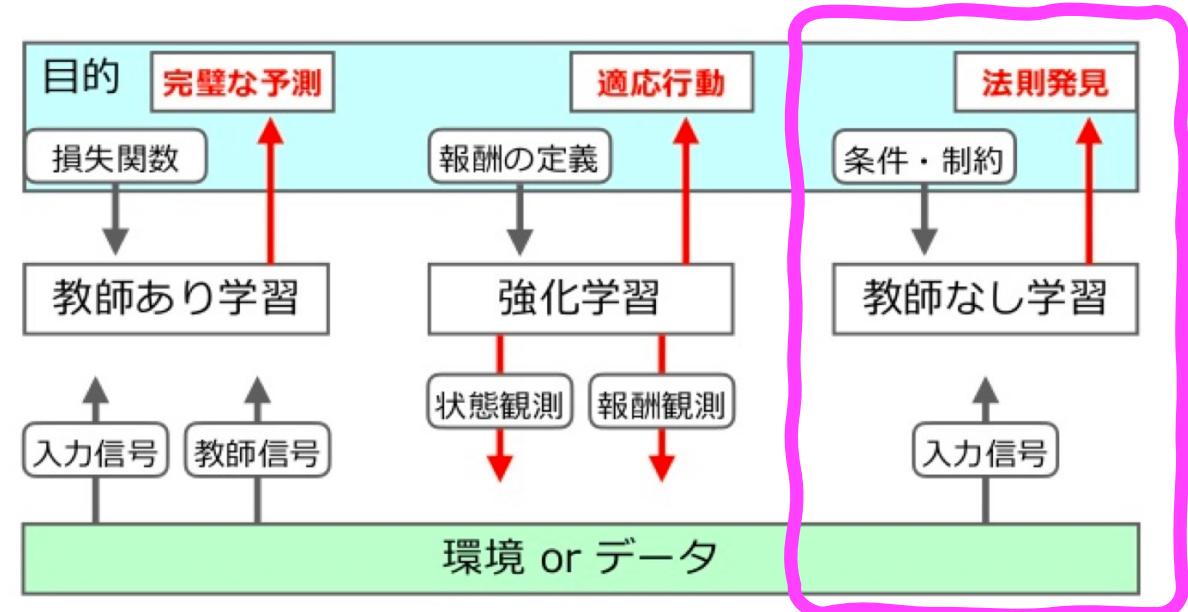
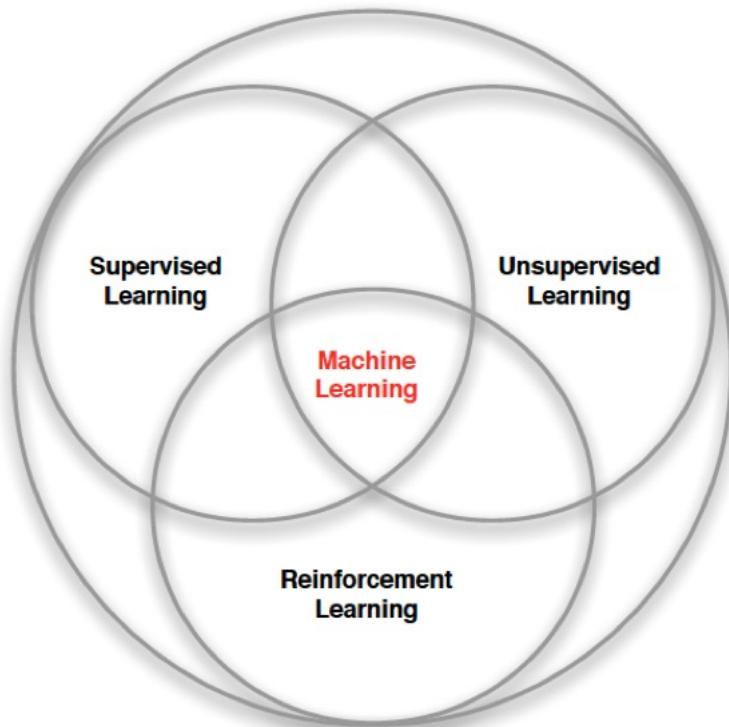


分類機  
→



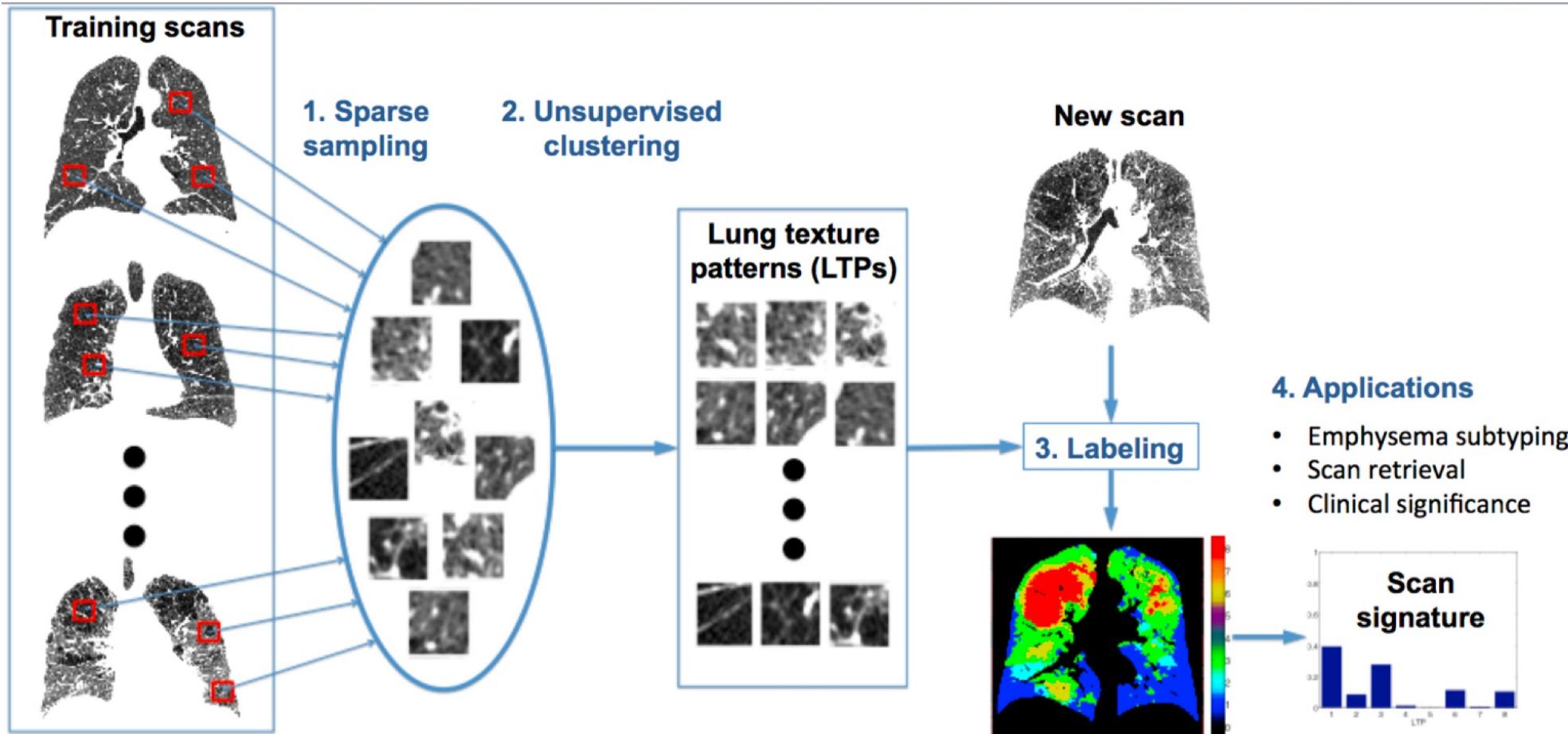
猫  
教師信号  
犬

# 機械学習

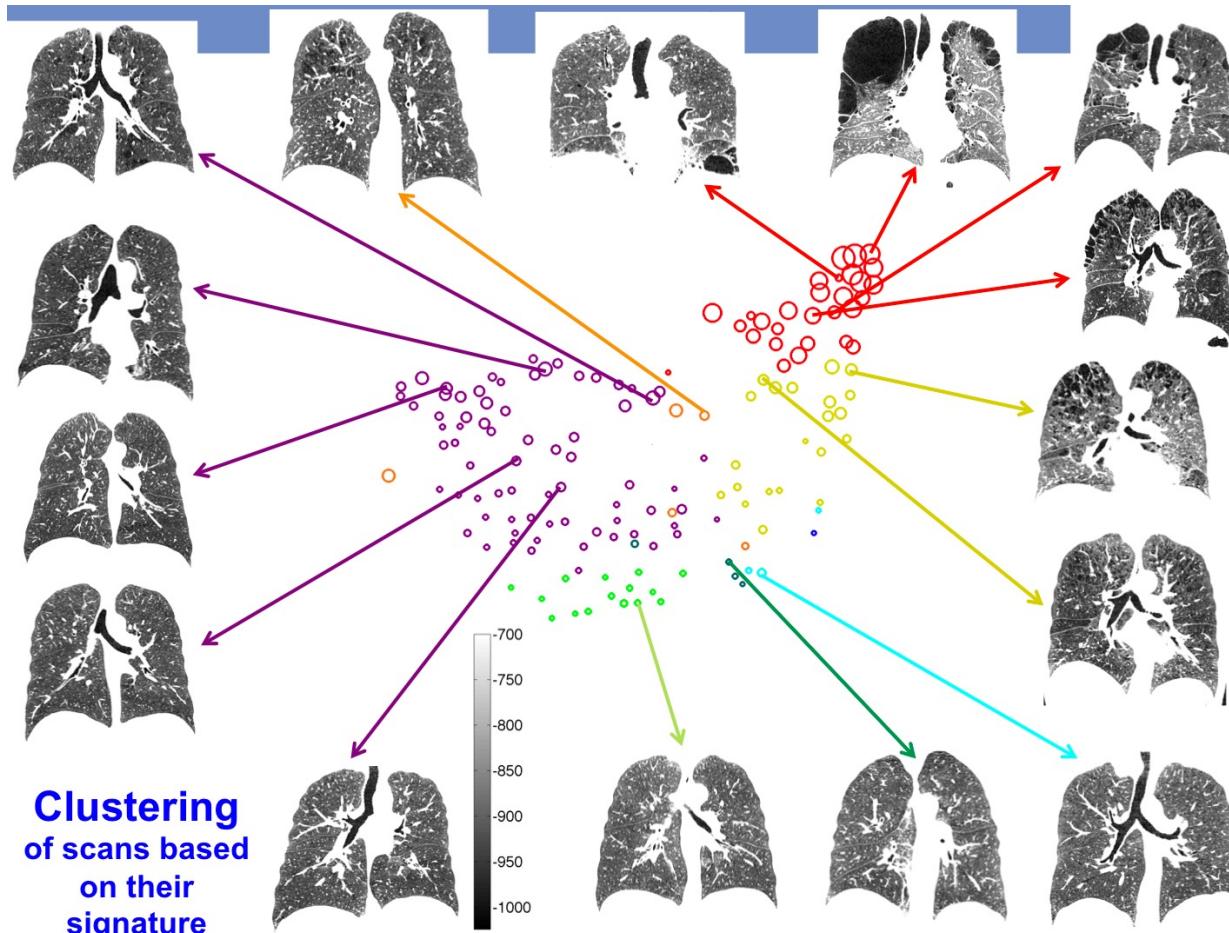


データ所与である前提と異なり、強化学習は  
環境を探索して主体的にデータを獲得しつつ行動方策を最適化

# 教師なし学習



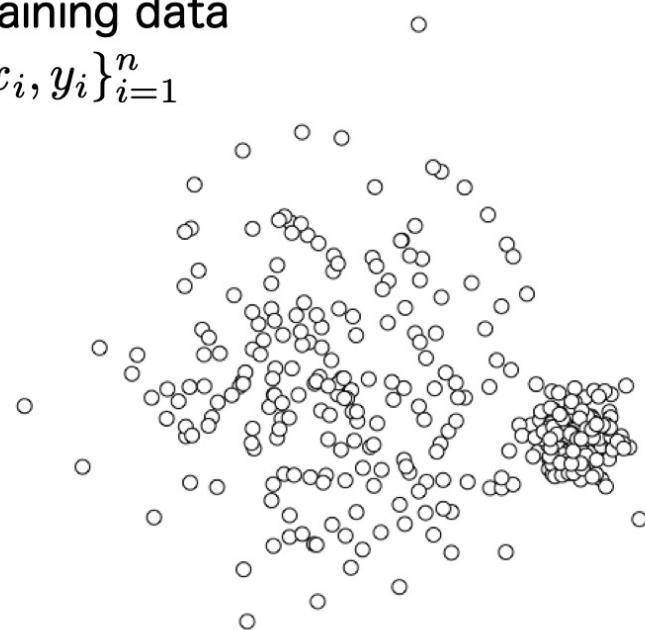
# 教師なし学習



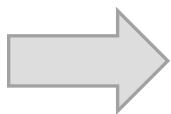
# 教師なし学習

Training data

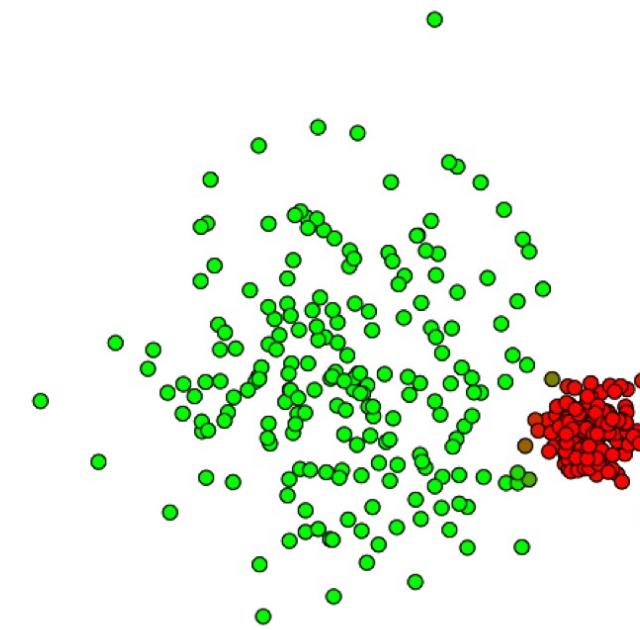
$$\{x_i, y_i\}_{i=1}^n$$



クラスターリング



入力信号だけ、教師信号なし

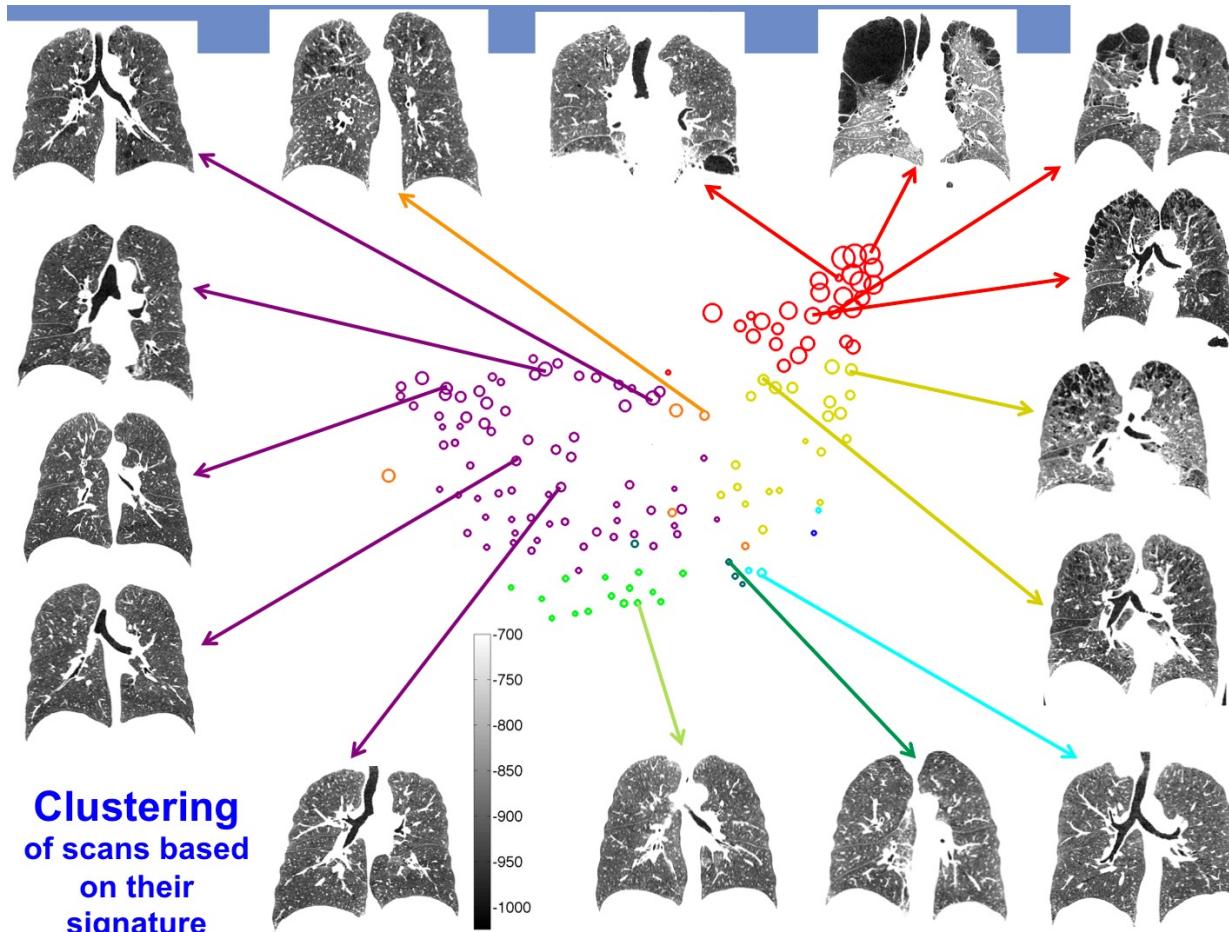


グループ1

グループ構造  
の発掘

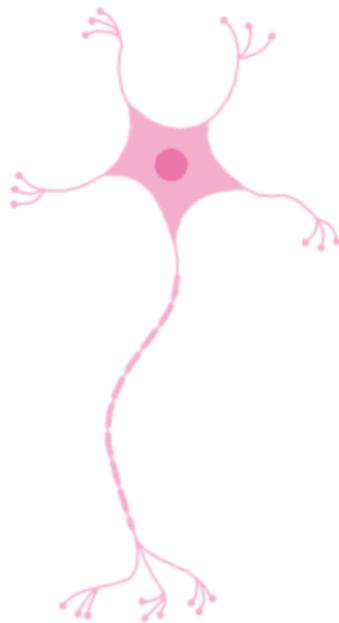
グループ2

# 教師なし学習

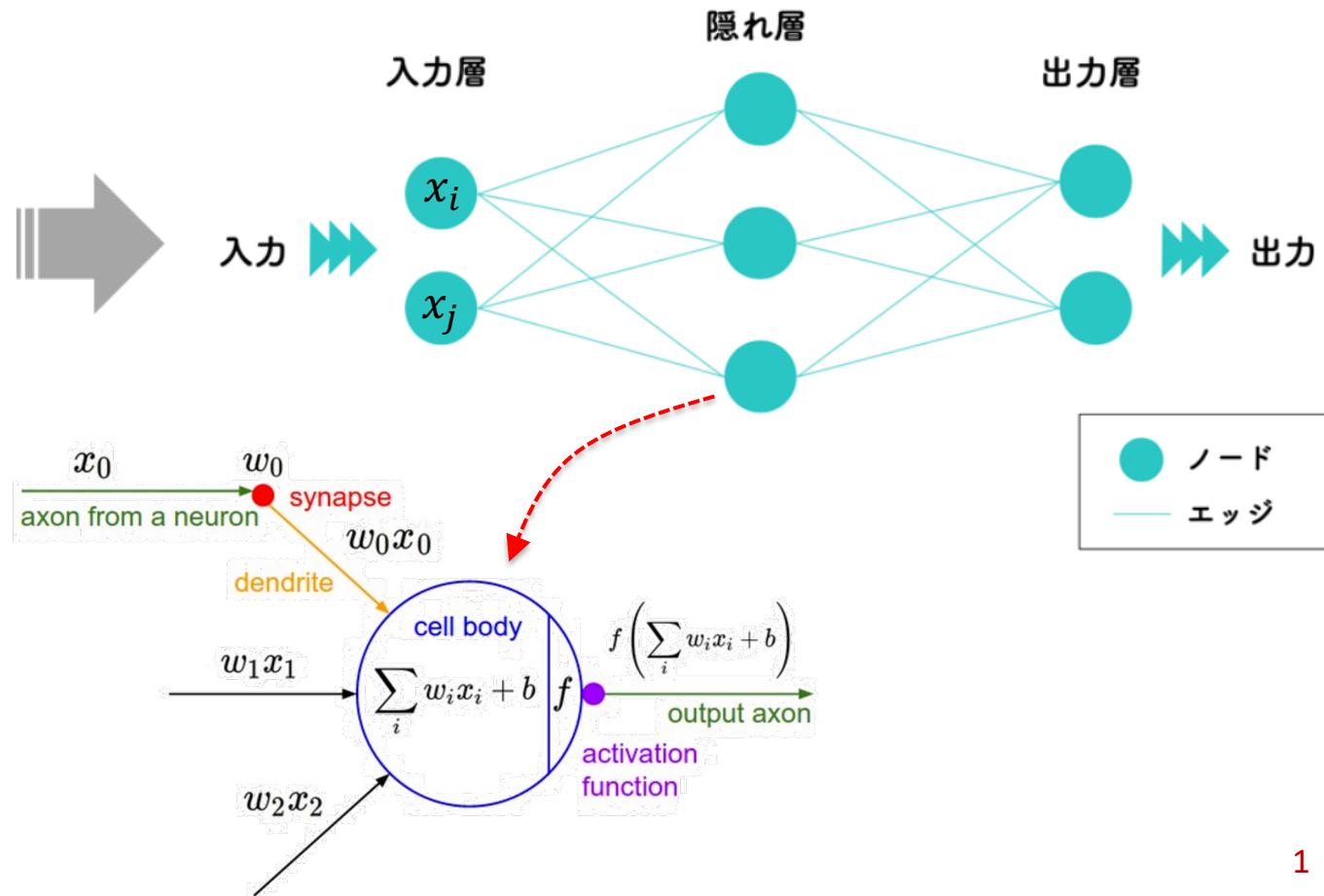


# ニューラルネットワーク構造

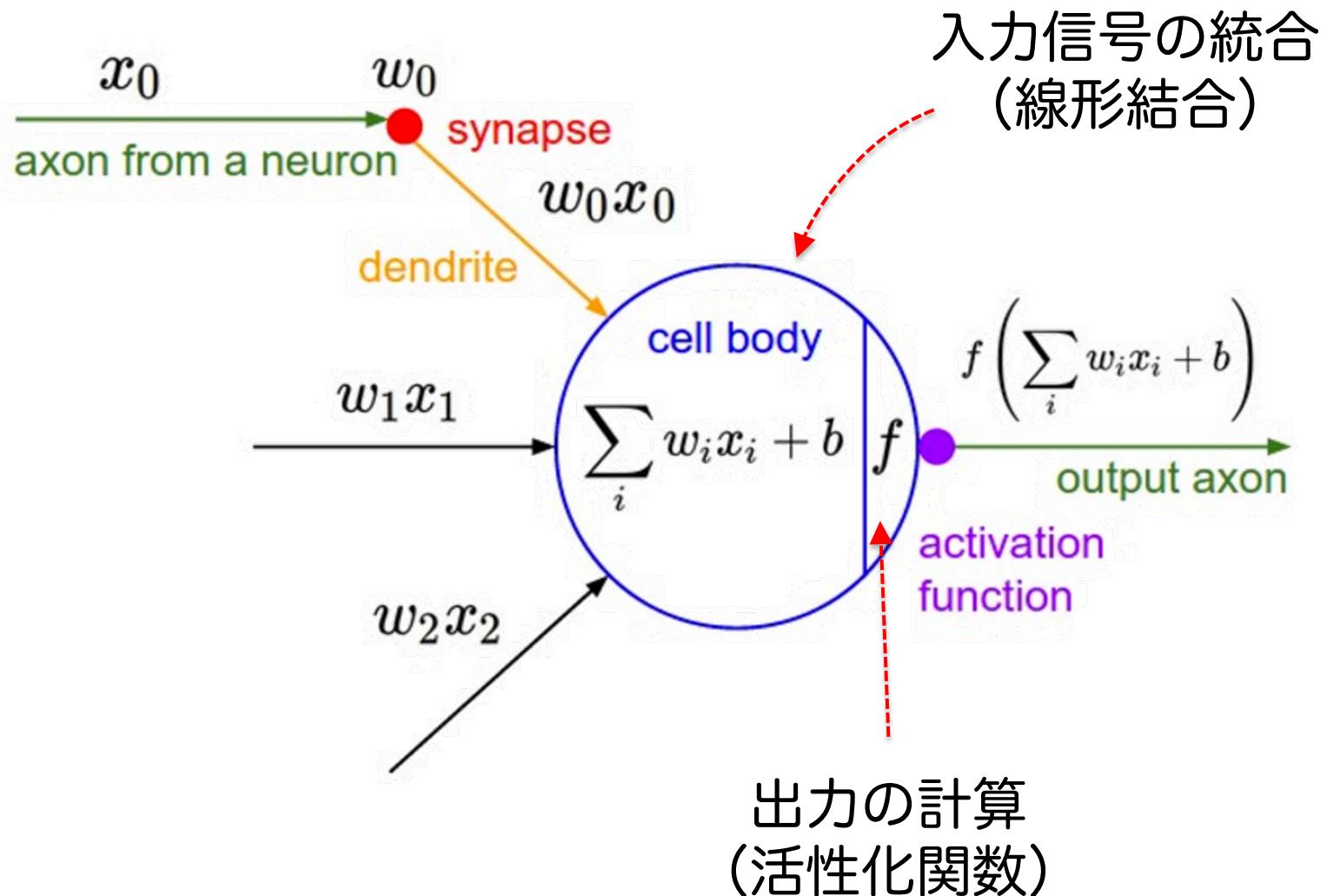
神経細胞(ニューロン)



ニューラルネットワーク

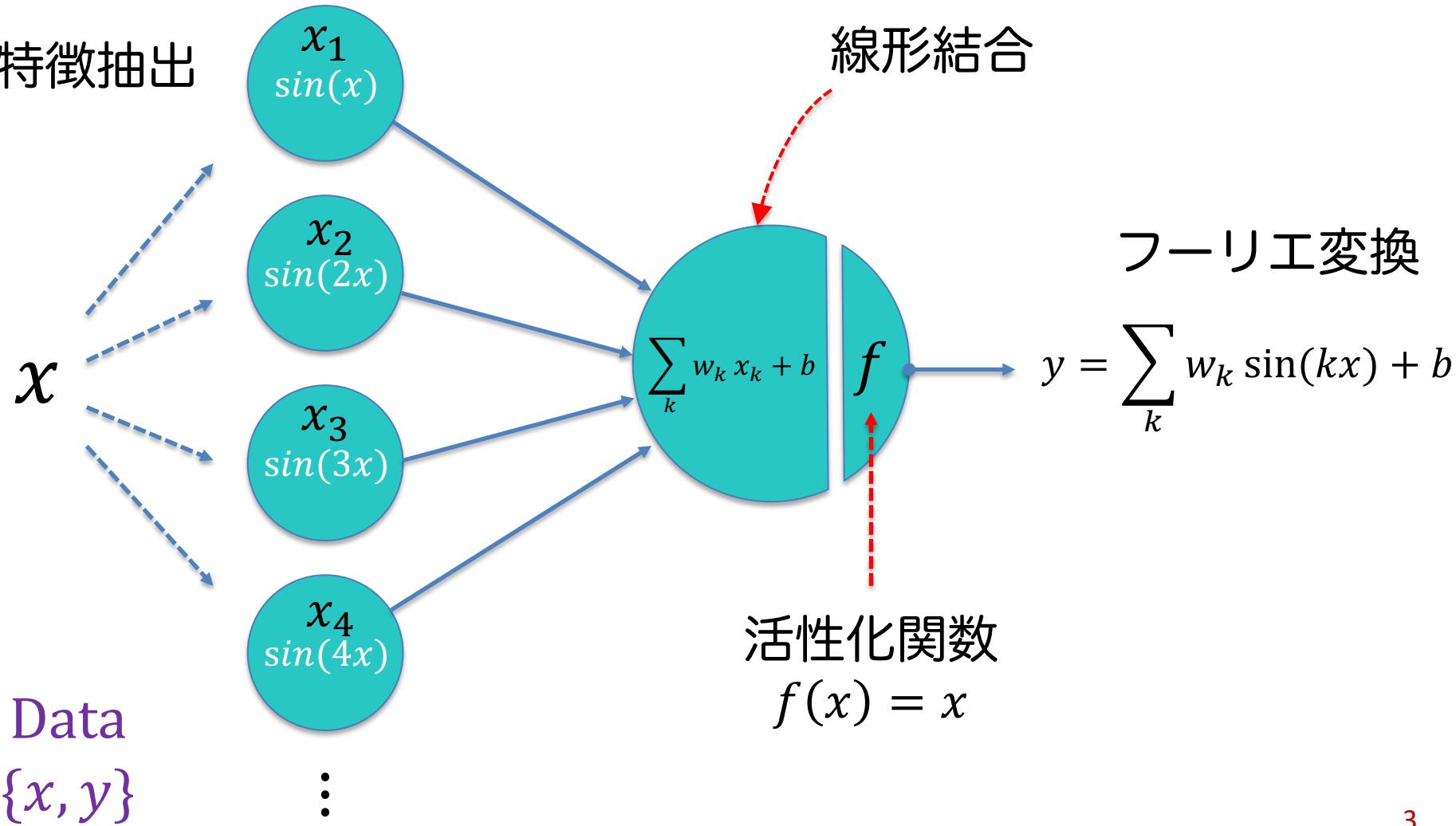


# ニューラルノード



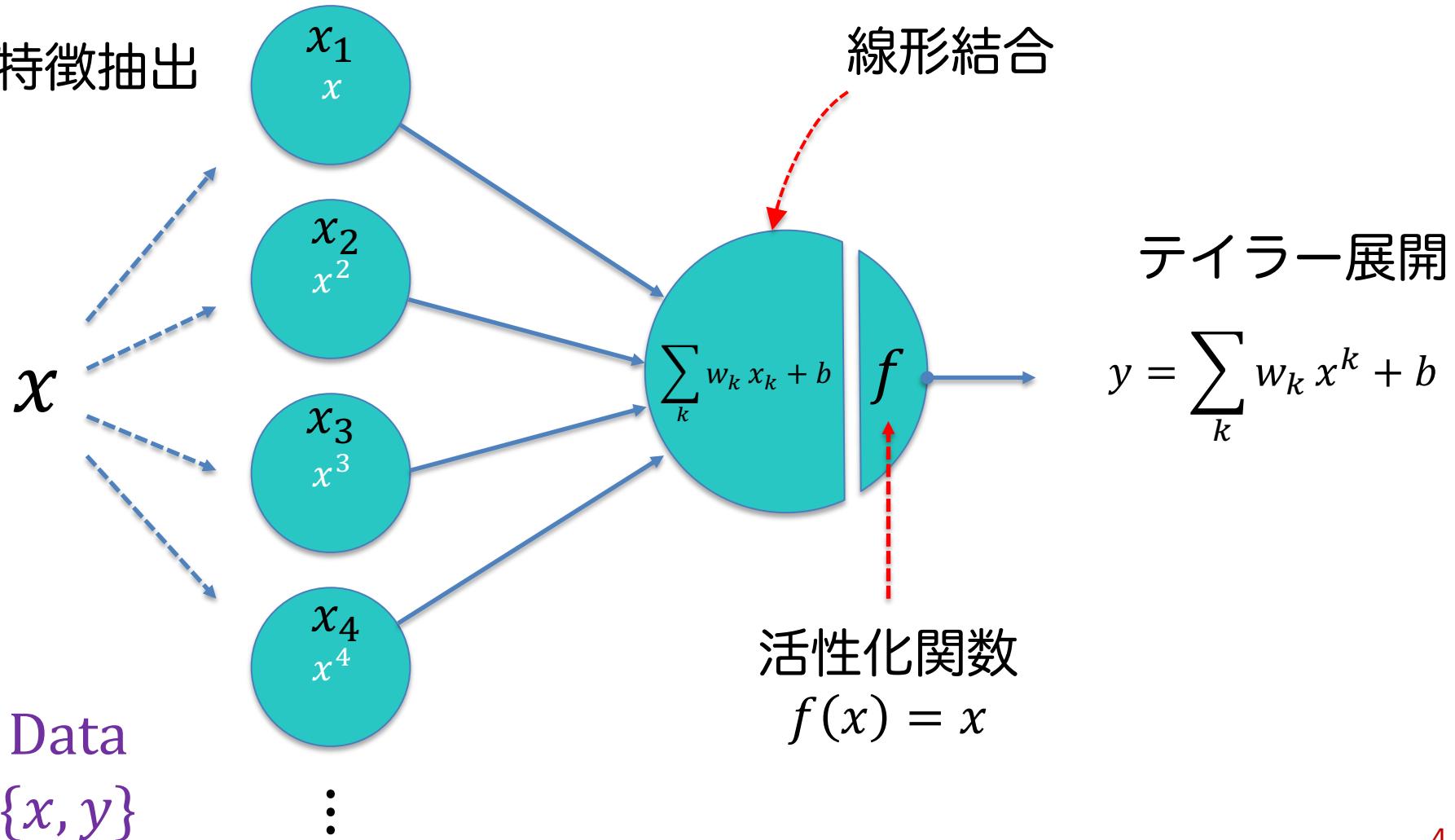
# ニューラルネットの事例 ①

特徴抽出

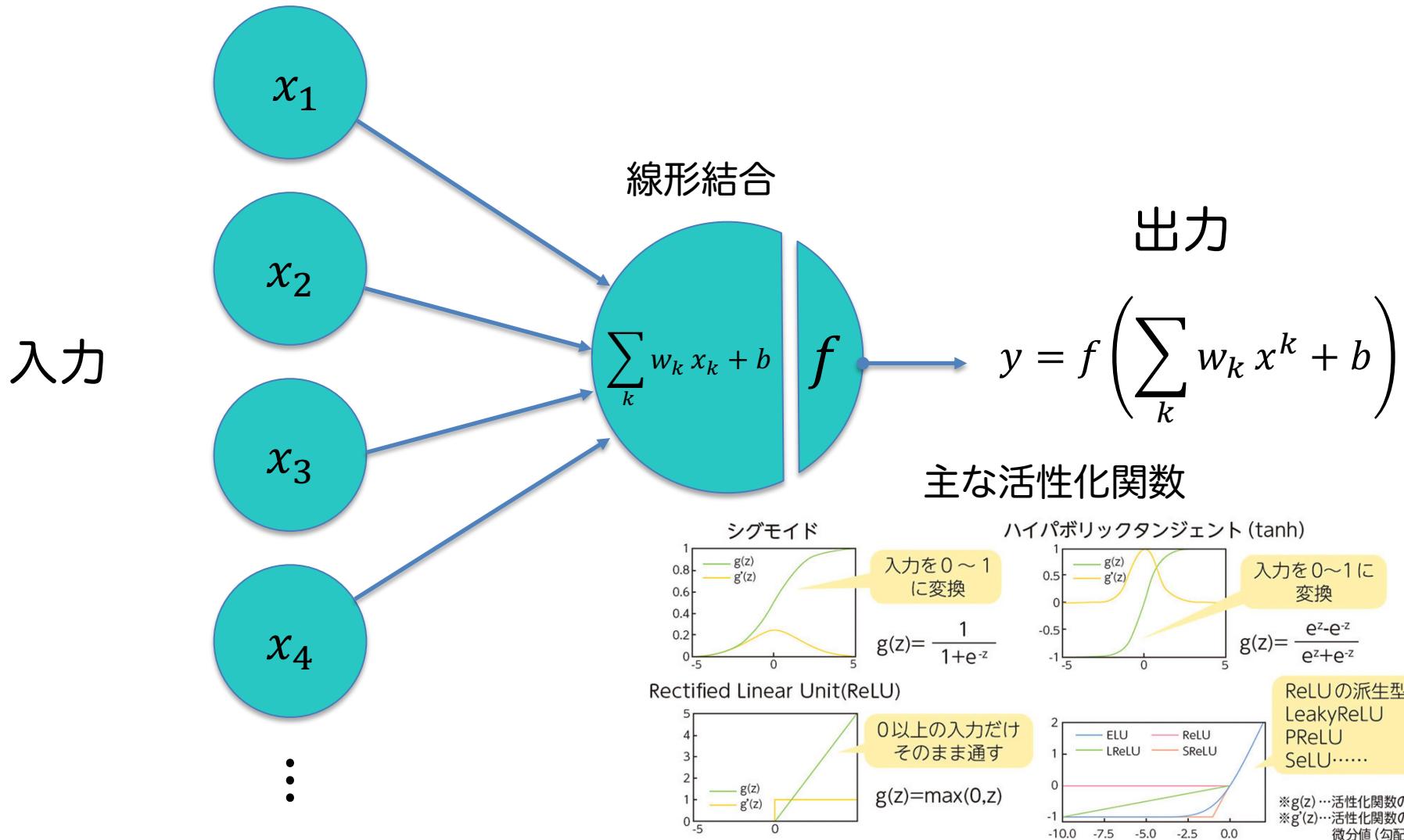


# ニューラルネットの事例 ②

特徴抽出

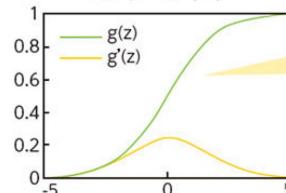


# 活性化関数



# 活性化関数の効果

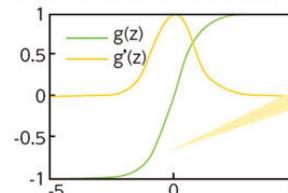
シグモイド



$$g(z) = \frac{1}{1+e^{-z}}$$

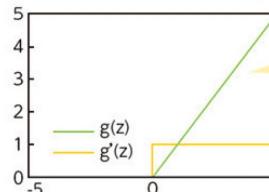
入力を0～1に変換

ハイパボリックタンジェント (tanh)



入力を0～1に変換

Rectified Linear Unit(ReLU)

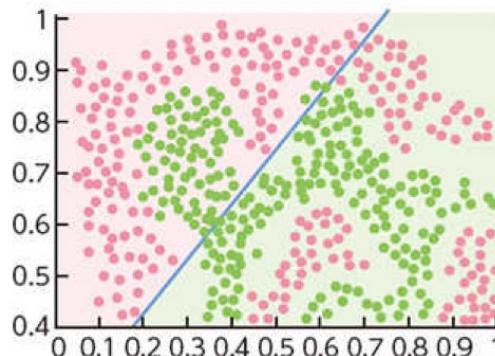


0以上の入力だけそのまま通す  
 $g(z) = \max(0, z)$

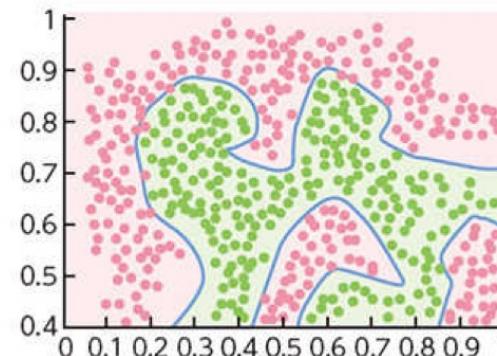
$$g(z) = \max(0, z)$$

ReLUの派生型  
LeakyReLU  
PReLU  
SELU.....

※ $g(z)$ …活性化関数の値  
※ $g'(z)$ …活性化関数の微分値(勾配)

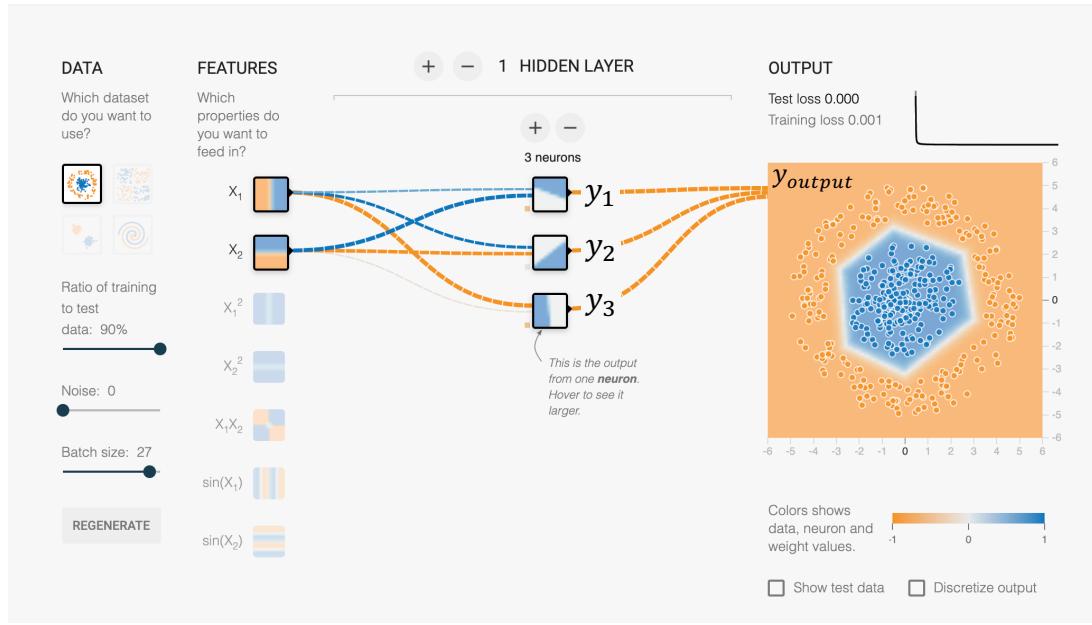


線形活性化関数では層をディープにしても線形

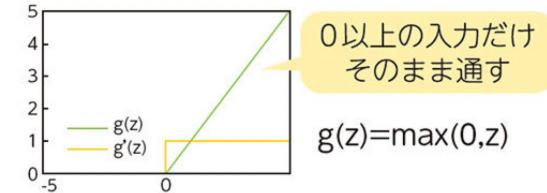


非線形活性化関数では層をディープにすると非常に複雑になる

# ニューラルネットを実感しましょう！



Rectified Linear Unit(ReLU)



$$y_1 = \max(0, w_1^1 x_1 + w_2^1 x_2 + b^1)$$

$$y_2 = \max(0, w_1^2 x_1 + w_2^2 x_2 + b^2)$$

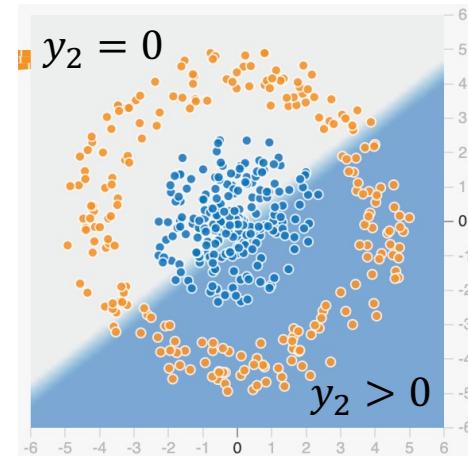
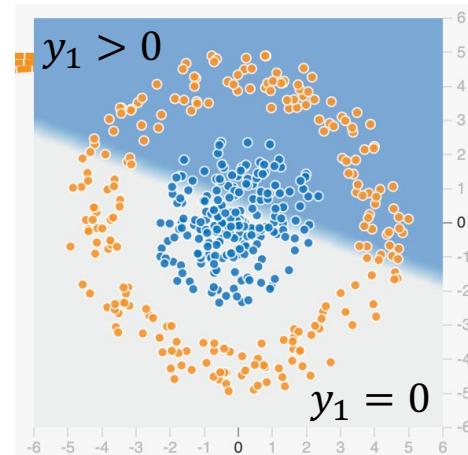
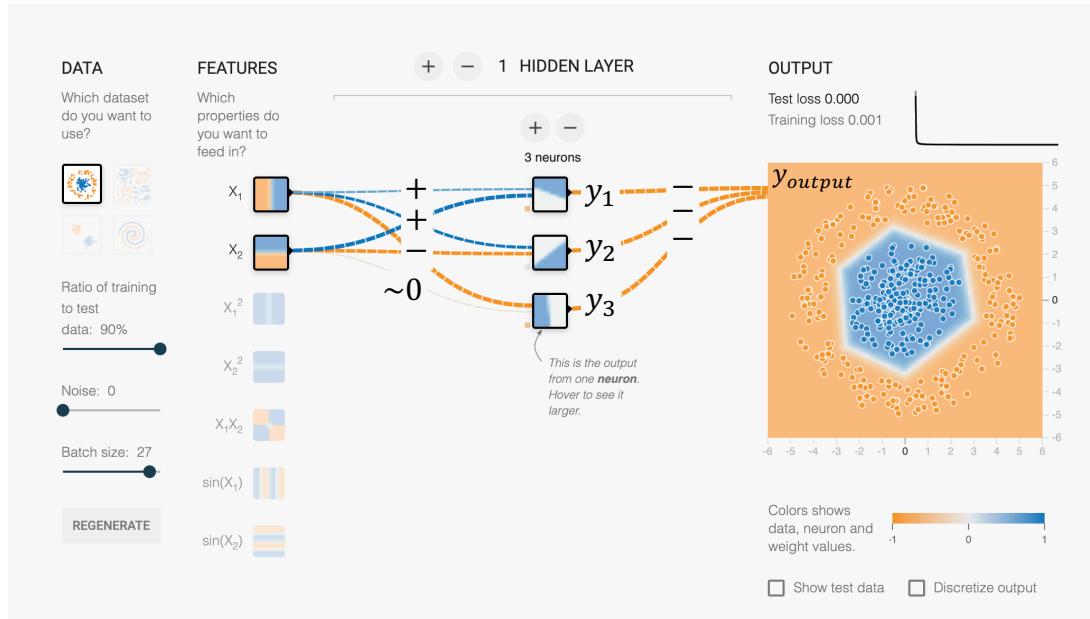
$$y_3 = \max(0, w_1^3 x_1 + w_2^3 x_2 + b^3)$$



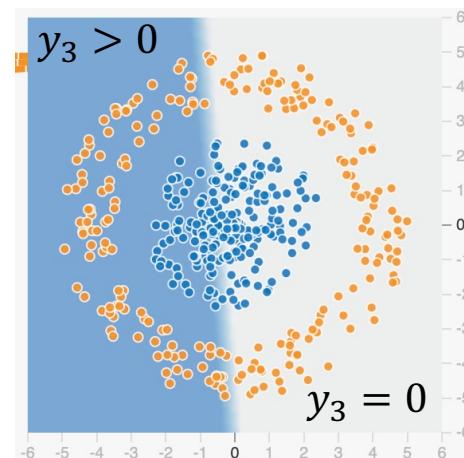
$$y_{output} = \max(0, w_1^4 y_1 + w_2^4 y_2 + w_3^4 y_3 + b^4)$$

<https://playground.tensorflow.org/>

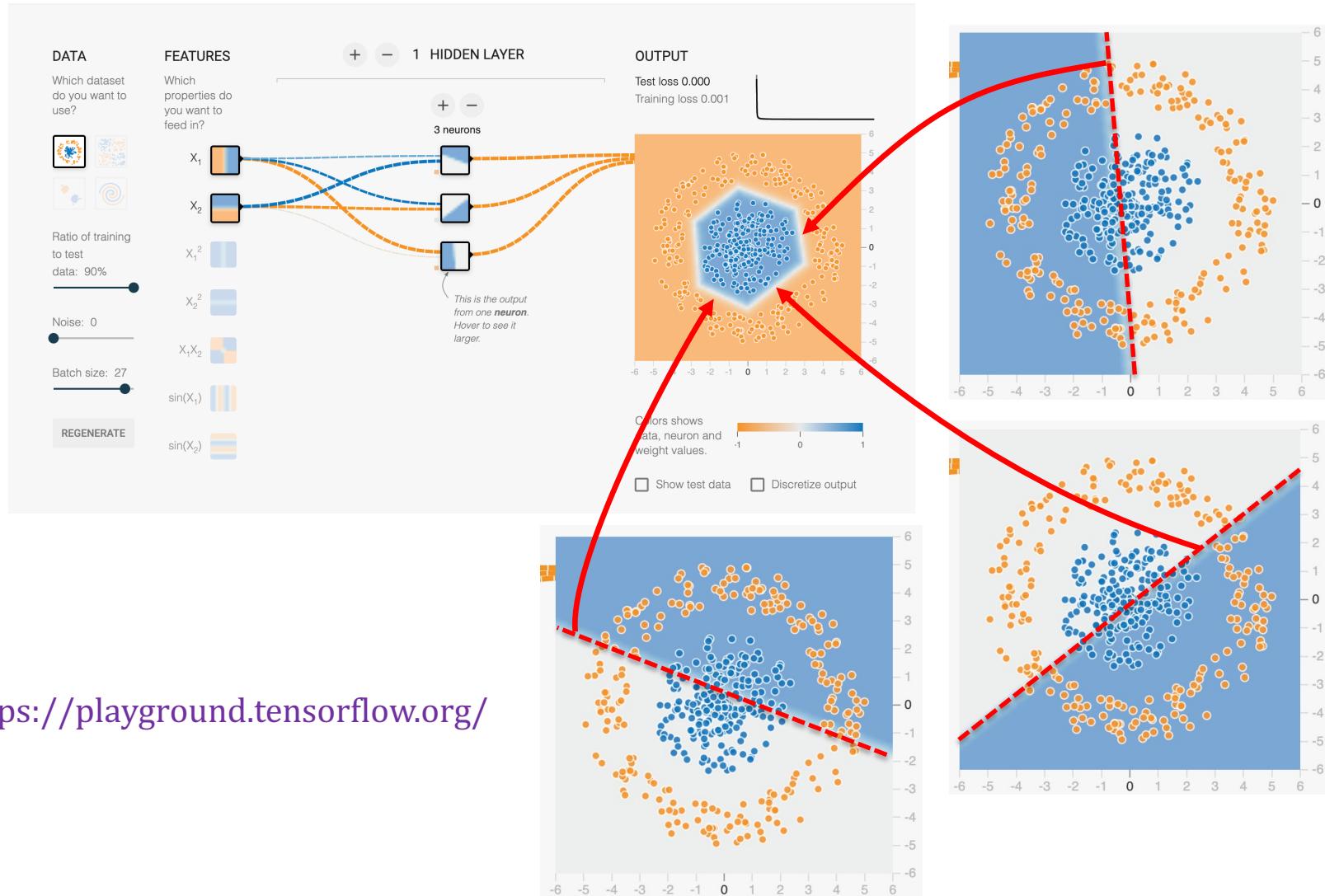
# ニューラルネットを実感しましょう！



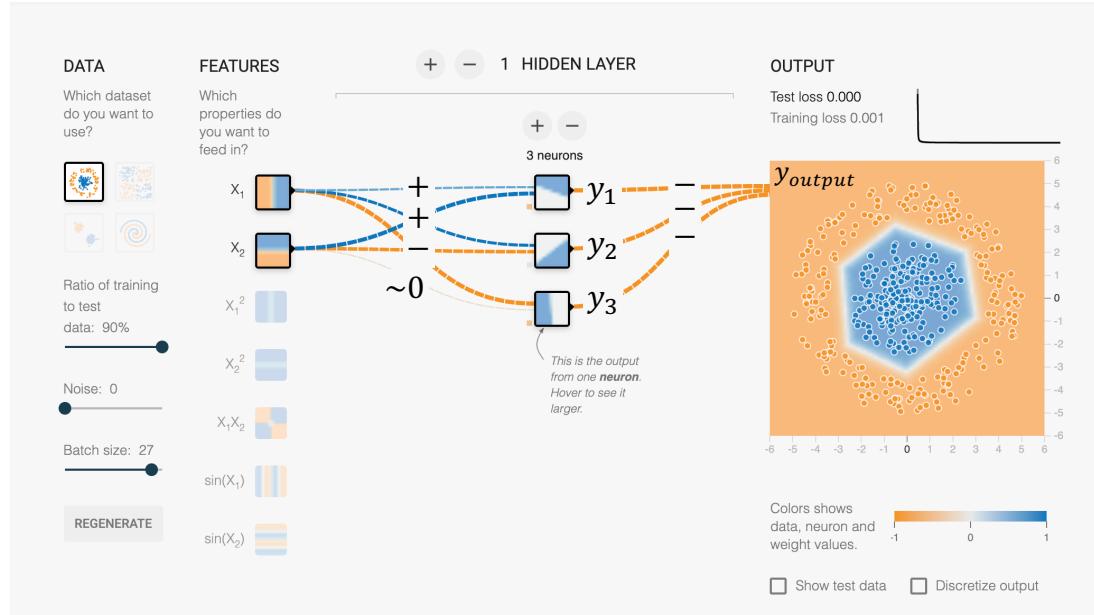
<https://playground.tensorflow.org/>



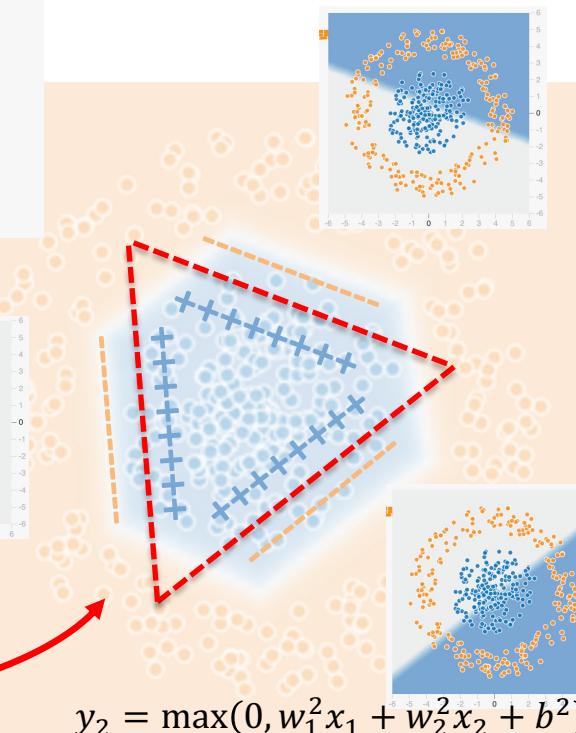
# ニューラルネットを実感しましょう！



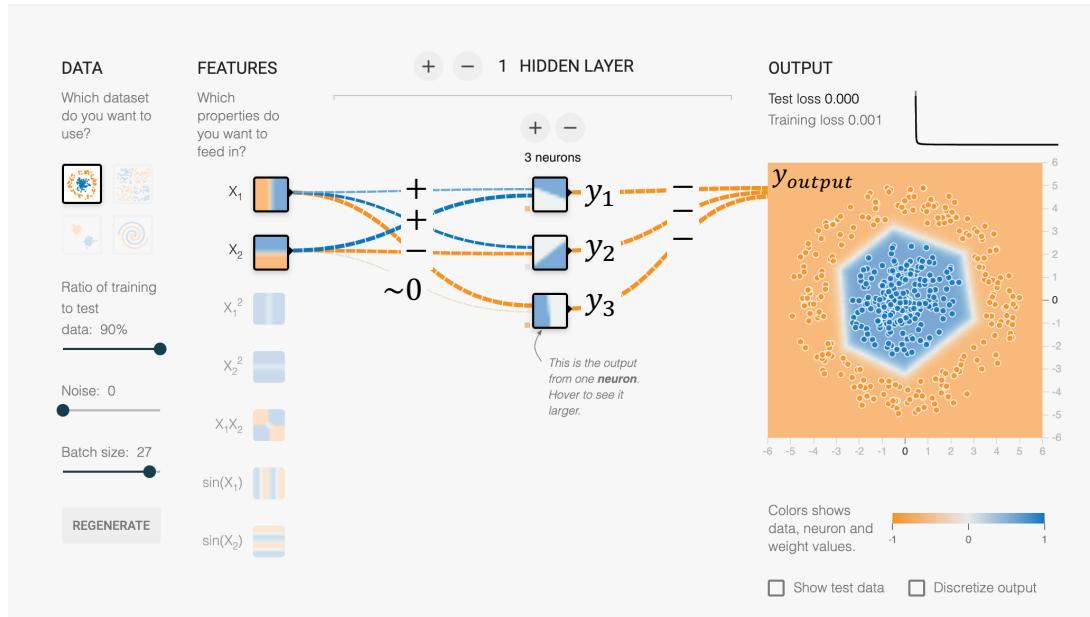
# ニューラルネットを実感しましょう！



$$y_1 = \max(0, w_1^1 x_1 + w_2^1 x_2 + b^1)$$



# ニューラルネットを実感しましょう！



$$y_1 = \max(0, w_1^1 x_1 + w_2^1 x_2 + b^1)$$

$$w_1^4 y_1 + w_2^4 y_2 + w_3^4 y_3 + b^4 < 0$$

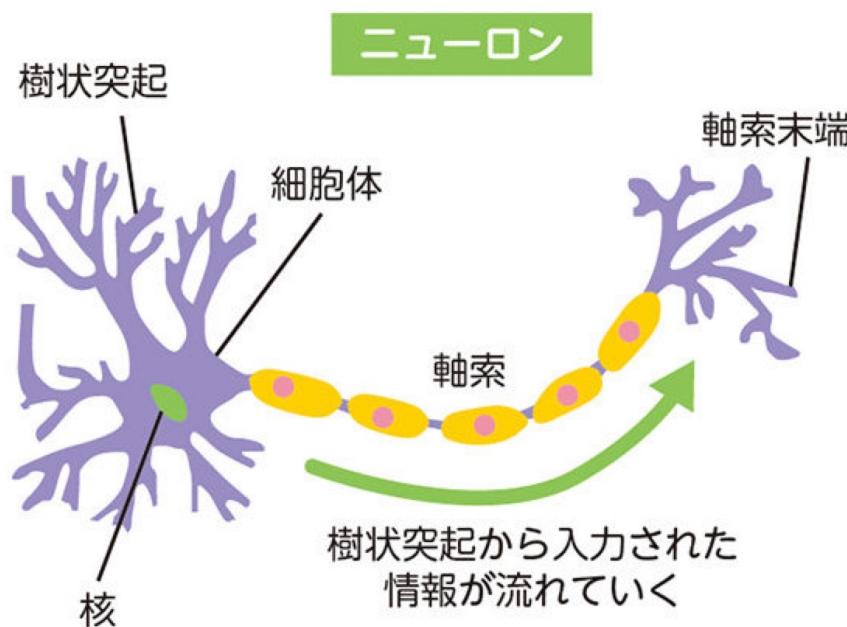
$$y_3 = \max(0, w_1^3 x_1 + w_2^3 x_2 + b^3)$$

$$y_{output} = \max(0, w_1^4 y_1 + w_2^4 y_2 + w_3^4 y_3 + b^4)$$

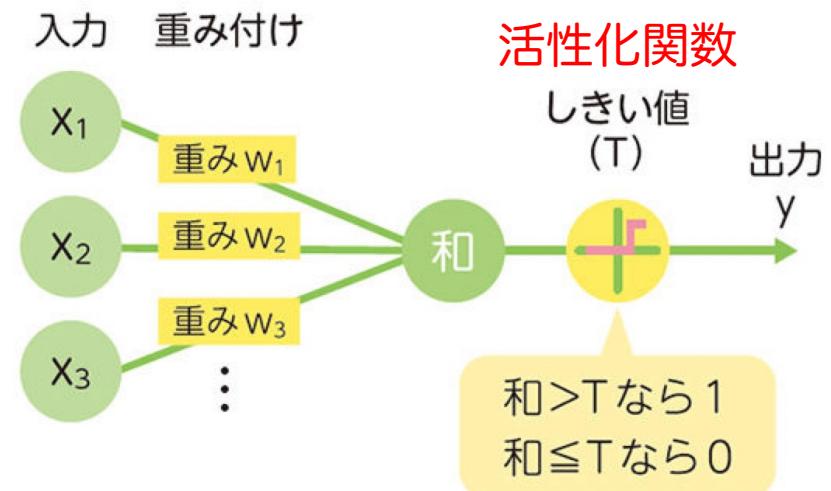
$$y_2 = \max(0, w_1^2 x_1 + w_2^2 x_2 + b^2)$$

# ニューロン・パーセプトロン

## ■ 1960年代までの研究



## 形式ニューロン・パーセプトロン



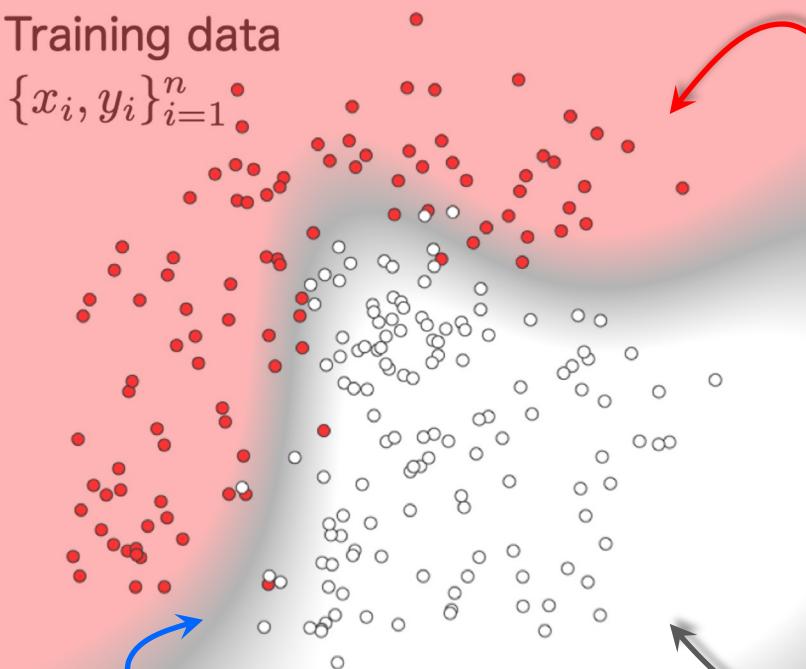
出典：図解即戦力 機械学習&ディープラーニングのしくみと技術がこれ1冊で  
しっかりわかる教科書 株式会社アイデミー(著), 山口 達輝(著), 松田 洋之(著)

# 教師あり学習（復習）

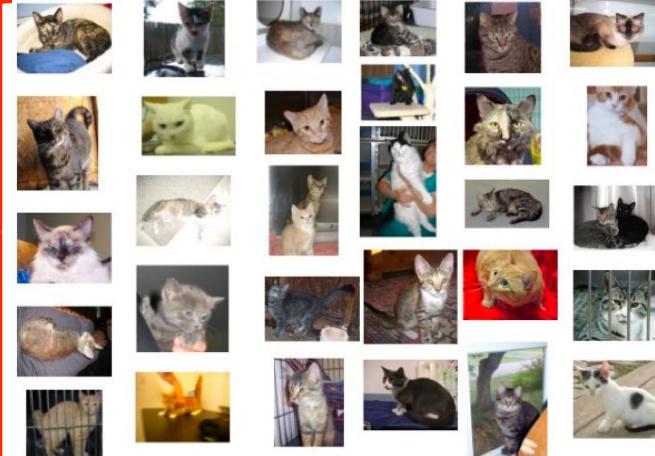
データ変換  
(属性, 特徴量など)

Training data

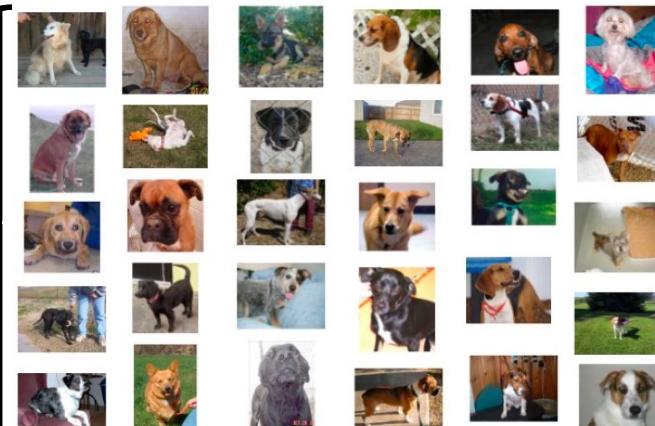
$$\{x_i, y_i\}_{i=1}^n$$



$$y = f \left( \sum_k w_k x^k + b \right)$$



Data source: Kaggle



猫

教師信号

犬

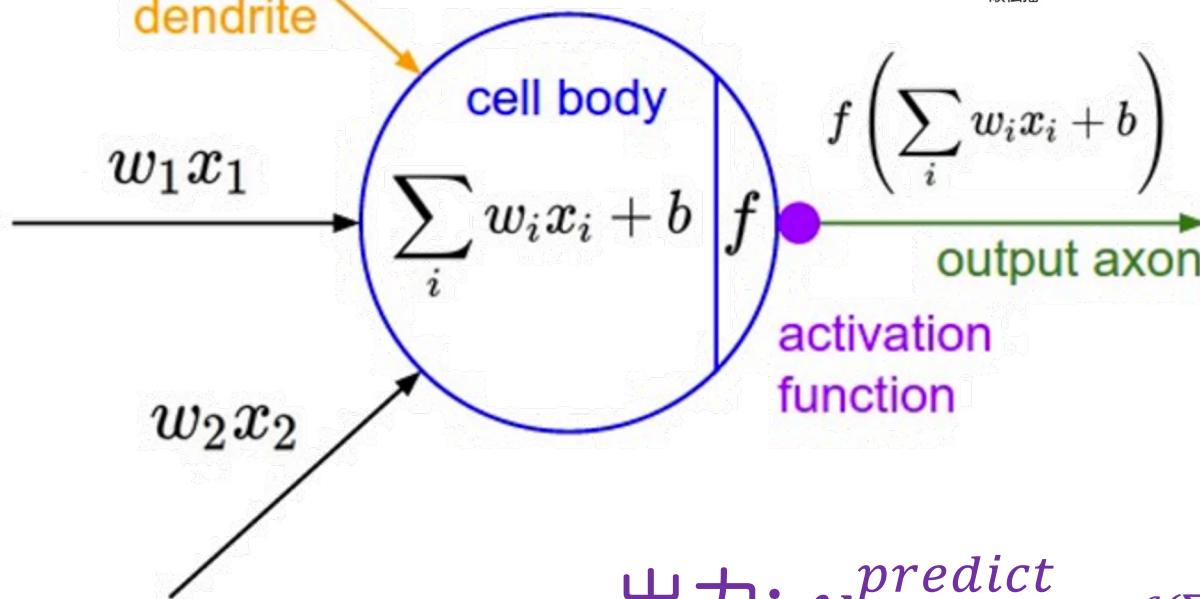
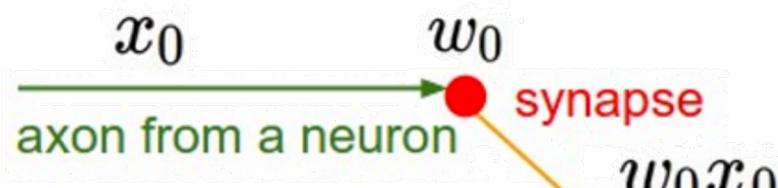
# ニューラルネットの学習（1）

訓練データ

$$\{\vec{x}_i, y_i\}$$

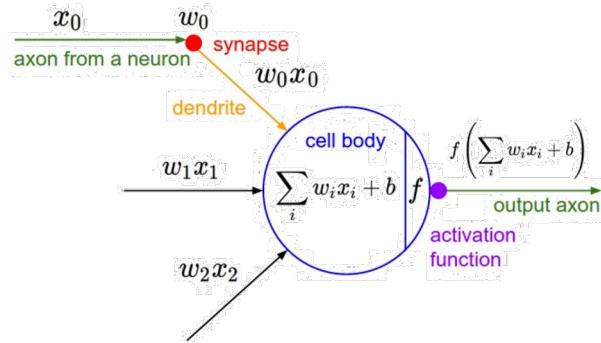
誤差の和 = 損失関数

(小さくなるように  $w_i$  と  $b$  を学習する)



出力:  $y_i^{predict} = f(\sum w_i x_i + b)$   
誤差:  $|y_i - y_i^{predict}|$

# ニューラルネットの学習（2）



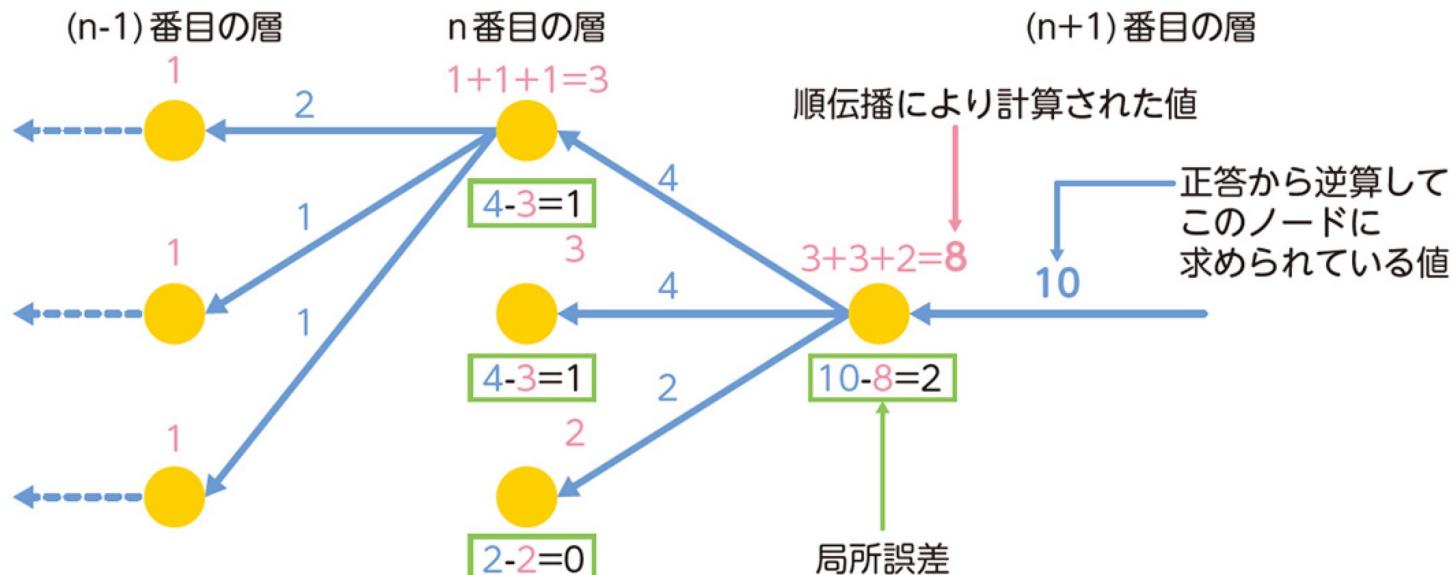
例:  $f(x) = x$

$$y_i^{predict} = f\left(\sum w_j x_i^j + b\right) = \sum w_i x_i^j + b$$

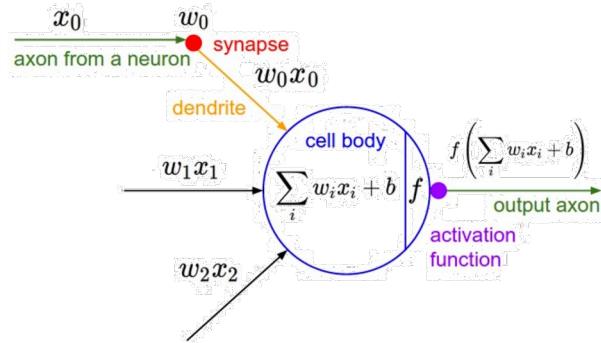
$$\Delta y_i = \sum \Delta w_j x_i^j + \Delta b$$

損失が小さくなるように  $w_j$  と  $b$  を調整する

## ■ 誤差逆伝播法



# ニューラルネットの学習（3）



■ 損失関数を最小化する

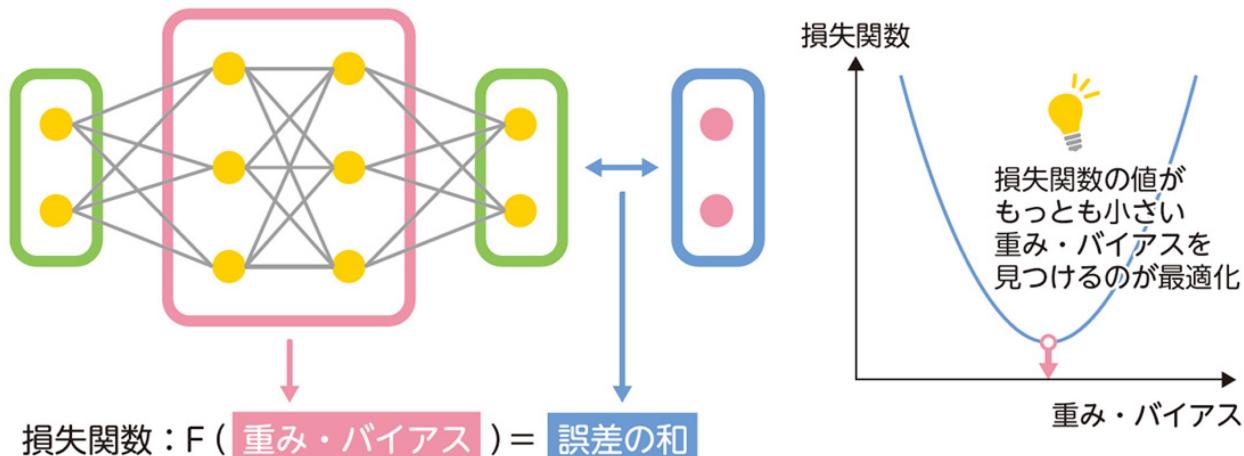
例 :  $f(x) = x$

$$y_i^{predict} = f\left(\sum_j w_j x_j + b\right) = \sum_i w_i x_i + b$$

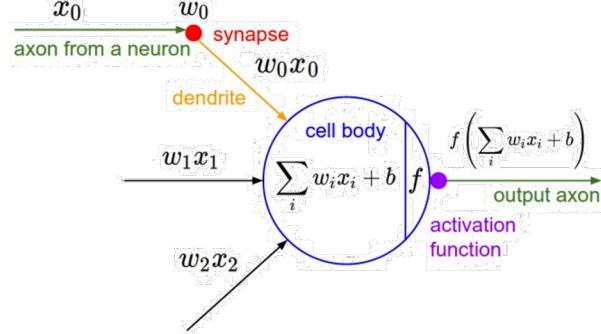
$$\Delta y_i = \sum_j \Delta w_j x_j + \Delta b$$

損失が小さくなるように  $w_j$  と  $b$  を調整する

重み・バイアス



# ニューラルネットの学習 (4)



$$y_i^{predict} = f\left(\sum w_j x_i^j + b\right)$$

$$\Delta y_i = |y_i - y_i^{predict}| = \left|y_i - f\left(\sum w_j x_i^j + b\right)\right|$$

活性化関数  $f(z)$

$$F_i(\vec{w}, b) = \Delta y_i = \left|y_i - f\left(\sum w_j x_i^j + b\right)\right|$$

$$\frac{\delta F_i}{\delta w_j} = \frac{\delta F_i}{\delta f} \frac{\delta f}{\delta z} \frac{\delta z}{\delta w_j} = \pm \frac{\delta f}{\delta z} x_i^j \rightarrow \delta F_i = \pm \frac{\delta f}{\delta z} x_i^j \delta w_j$$

$$\frac{\delta F_i}{\delta b} = \frac{\delta F_i}{\delta f} \frac{\delta f}{\delta z} \frac{\delta z}{\delta b} = \pm \frac{\delta f}{\delta z} \rightarrow \delta F_i = \pm \frac{\delta f}{\delta z} \delta b$$

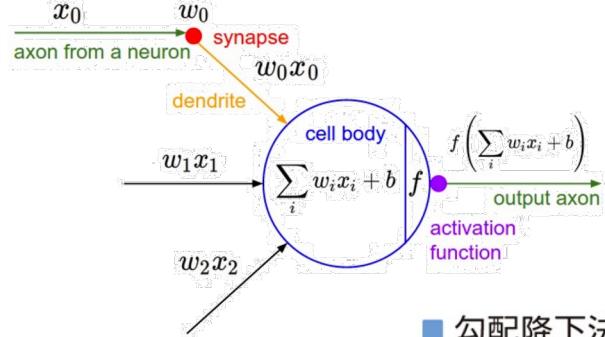
訓練データ  
 $\{\vec{x}_i, y_i\}$

効率的に重みの調整 (小  $\delta w_j$ ) で  
 損失減らしたい (大  $\delta F_i$ )



重み・バイアス  
 $w_i$  を  $\delta w_j$  や  $b$  を  $\delta b$  調整する際に  
 損失が変化する量

# ニューラルネットの学習（5）

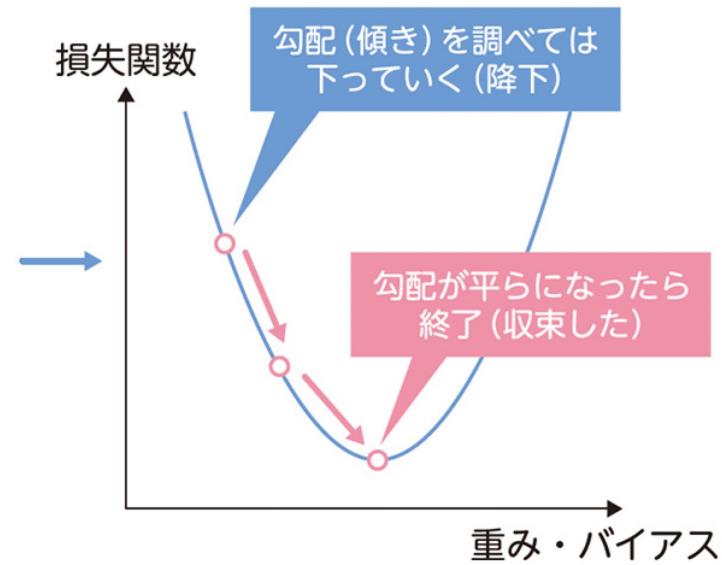


■ 勾配降下法は「山下り」

訓練データ  
 $\{\vec{x}_i, y_i\}$

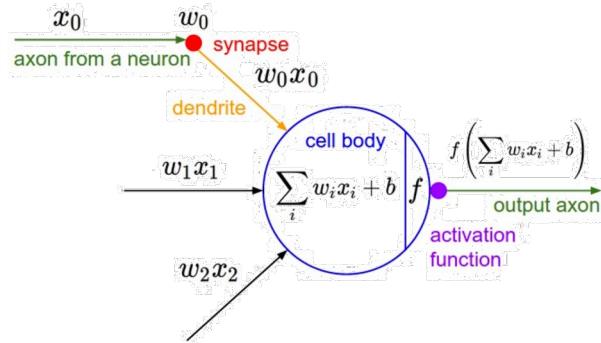


効率的に重みの調整（小  $\delta w_j$ ）で  
 損失減らしたい（大  $\delta F_i$ ）



$w_i$  を  $\delta w_j$  や  $b$  を  $\delta b$  調整する際に  
 損失が変化する量

# ニューラルネットの学習（6）



例:  $f(x) = x$

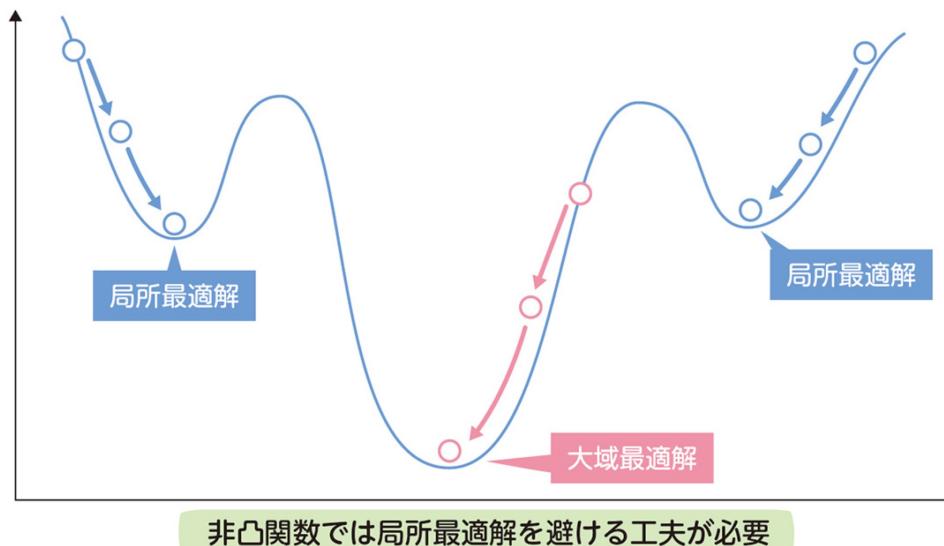
$$y_i^{predict} = f(\sum w_i x_i + b) = \sum w_i x_i + b$$

$$\Delta y = y_i - y_i^{predict} = \sum \Delta w_i x_i + \Delta b$$

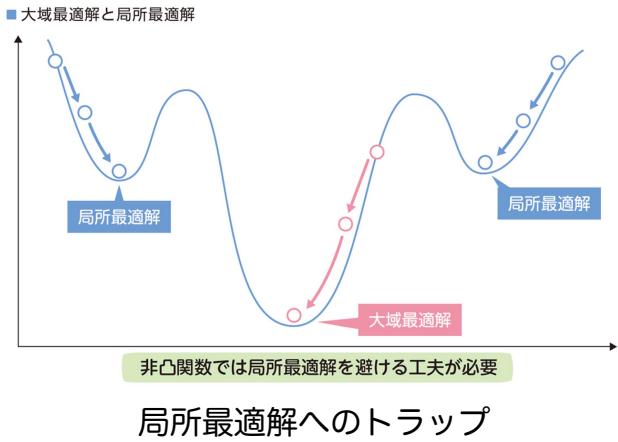
損失が小さくなるように  $w_i$  と  $b$  を調整する

重み・バイアス

■ 大域最適解と局所最適解

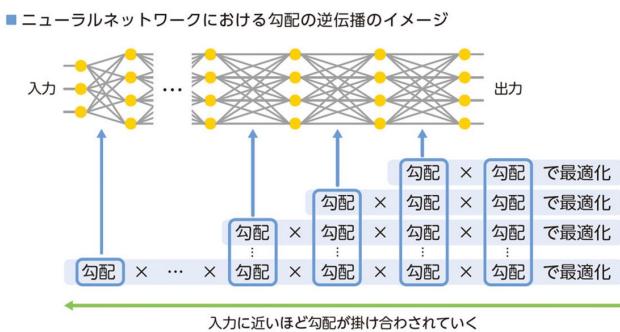


# ニューラルネットの学習（7）



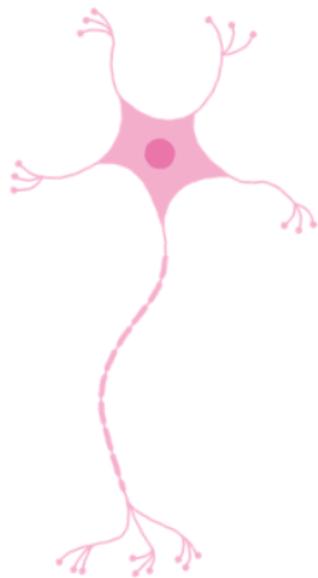
#### ■ 最適化アルゴリズム

- SGD (確率勾配降下法)  
一つ一つの学習データの順番を入れ替えながら勾配降下法を適用する
  - Momentum SGD  
勾配を降りている方向に「運動量」をもたせる
  - Adagrad  
パラメターごとに適切に学習率を変化させる
  - RMSprop  
勾配の合計を指数移動平均を利用してAdagradを改良する
  - Adam  
RMSpropとMomentum SGDを組み合わせたアルゴリズム

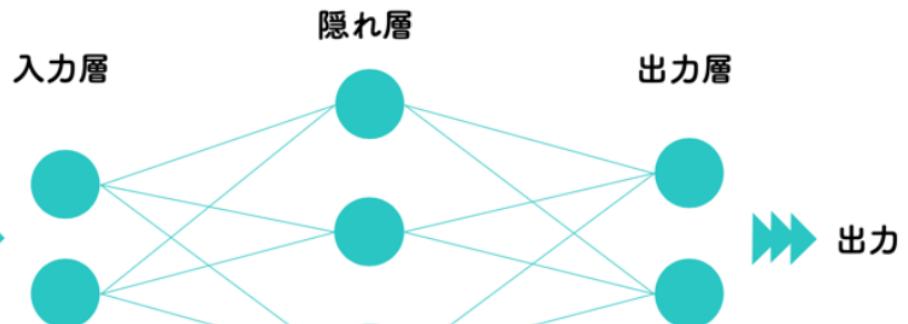


# ニューラルネットのまとめ

## 神経細胞(ニューロン)

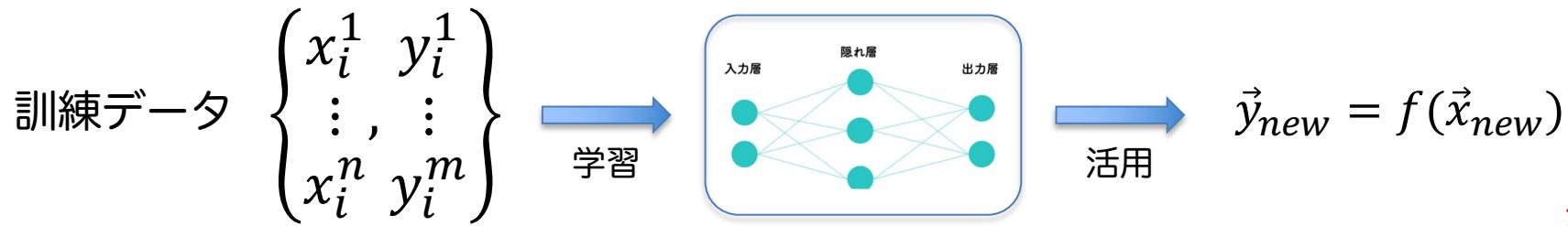


## ニューラルネットワーク



## 基本機能

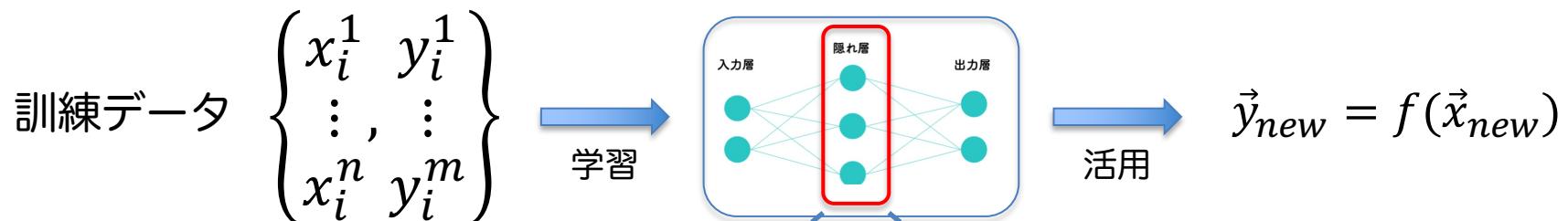
$$\vec{y} = f(\vec{x})$$



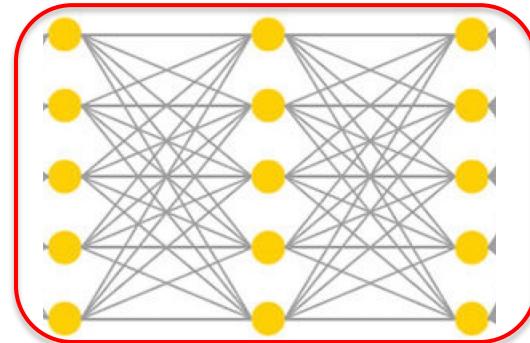
# 深層学習への展開

ニューラルネット

## 基本機能



②  
基本機能の繋がりの設計により柔軟に問題を提起できる



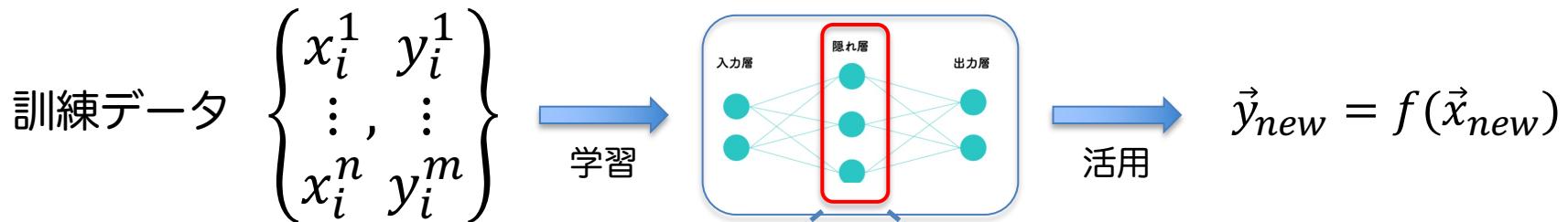
③  
 $\vec{y}$ の定義仕方で適切に問題を提起できる

①  
隠れ層の数を増やし、多種の活性化関数の使用によって表現できる関数を複雑化できる

# 深層学習：複雑関数の表現

ニューラルネット

## 基本機能

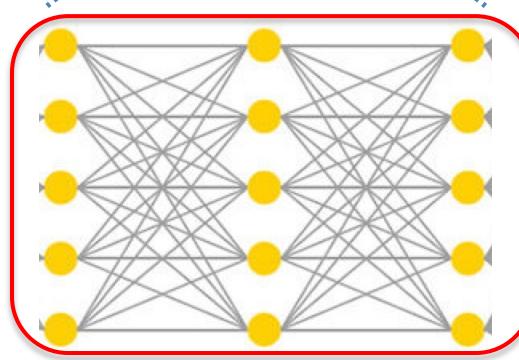


$k$  層目の出力  $x_j^k = f_j^k \left( \sum w_j^k x_j^{k-1} + b_k \right)$

活性化関数  $f_j^k$

$k$  層目の重み  $w_j^k$

$k$  層目のバイアス  $b_k$



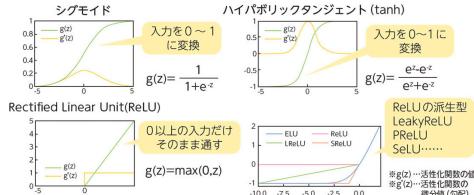
$$y = f^{\textcolor{red}{n}} \left( b_{\textcolor{red}{n}} + \sum w_j^{\textcolor{red}{n-1}} f_j^{\textcolor{red}{n-1}} \right)$$

$$\left( b_{\textcolor{red}{n-1}} + \sum w_j^{\textcolor{red}{n-2}} f_j^{\textcolor{red}{n-2}} \right)$$

$$\left( b_{\textcolor{red}{n-2}} + \sum w_j^{\textcolor{red}{n-3}} f_j^{\textcolor{red}{n-3}} \right) \dots$$

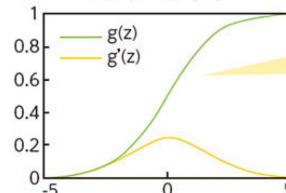
①

隠れ層の数を増やし、多種の活性化関数の使用によって表現できる関数を複雑化できる



# 活性化関数の効果

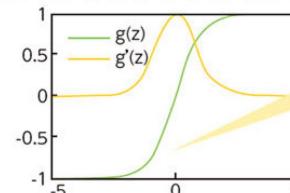
シグモイド



入力を0～1に  
変換

$$g(z) = \frac{1}{1+e^{-z}}$$

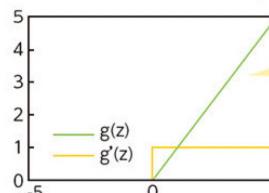
ハイパボリックタンジェント (tanh)



入力を0～1に  
変換

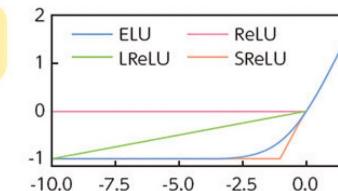
$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Rectified Linear Unit(ReLU)



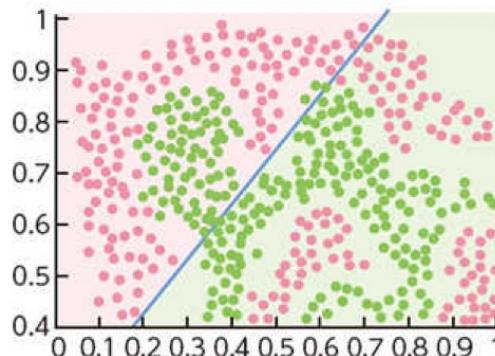
0以上の入力だけ  
そのまま通す

$$g(z) = \max(0, z)$$

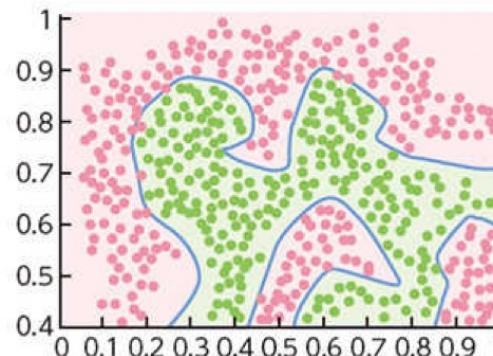


ReLUの派生型  
LeakyReLU  
PReLU  
Selu.....

※ $g(z)$ …活性化関数の値  
※ $g'(z)$ …活性化関数の  
微分値(勾配)



線形活性化関数では  
層をディープに  
しても線形



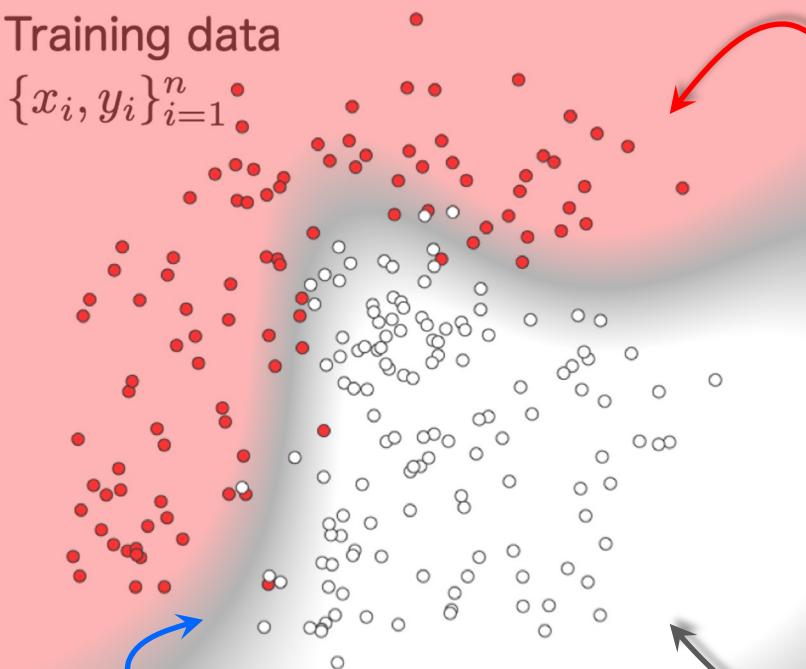
非線形活性化関数では  
層をディープに  
すると非常に  
複雑になる

# 教師あり学習（復習）

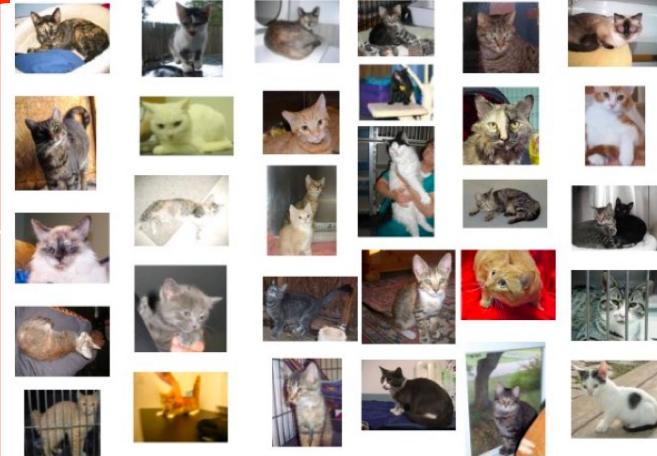
データ変換  
(属性, 特徴量など)

Training data

$$\{x_i, y_i\}_{i=1}^n$$

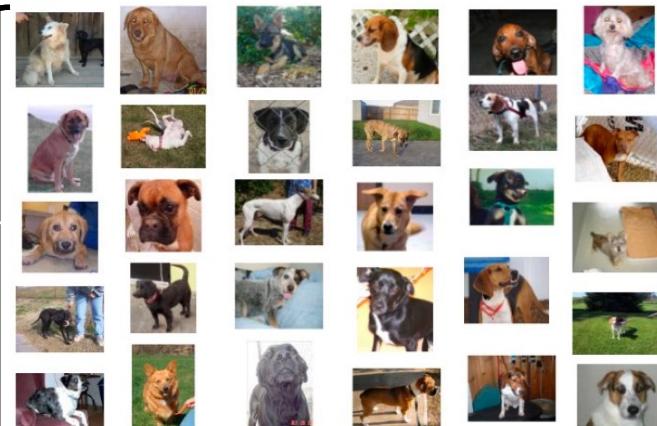


$$y = f \left( \sum_k w_k x^k + b \right)$$



猫

Data source: Kaggle



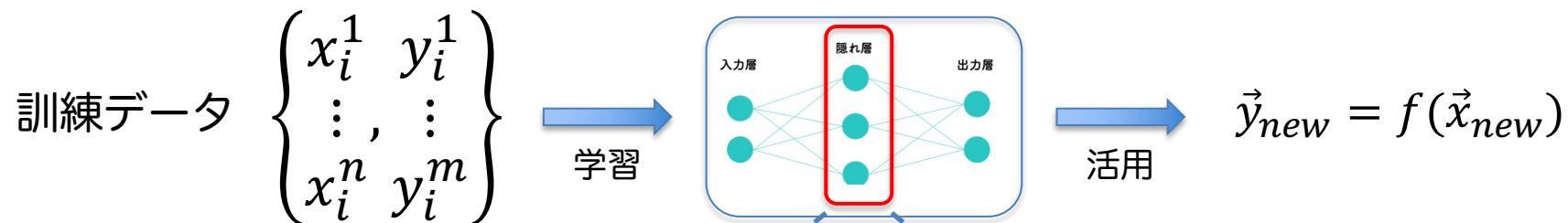
教師信号

犬

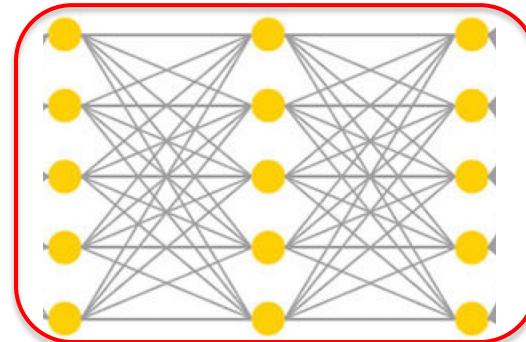
# 深層学習への展開

ニューラルネット

## 基本機能



②  
基本機能の繋がりの設計により柔軟に問題を提起できる



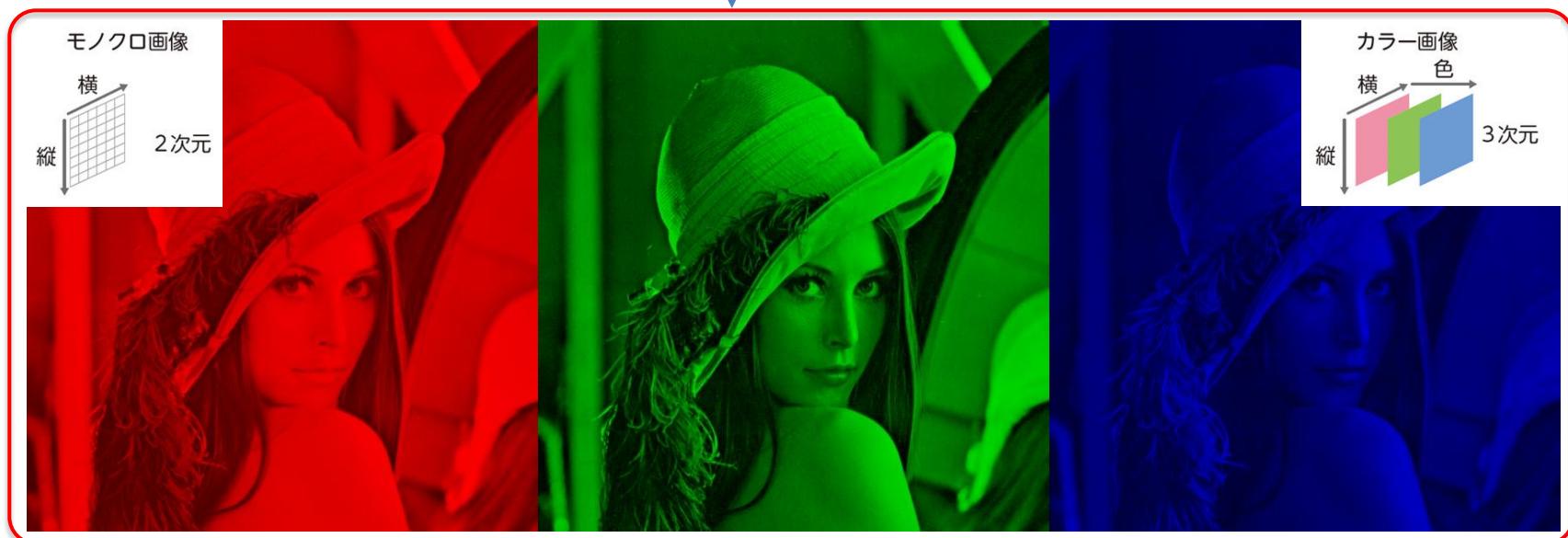
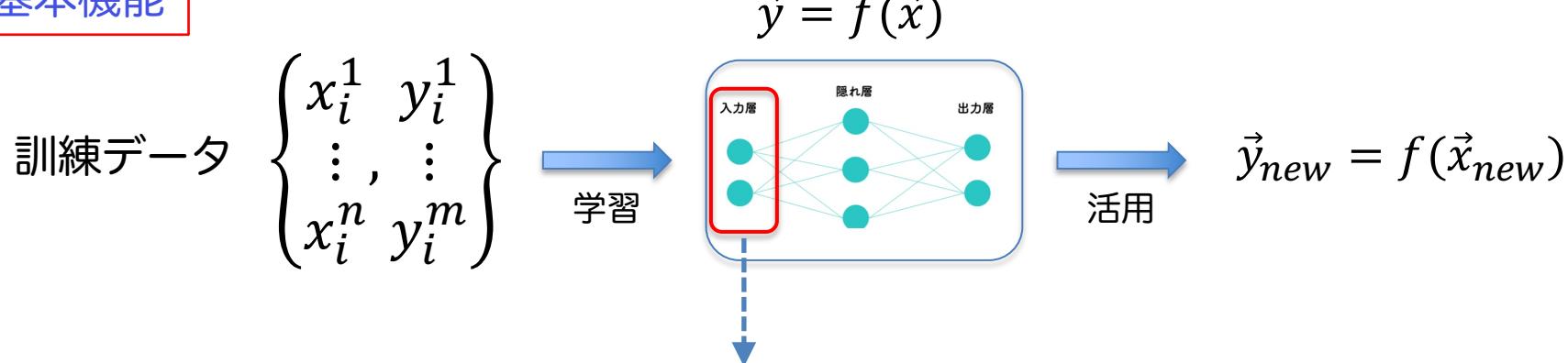
①  
隠れ層の数を増やし、多種の活性化関数の使用によって表現できる関数を複雑化できる

③  
 $\vec{y}$ の定義仕方で適切に問題を提起できる

# 深層学習の画像処理への展開（1）

ニューラルネット

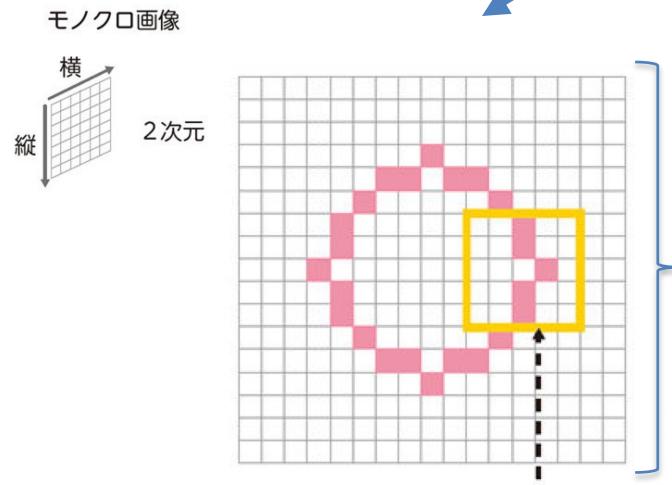
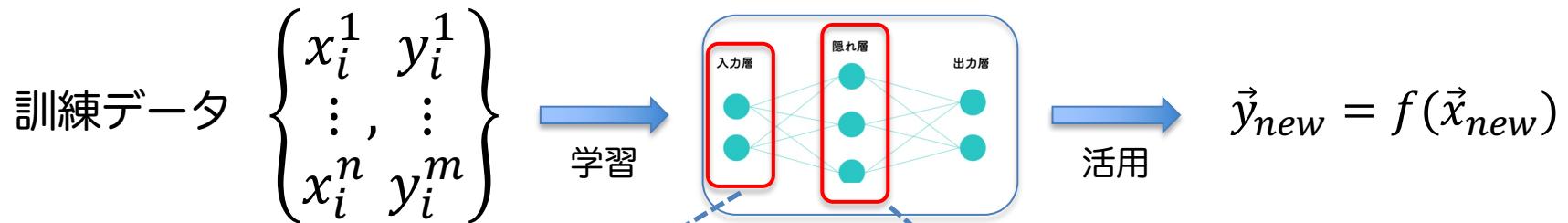
## 基本機能



# 深層学習の画像処理への展開（2）

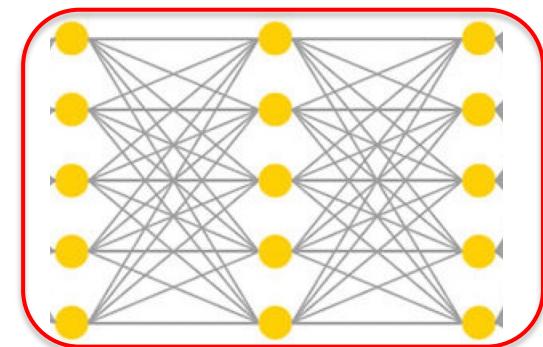
ニューラルネット

## 基本機能



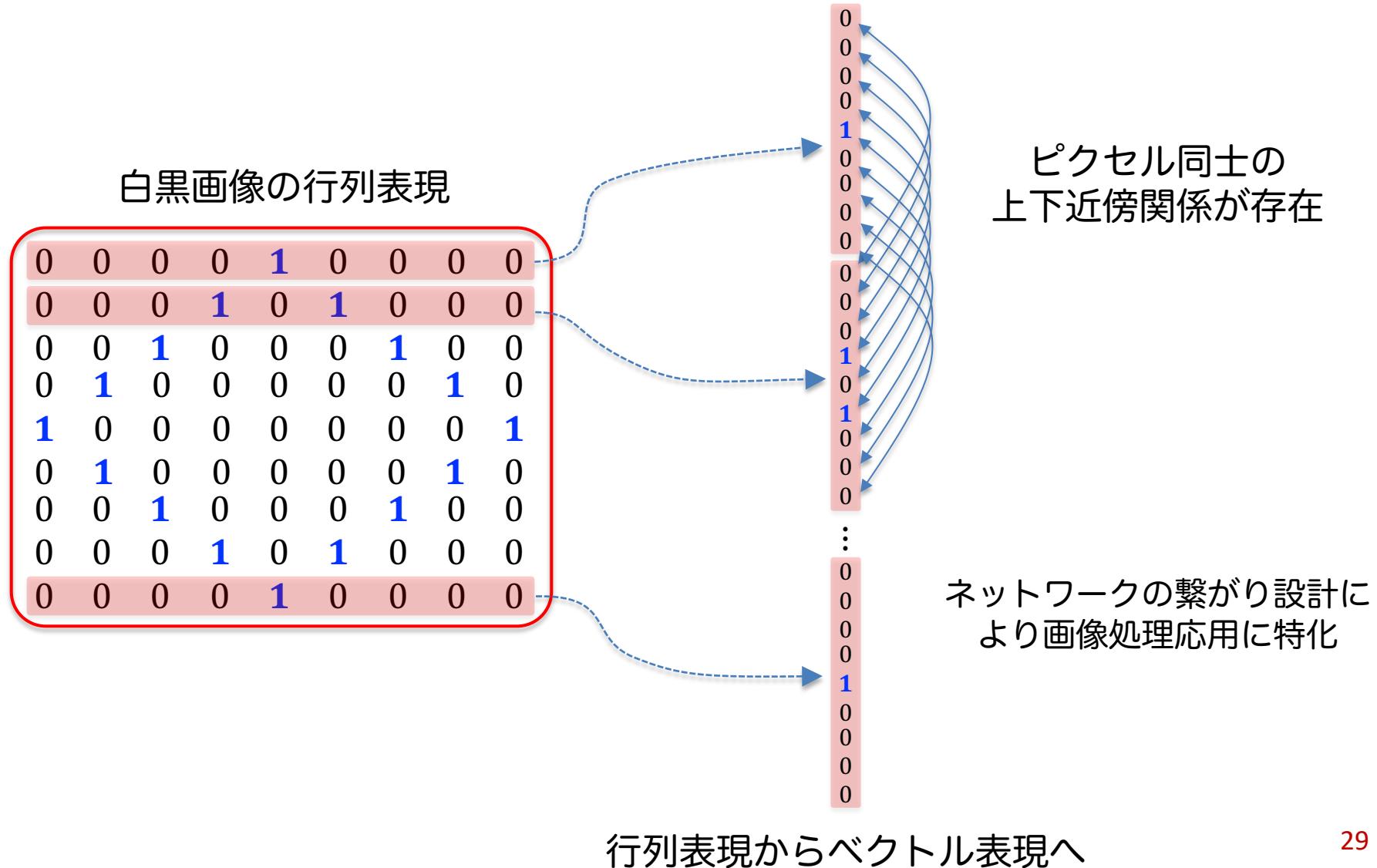
特化した  
ネットワークの設計

固定グリッド  
数に変換

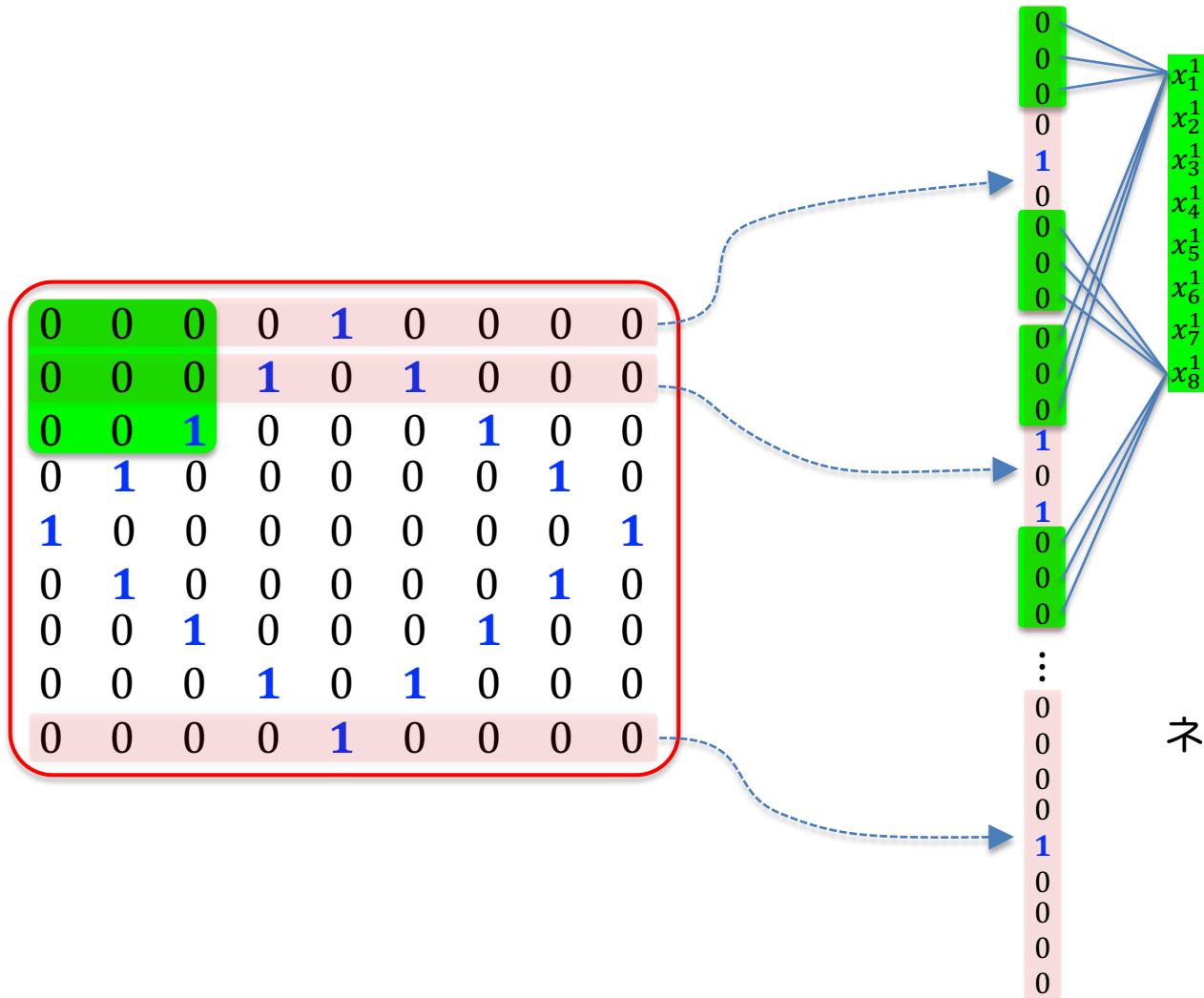


ピクセル同士の「位置関係」が重要

# 画像処理用ネットワークの設計（1）



# 画像処理用ネットワークの設計（2）

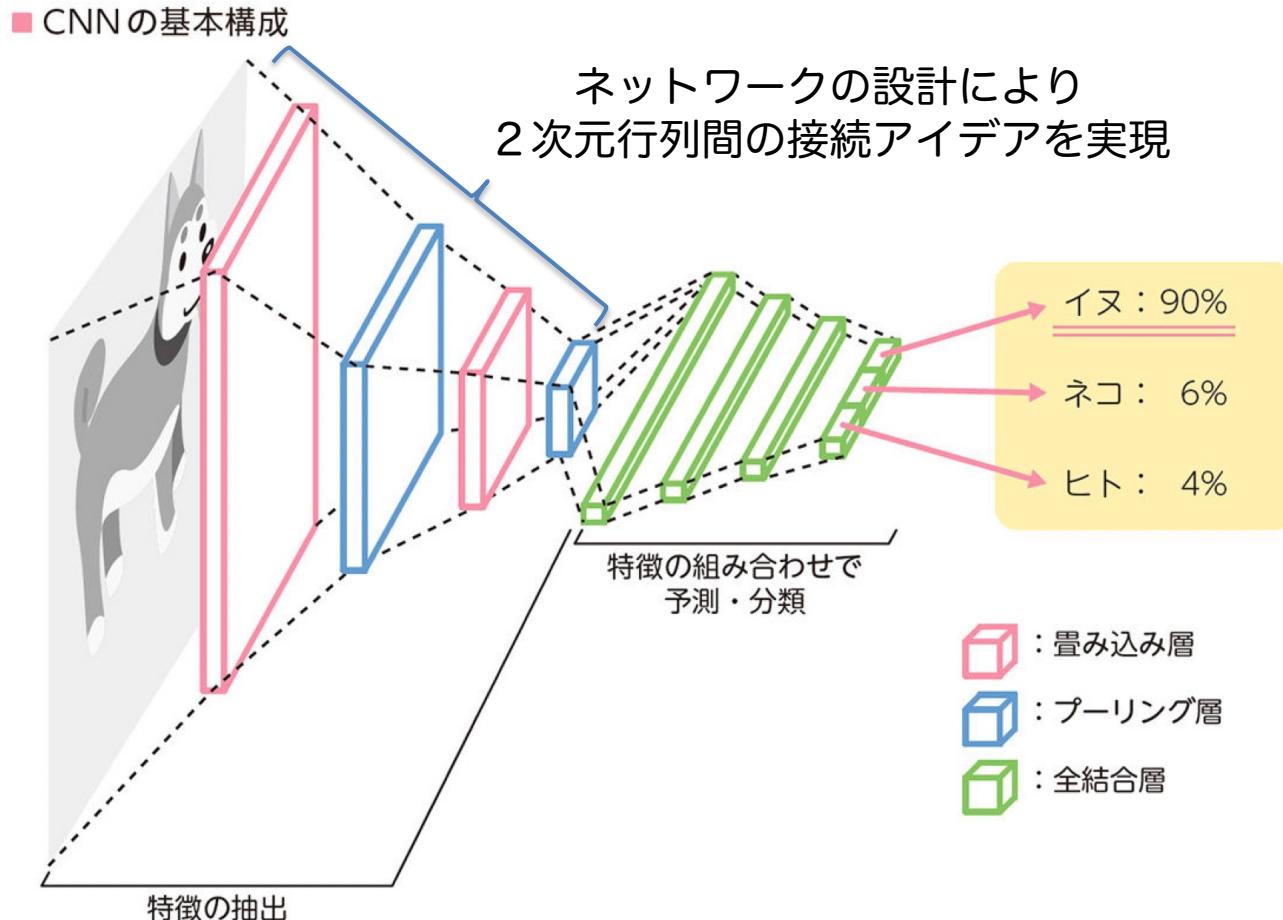


ピクセル同士の  
上下近傍関係の考慮

ネットワークの繋がり設計に  
より画像処理応用に特化

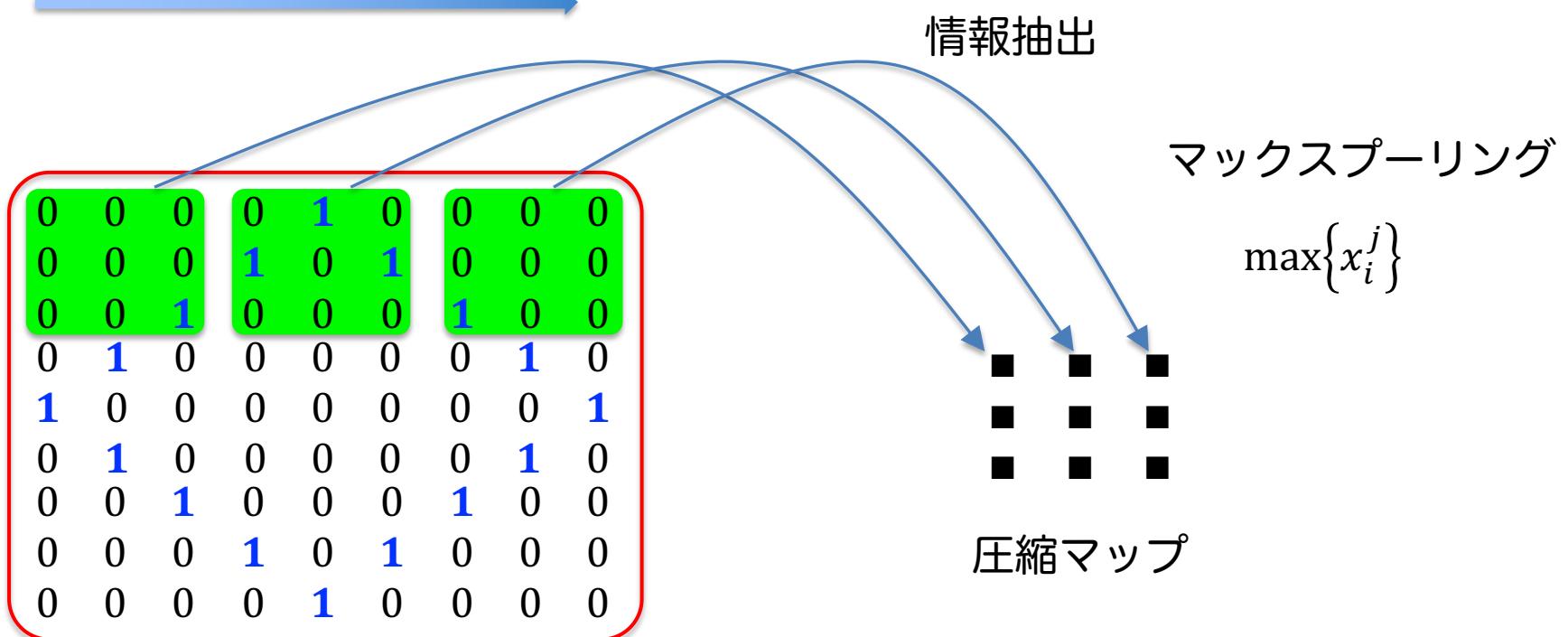
行列表現からベクトル表現へ

# 畳み込みニューラルネットワーク（1）

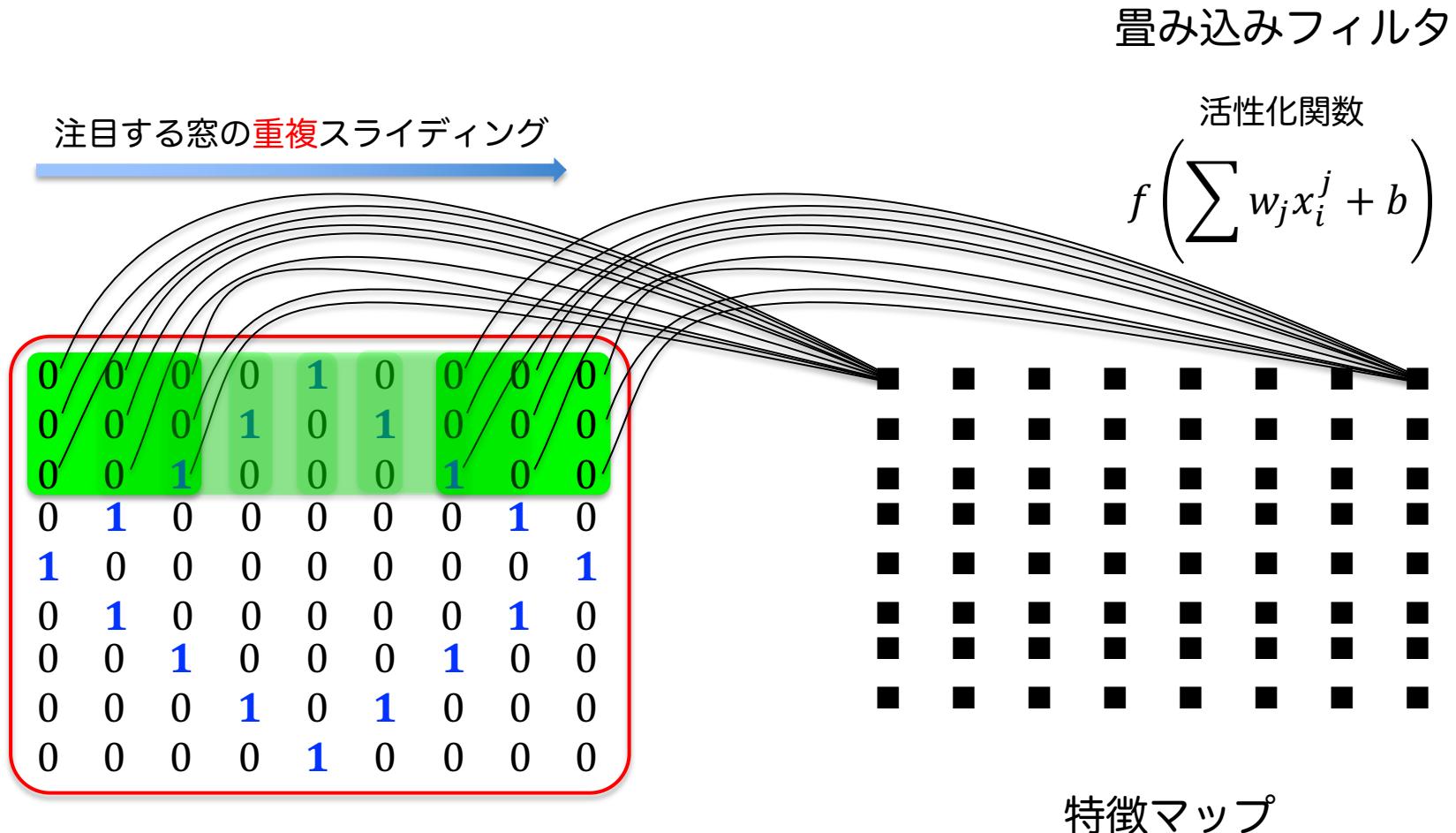


# プーリング層

注目する窓の無重複スライディング



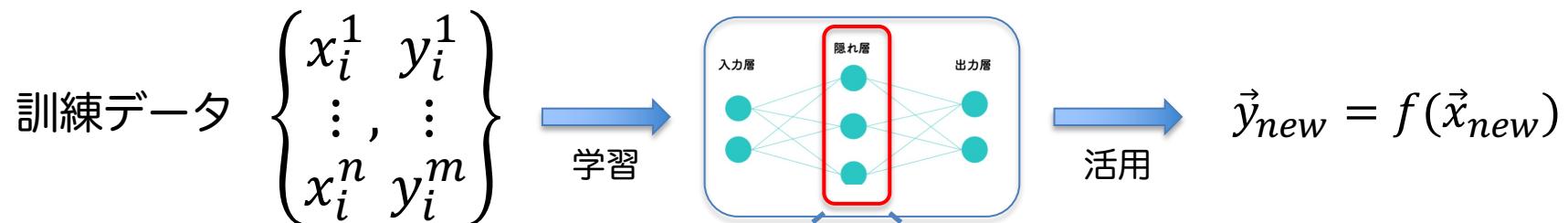
# 畳み込み層



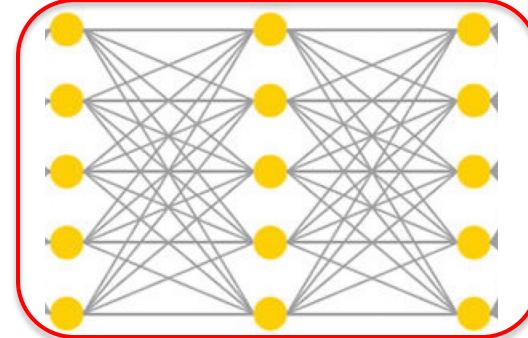
# 深層学習への展開

ニューラルネット

## 基本機能



②  
基本機能の繋がりの設計により柔軟に問題を提起できる



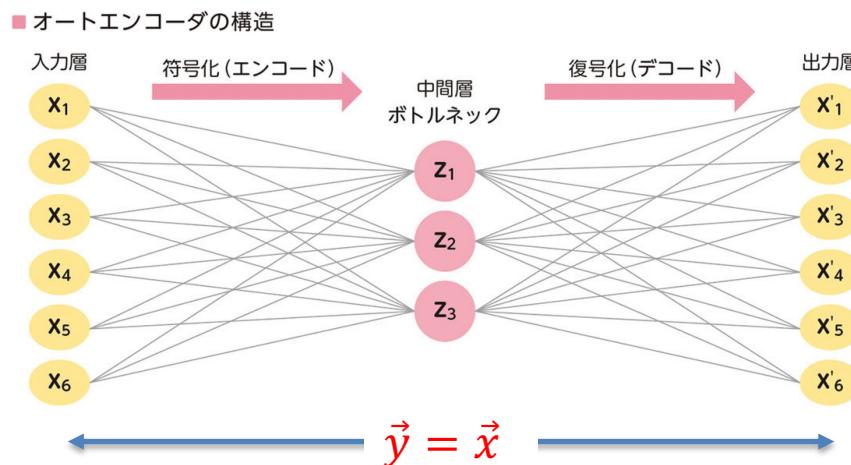
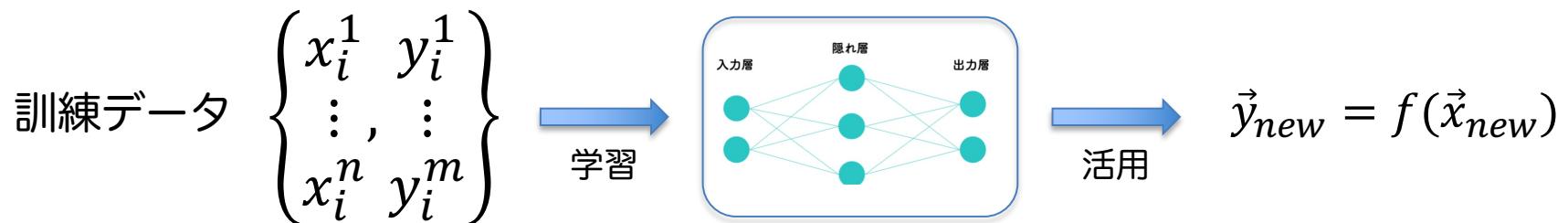
①  
隠れ層の数を増やし、多種の活性化関数の使用によって表現できる関数を複雑化できる

③  
 $\vec{y}$ の定義仕方で適切に問題を提起できる

# オートエンコーダ（1）

ニューラルネット

## 基本機能



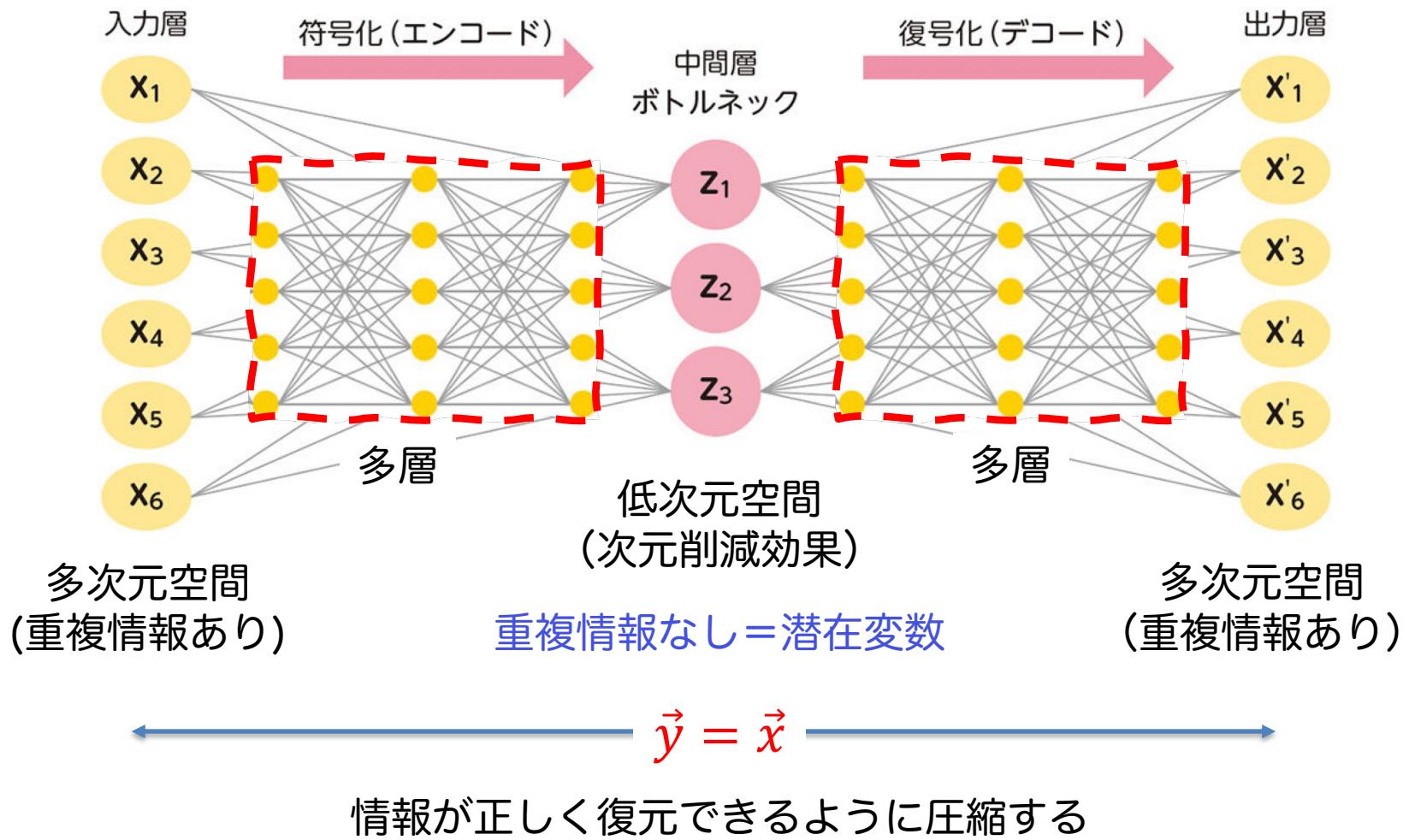
$\vec{y}$  の柔軟な定義による  
教師なし学習

③ $\vec{y}$ の定義仕方で適切に問  
題を提起できる

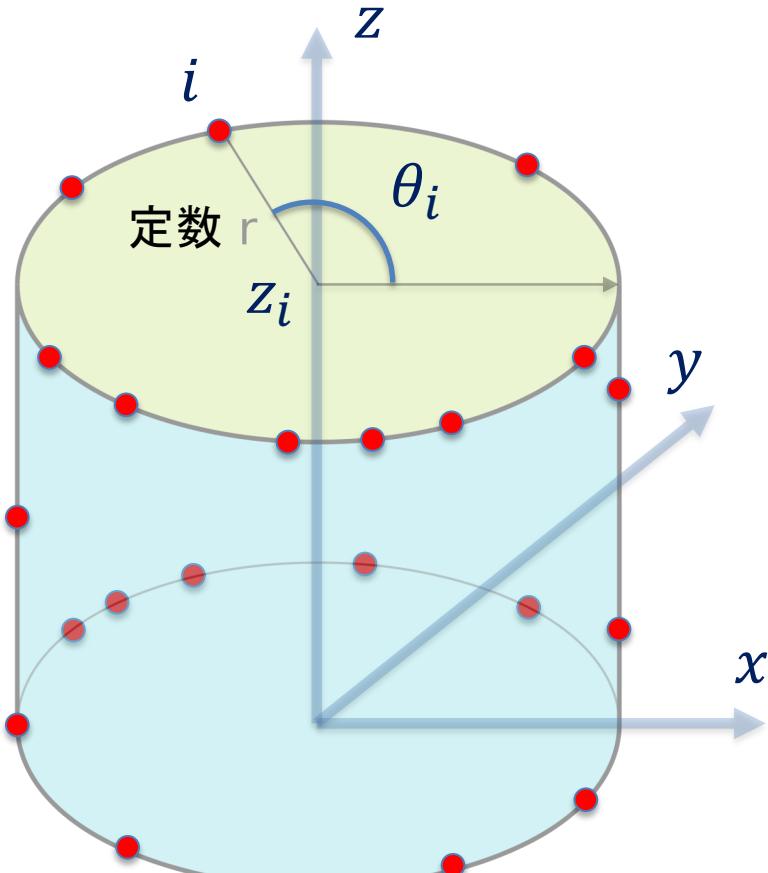
入力したデータと同じデータが出力されるニューラルネット

# オートエンコーダ (2)

## ■ オートエンコーダの構造



# 潜在変数の例（1）



半径が定数の円柱の表面

直交座標系

円柱座標系

座標変換

データ

$$\begin{cases} x_i \\ y_i \\ Z_i \end{cases}$$

3次元

半径が定数の円柱表面にある束縛条件が利用されていない

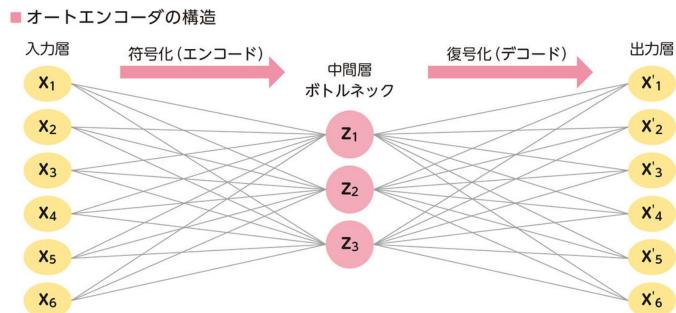
$$\begin{cases} \theta_i \\ Z_i \end{cases}$$

2次元

半径が定数の円柱表面にある束縛条件が利用する潜在変数

$$\begin{cases} x_i = r \cos \theta_i \\ y_i = r \sin \theta_i \\ Z_i = Z_i \end{cases}$$

# 潜在変数の例（2）



深層学習によつ  
て自動的にでき  
ると期待する  
(それほど簡単にうま  
くいきませんが！！)

直交座標系

符号化  
(エンコード)

データ

$$\begin{Bmatrix} x_i \\ y_i \\ z_i \end{Bmatrix}$$

3次元

半径が定数の円柱表  
面にある束縛条件が  
利用されていない

円柱座標系

復号化  
(デコード)

$$\begin{Bmatrix} \theta_i \\ z_i \end{Bmatrix}$$

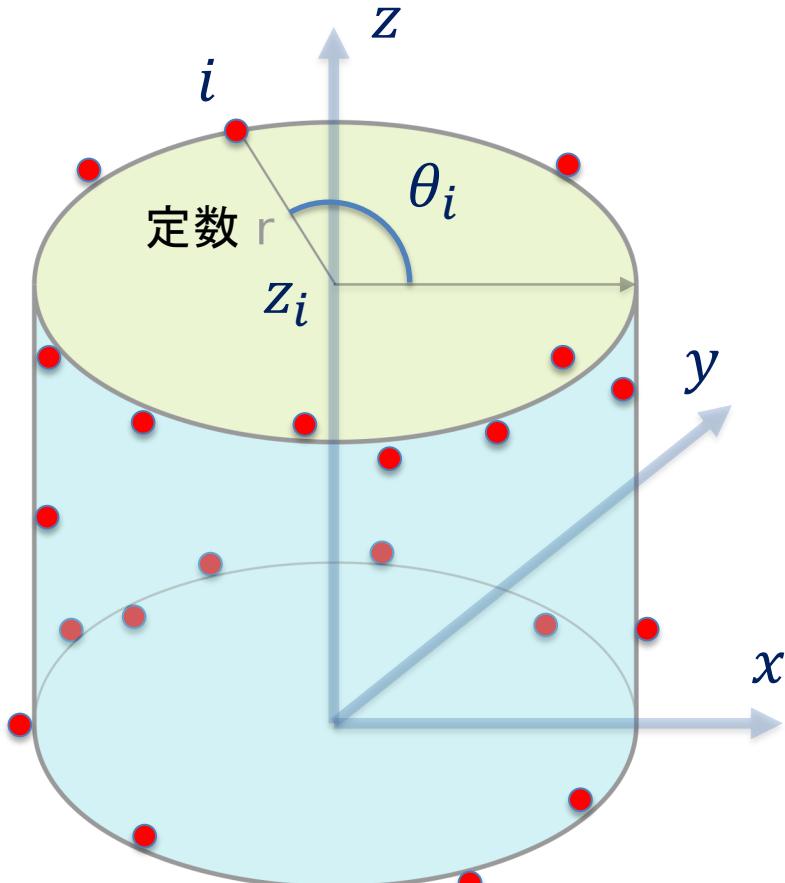
2次元

半径が定数の円柱表  
面にある束縛条件が  
利用する潜在変数

座標変換

$$\begin{Bmatrix} x_i = r \cos \theta_i \\ y_i = r \sin \theta_i \\ z_i = z_i \end{Bmatrix}$$

# 潜在変数の例 (3)



半径が定数の円柱の表面

直交座標系

円柱座標系

座標変換

データ

$$\begin{Bmatrix} x_i \\ y_i \\ Z_i \end{Bmatrix}$$

3次元

半径が定数の円柱表面にある束縛条件が利用されていない

$$\begin{Bmatrix} \theta_i \\ Z_i \end{Bmatrix}$$

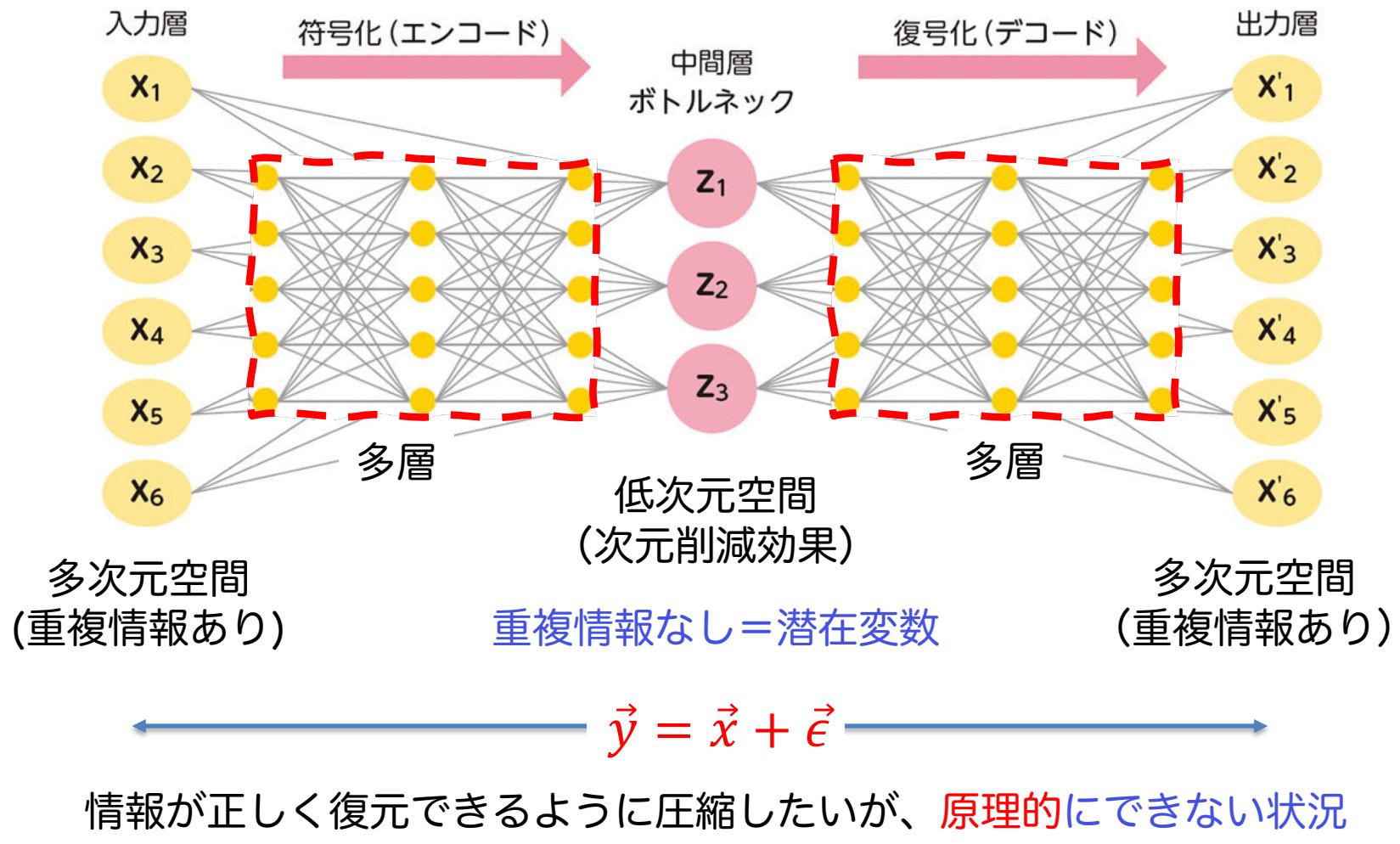
2次元

半径が定数の円柱表面にある束縛条件が利用する潜在変数

$$\begin{cases} x_i = r \cos \theta_i + \epsilon_x \\ y_i = r \sin \theta_i + \epsilon_y \\ z_i = z_i \end{cases}$$

# オートエンコーダ (3)

## ■ オートエンコーダの構造

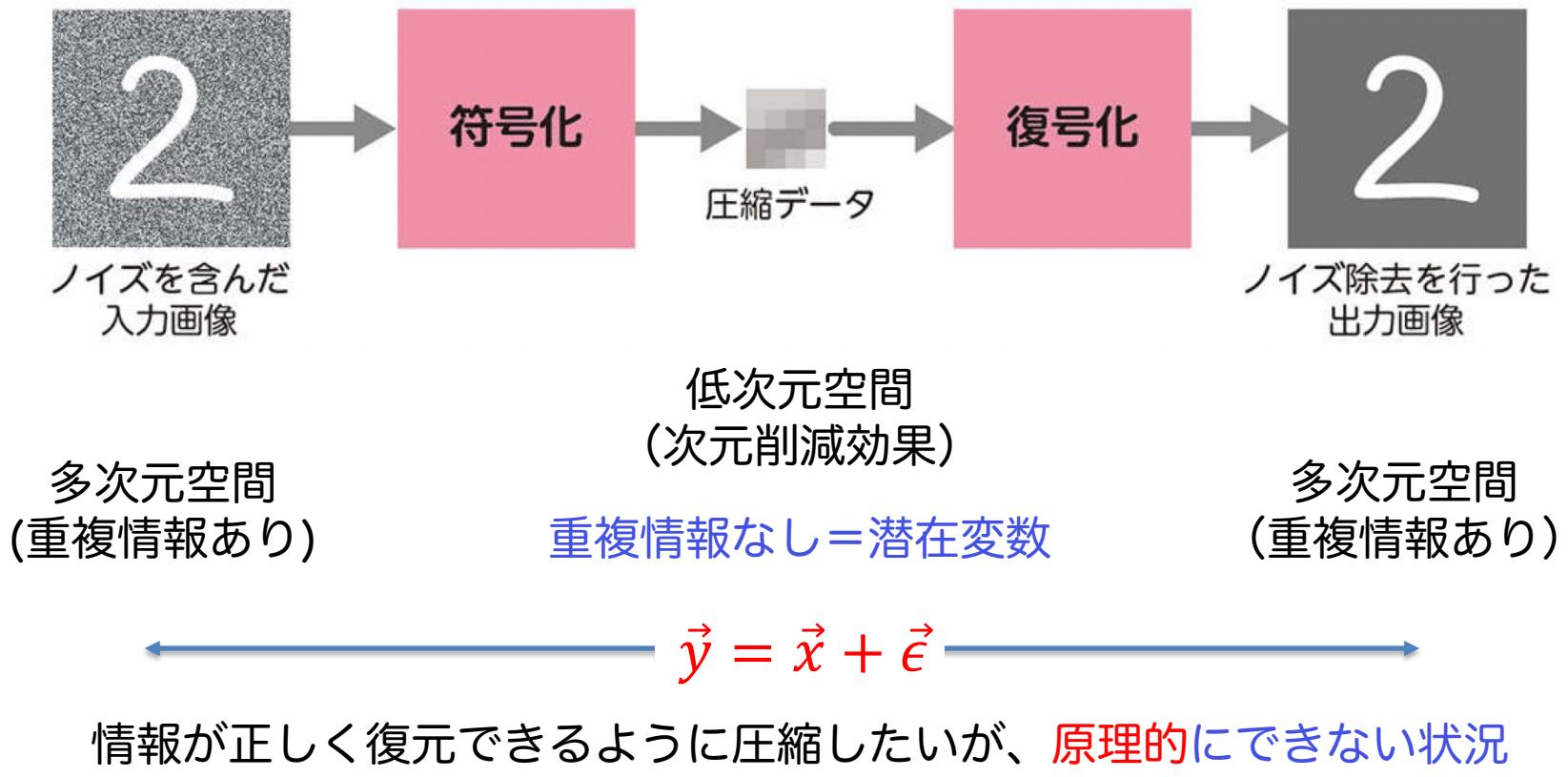


# オートエンコーダの例

Deep Clustering with Convolutional Autoencoders

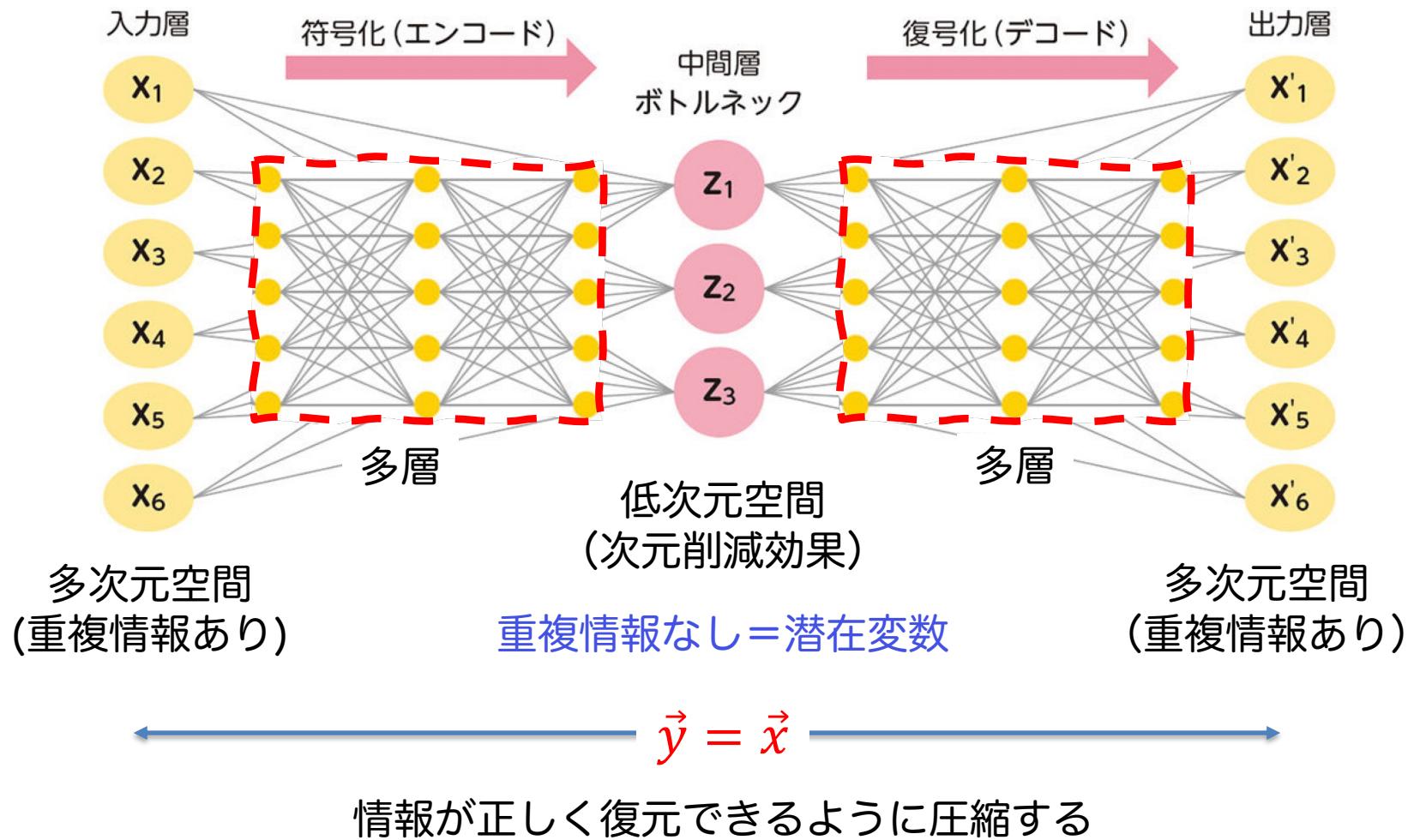
Xifeng Guo, Xinwang Liu, En Zhu, and Jianping Yin

Neural Information Processing. ICONIP 2017. Lecture Notes in Computer Science, vol 10635. Springer.



# オートエンコーダ (2)

## ■ オートエンコーダの構造

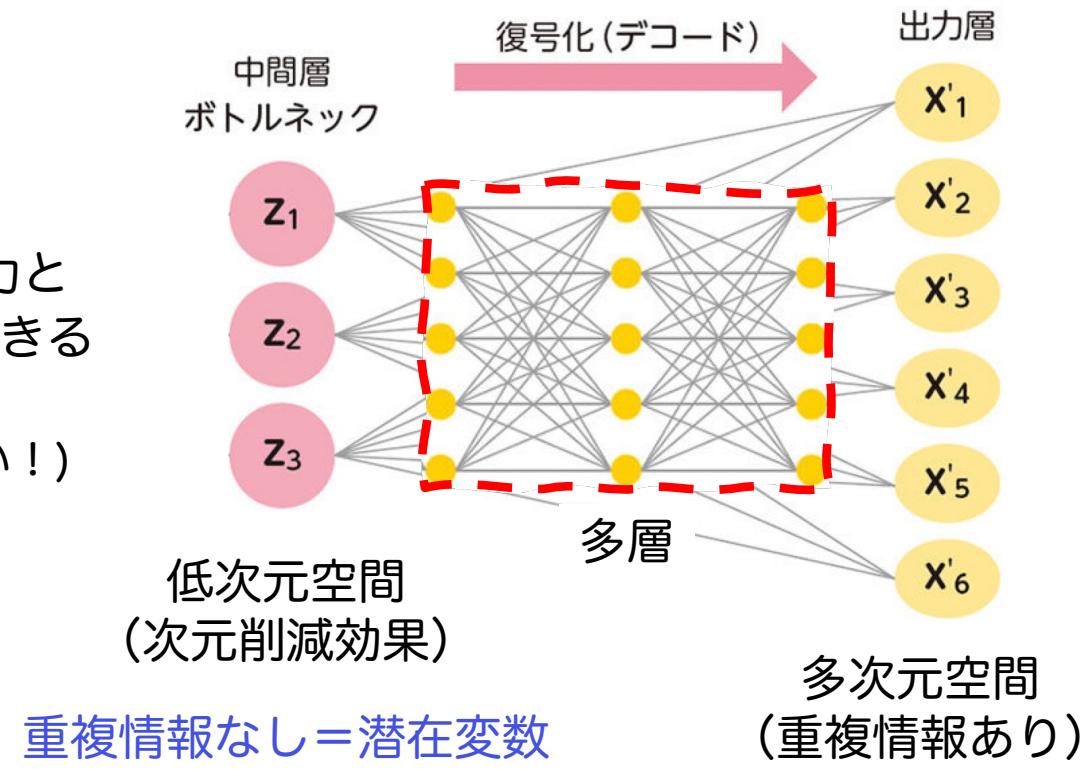


# オートエンコーダによるデータ生成

符号化情報を入力すれば入力と同じ性質を持つ出力を生成できる  
**データ生成モデル**  
(期待するが、そうとは限らない！)



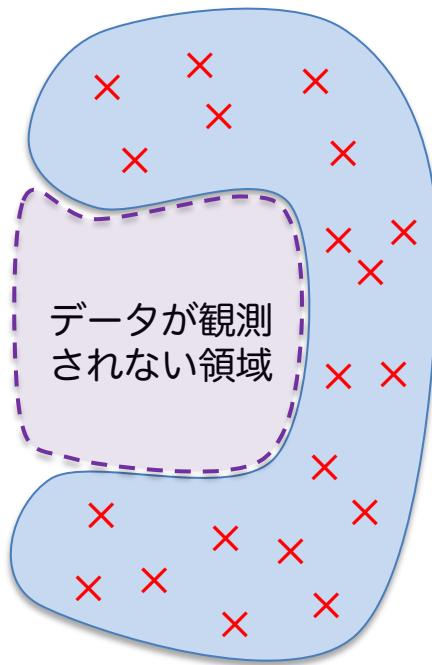
出典：  
東京大学  
松尾研究ホームページ



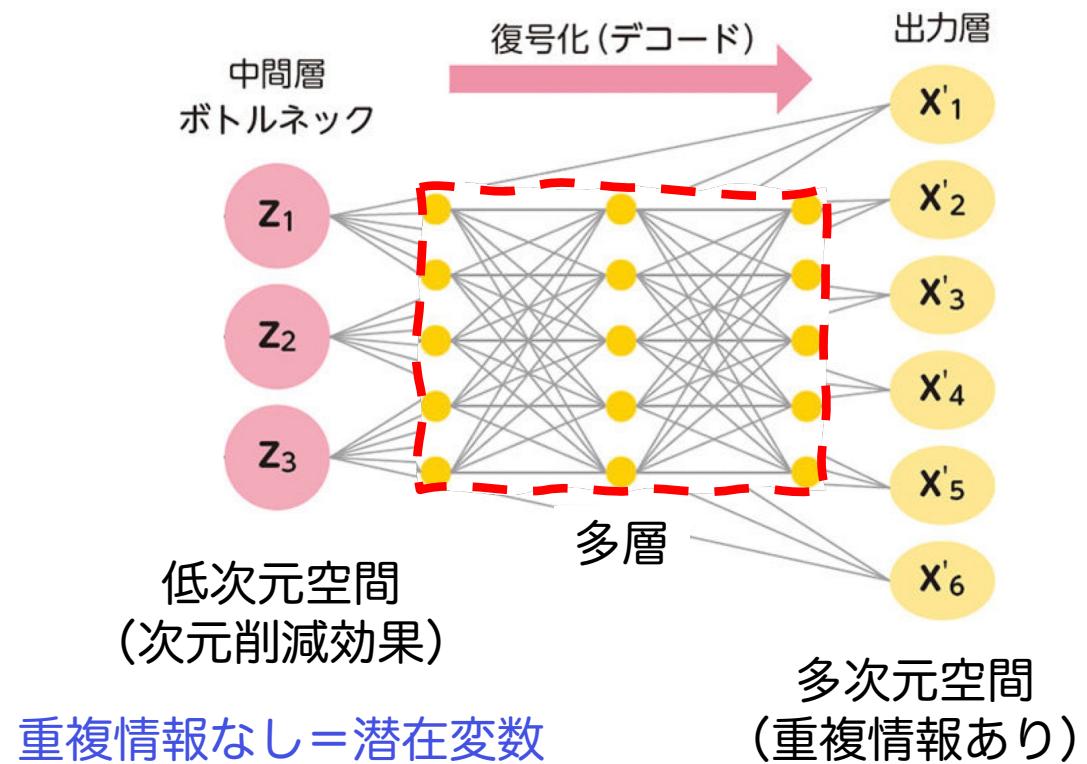
フェイクデータなどを生成できる

# オートエンコーダによるデータ生成

潜在変数  $\vec{z}$  空間

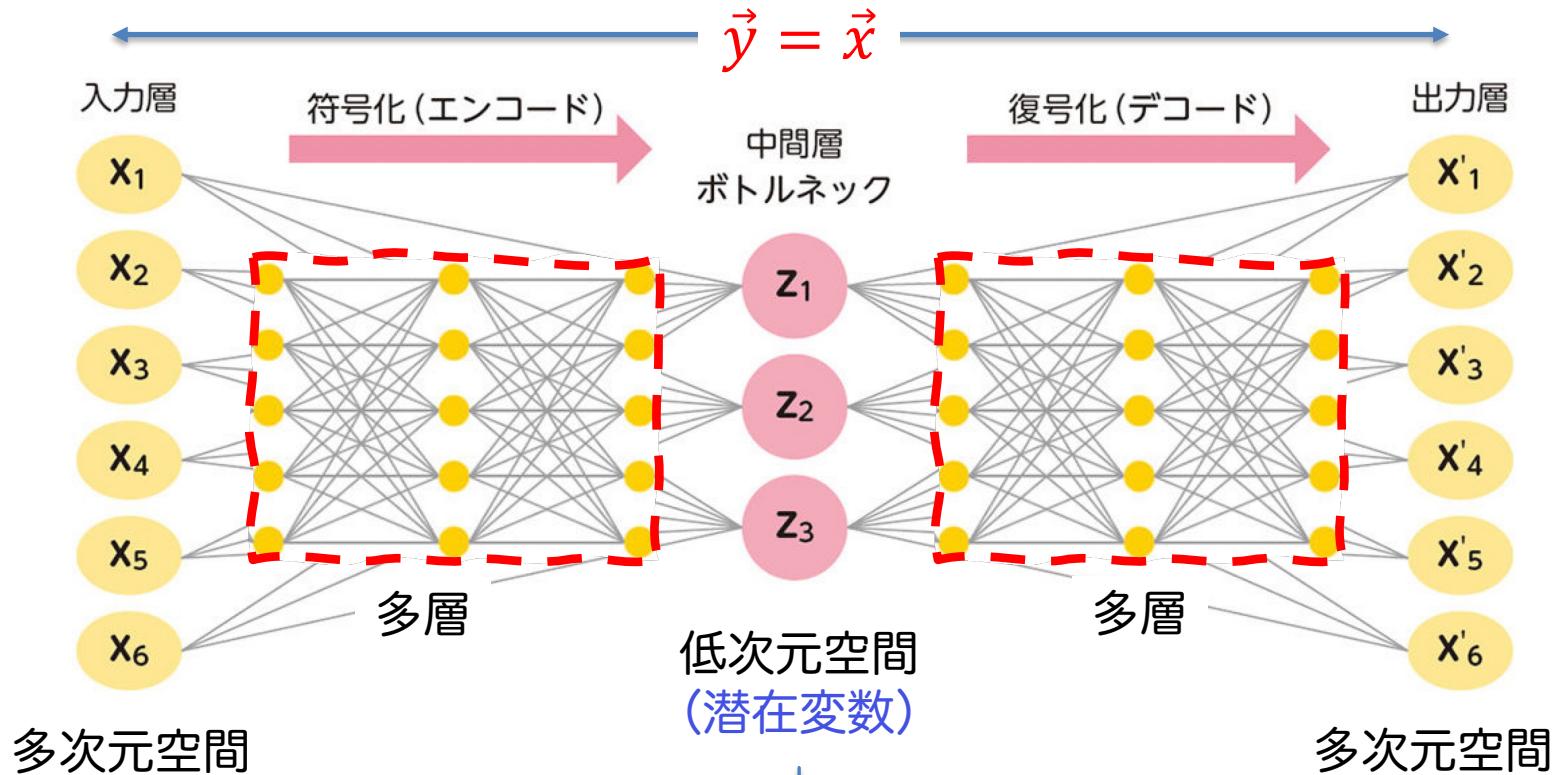


✗ 訓練データ点

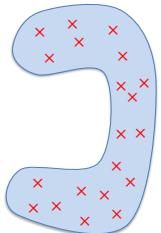


外挿問題：ありえないデータの生成？

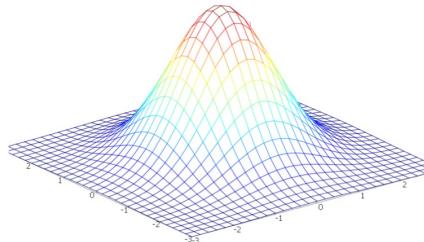
# 変分オートエンコーダ (VAE)



潜在変数空間  
分布の形を強制  
して学習する



強制



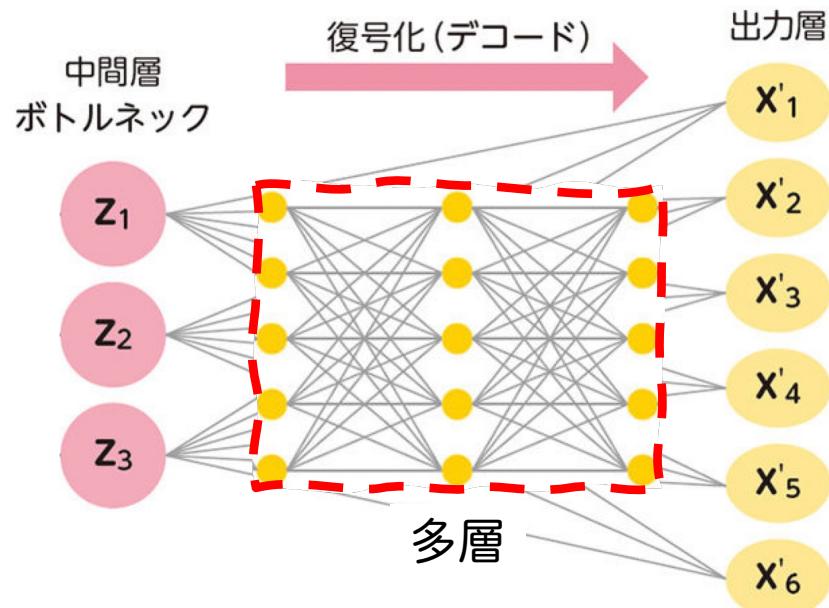
外挿問題：  
ありえないデータの  
生成をさせたい

# 敵対的生成ネットワークGAN

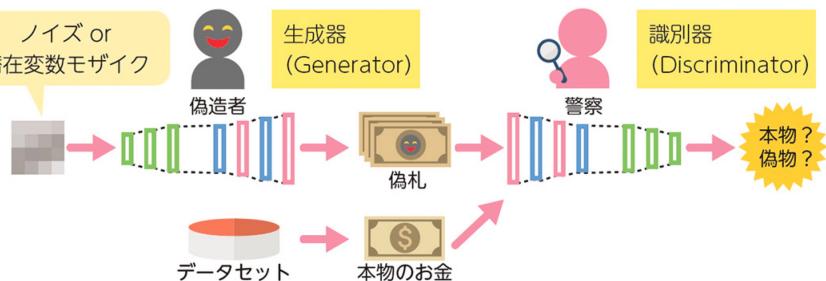
符号化情報を入力すれば入力と同じ性質を持つ出力を生成できる  
データ生成モデル



実データとフェイクデータを識別するモデルを同時に訓練する



■ GANは偽造者と警察？



# データサイエンスの基礎

---

北陸先端科学技術大学院大学  
Hieu-Chi Dam

2024年7月8日（第2部 19:30 ~ 20:45）

# 目次

---

1. イントロ
  - i. データサイエンスの紹介
  - ii. 推論方式の紹介
  - iii. 数理統計と機械学習の紹介
2. データサイエンスの材料科学への応用事例紹介



# Data analytics process

## 「5 iterative steps」

# Knowledge discovery and Data mining

可視化

The automatic extraction of non-obvious, hidden knowledge from large volumes/complicated data



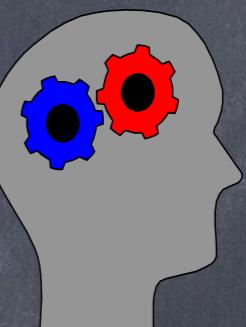
① Data collection and representation

② Learning/mining from data  
machine learning and data mining algorithms

③ Knowledge representation and evaluation

可視化

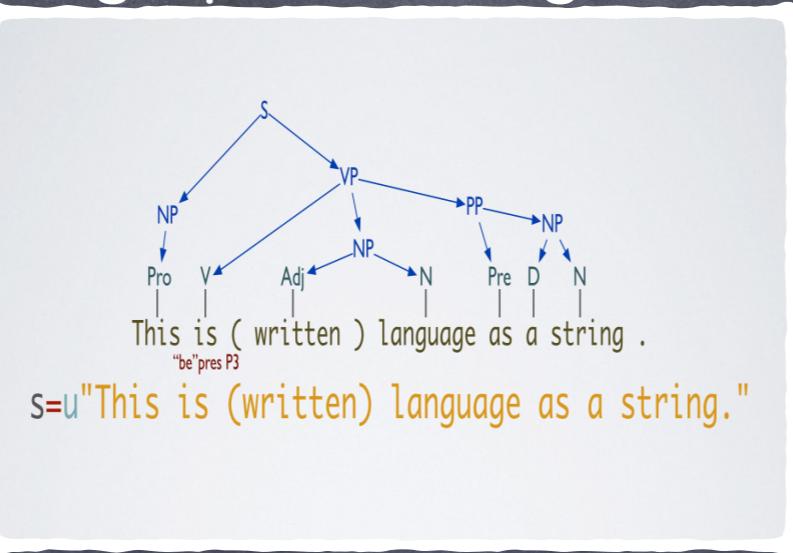
潜む知見  
人間に理解可能



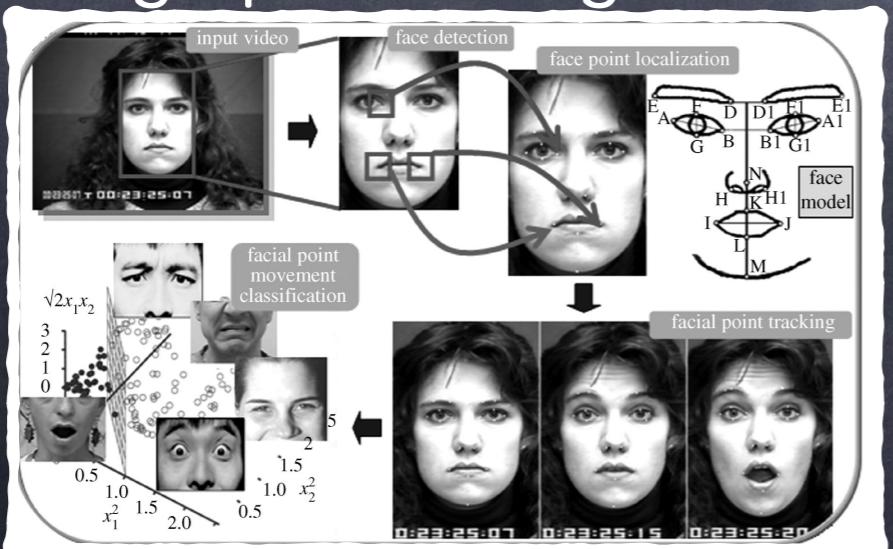
# ① Data

# collection and representation

# Text and Natural language processing



# Image processing



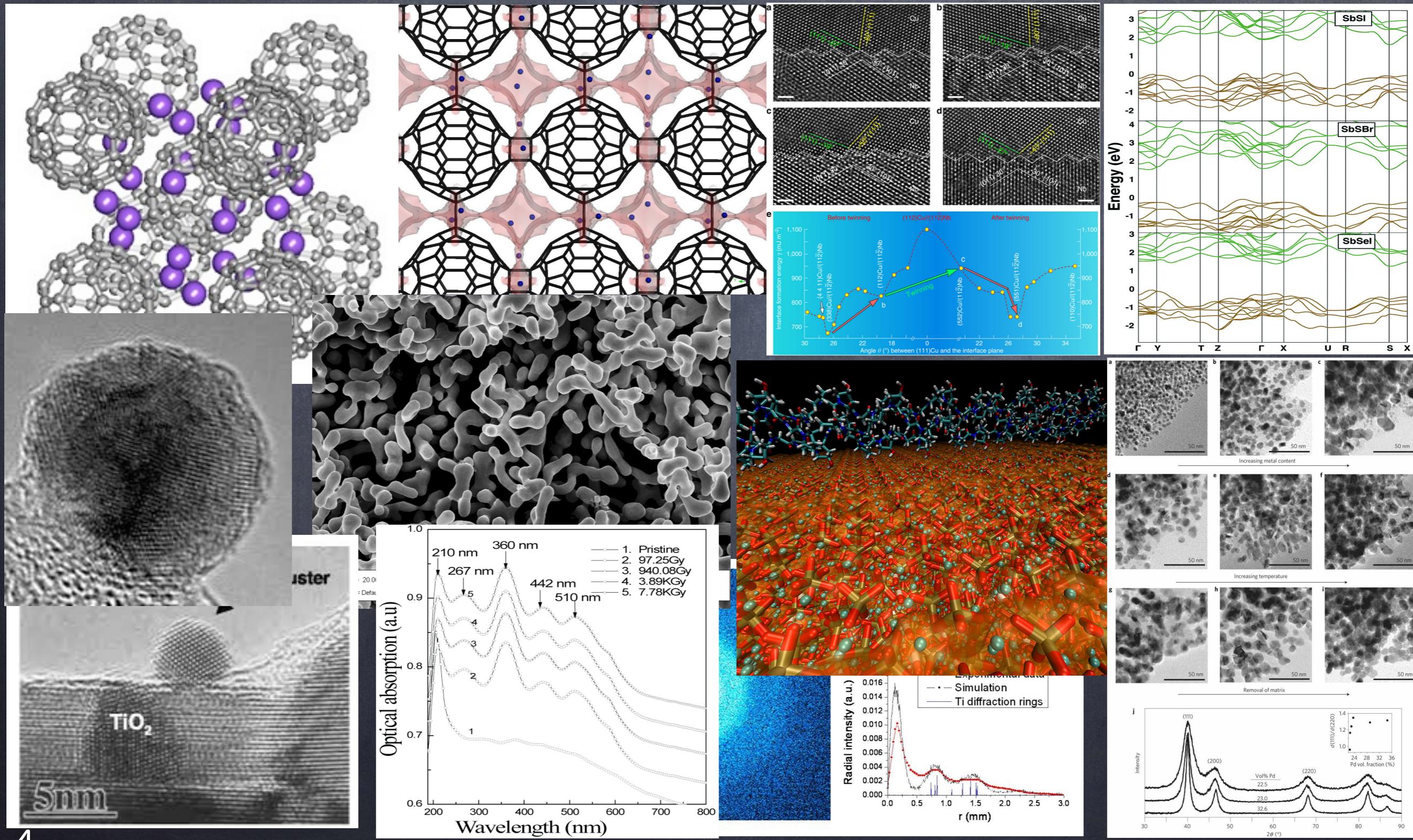
# Problem setting

# Data instance

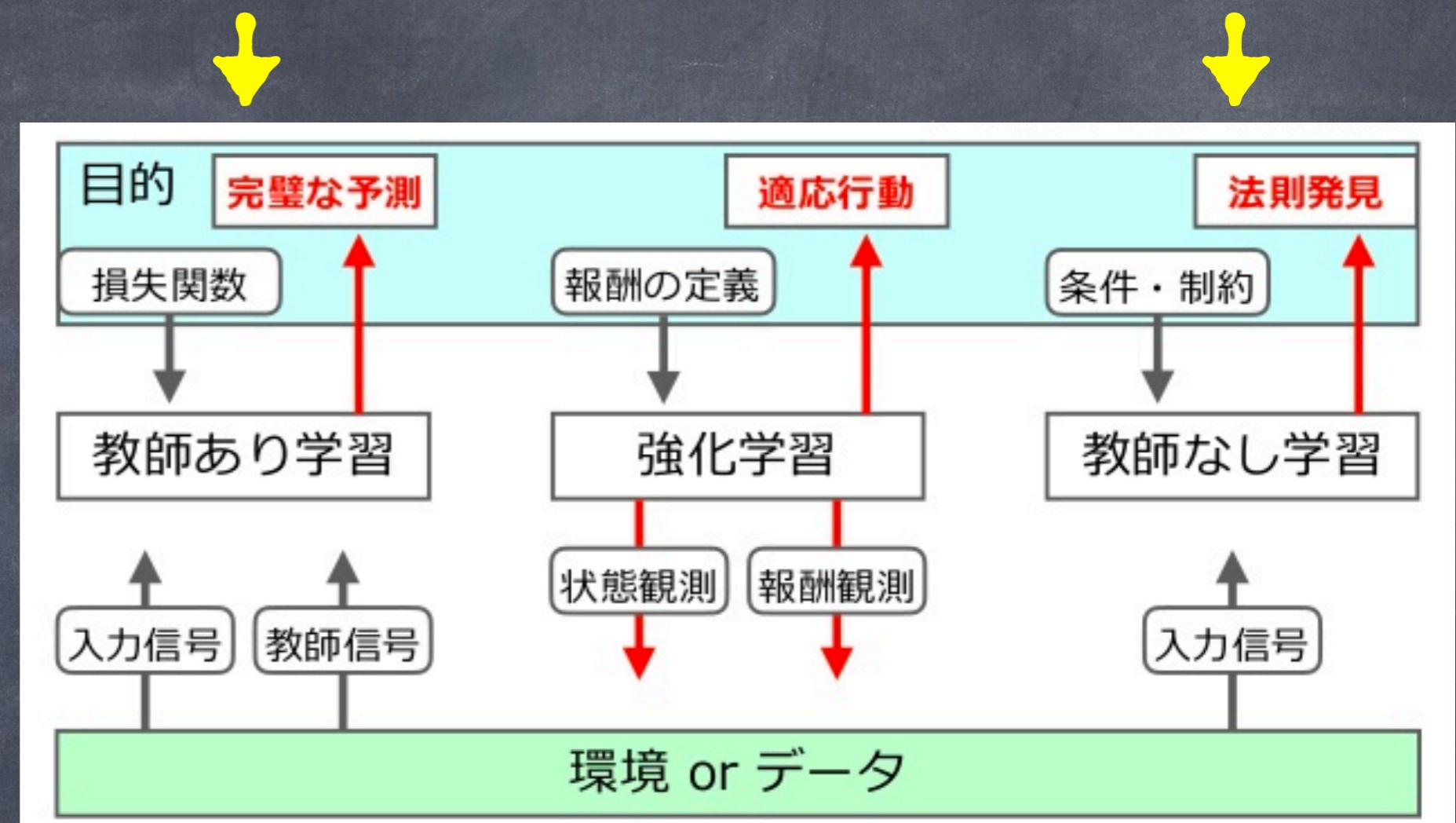
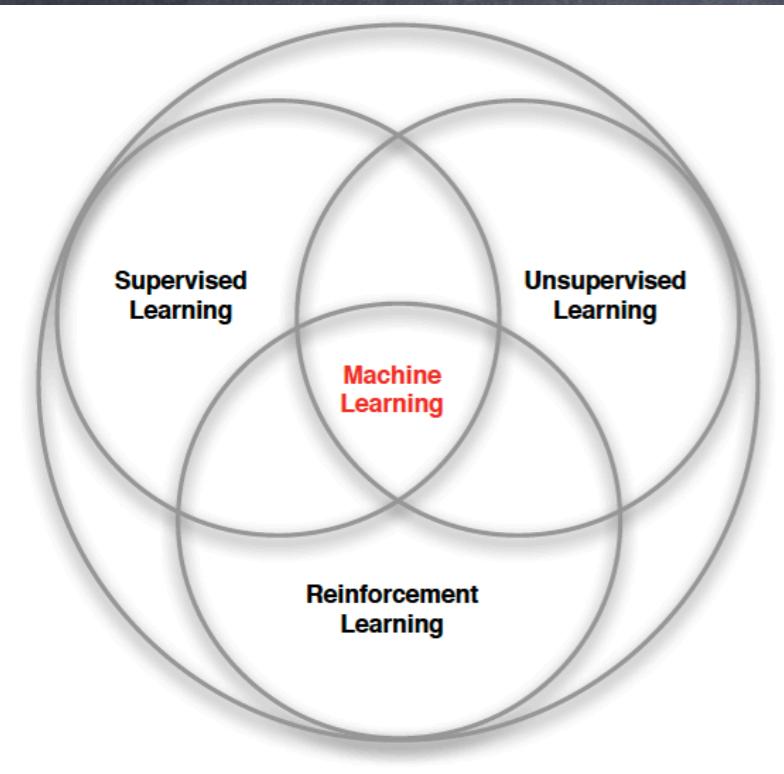
Raw data : characters  
character?  
word?  
sentence?  
paragraph?  
document?

Raw data : pixels  
pixel?  
pixel window?  
gradient?  
image?

# Data in Materials science



## Branches of Machine learning



How to “translate” a materials science problem  
into machine learning/data mining language?

Representation, similarity measure, and learning



## 磁石材料のキュリ温度データの事例

## 磁力の強い磁石ランキング

1位 ネオジム

2位 サマリウムコバルト

3位 アルニコ

4位 フェライト

## 熱に強い磁石ランキング

1位 アルニコ (500°C)

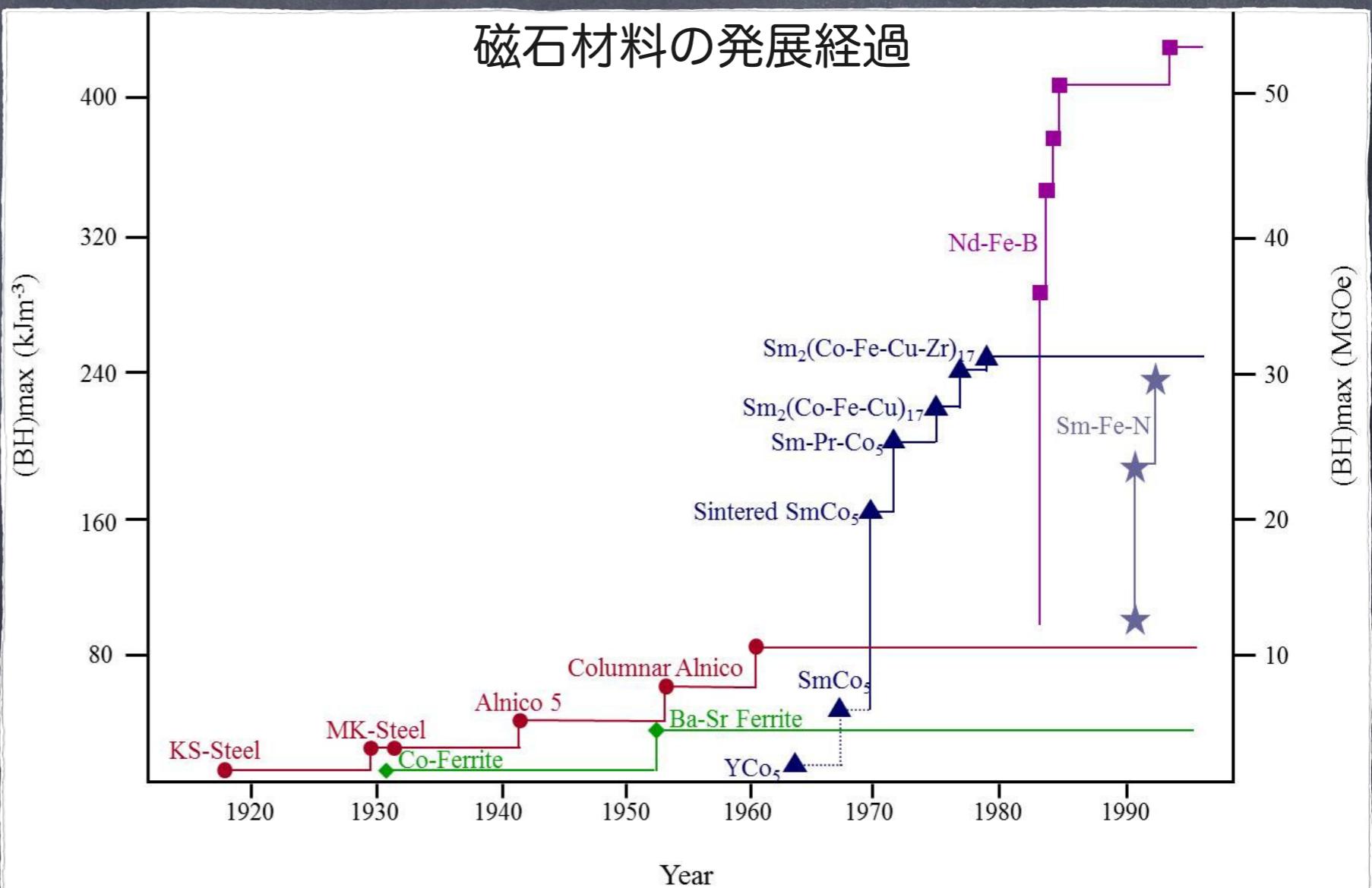
2位 サマリウム・コバルト (300°C)

3位 フェライト (200°C)

4位 ネオジム (80°C)



## 磁石材料の発展経過



1	H
3	T+
4	B+
11	Mg

## The Periodic Table of the Chemical Elements

### Transition metals

19	20	21	22	23	24	25	26	27	28	29	30
K	Cx	S+	T+	V	CR	Mn	T+	C+	M+	Cy	Zn
37	38	39	40	41	42	43	44	45	46	47	48
Rb	Fr	Y	Zr	Nb	Mn	Tc	Ru	Rh	Pd	Wg	Cd
55	56	57	58	59	60	61	62	63	64	65	66
Cs	B+	T+	H+	T+	U	Re	O+	H+	Pt	Wg	Hg
87	88										
Fr	R+	T+	R+	D+	S+	B+	H+	M+	D+	P+	

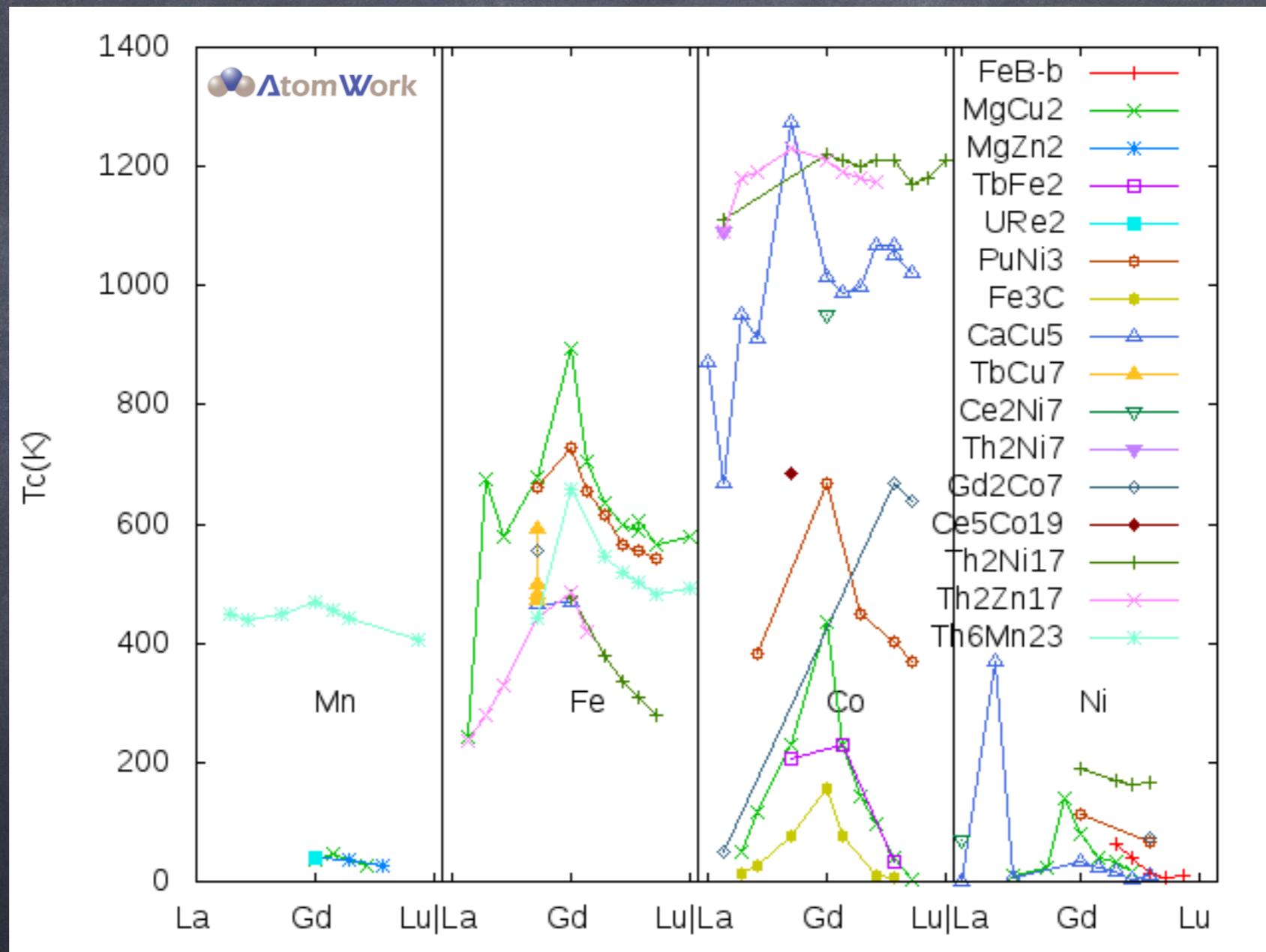
5	6	7	8	9	10
B	C	N	O	F	N+
X+	S+	P	S	C+	XR
G+	G+	X+	S+	B+	K+
49	50	51	52	53	54
Hg	S+	S+	T+	T+	X+
81	82	83	84	85	86
T+	Pb	B+	Po	X+	R+

### Rare-earth metals

57	58	59	60	61	62	63	64	65	66	67	68	69	70	71
T+	C+	Pr	M+	Pm	S+	E+	G+	T+	D+Y	H+	C+	ER	T+	U
X+	T+	P+	Y	M+	P+	Wg	M+	C+	B+	C+	C+	F+	M+	LR

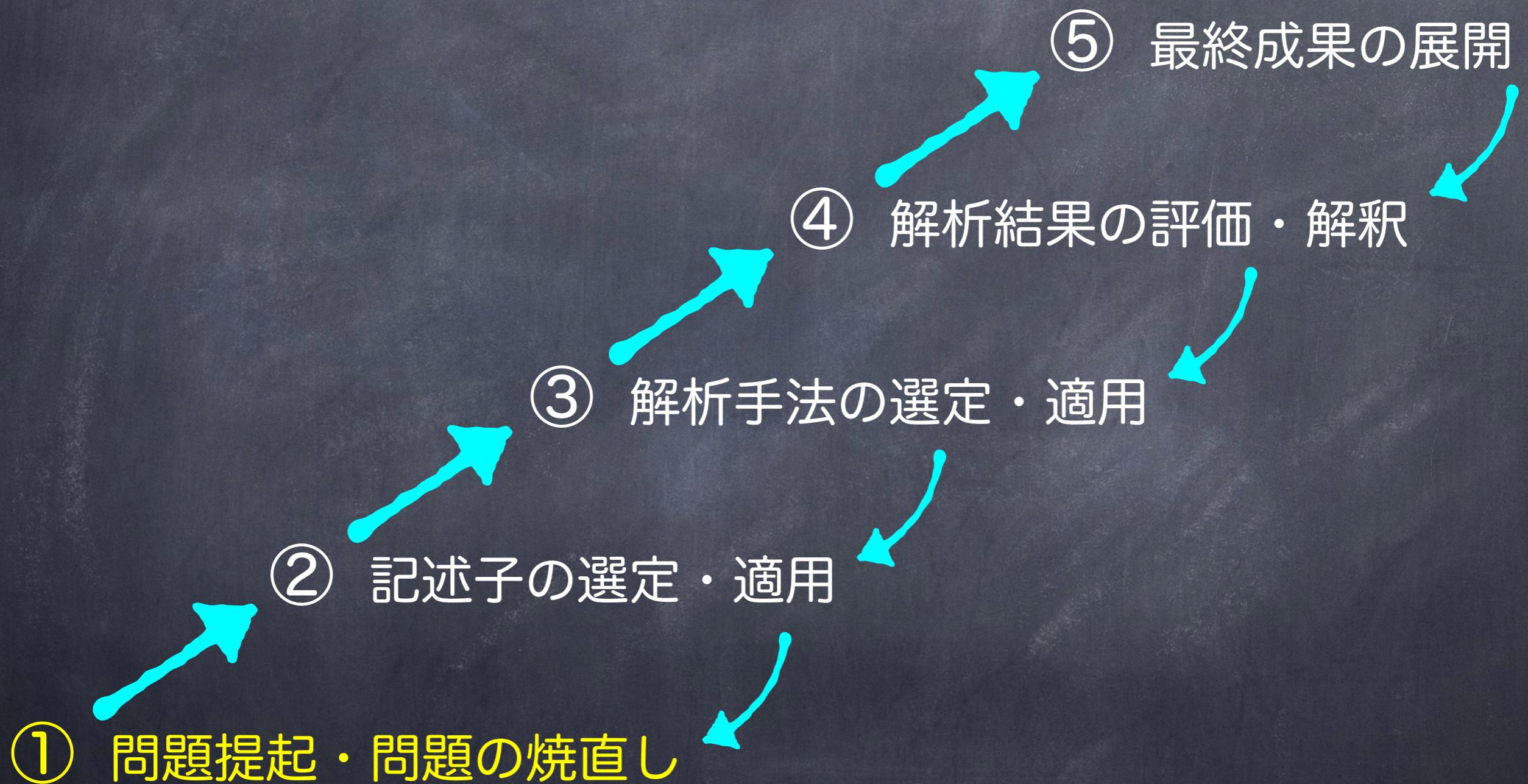
# Curie temperature of 3d-4f binary alloys

d-f bimetal alloys from literatures

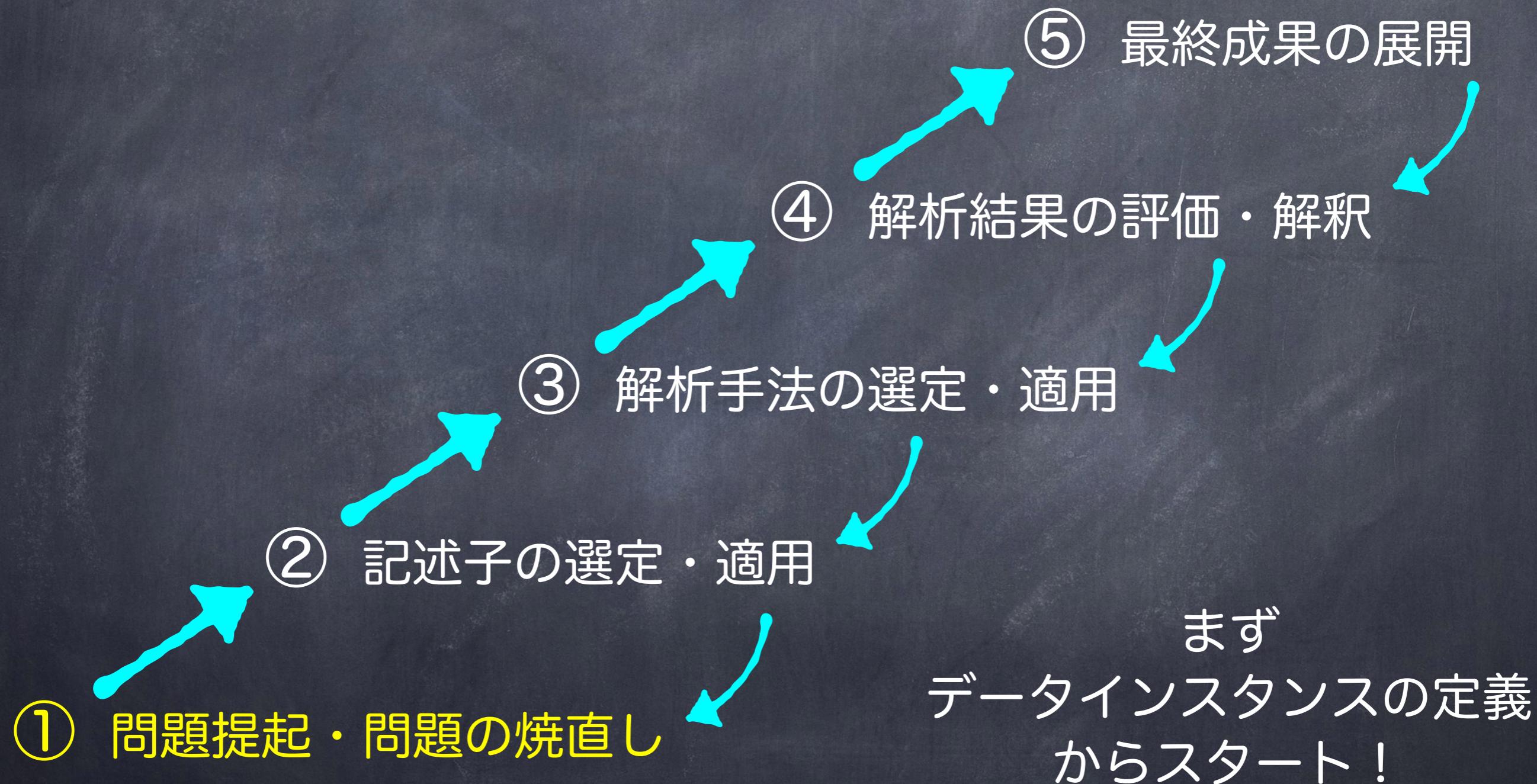


生データを集めました！

# MIの反復5ステップ

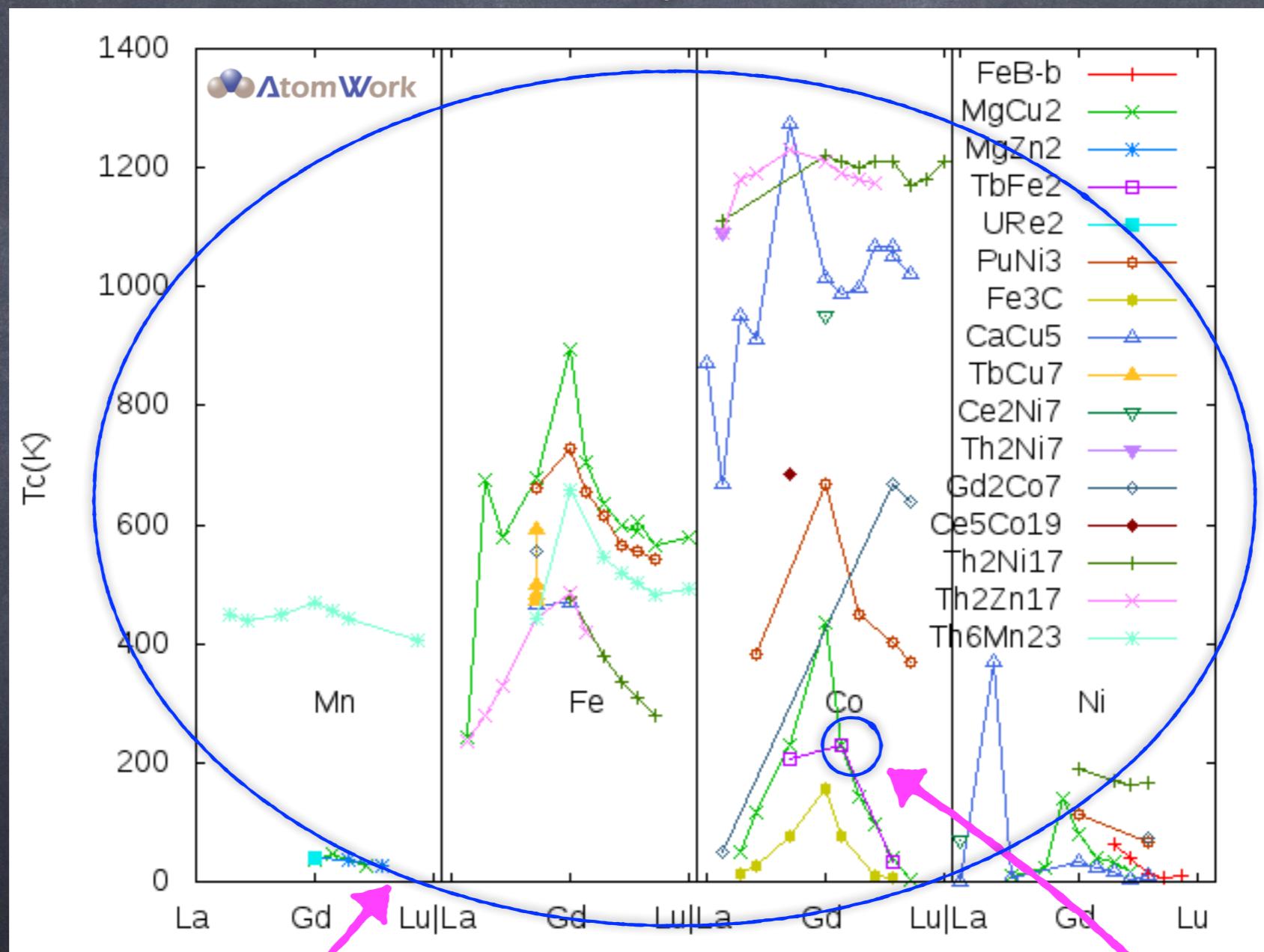


# MIの反復5ステップ



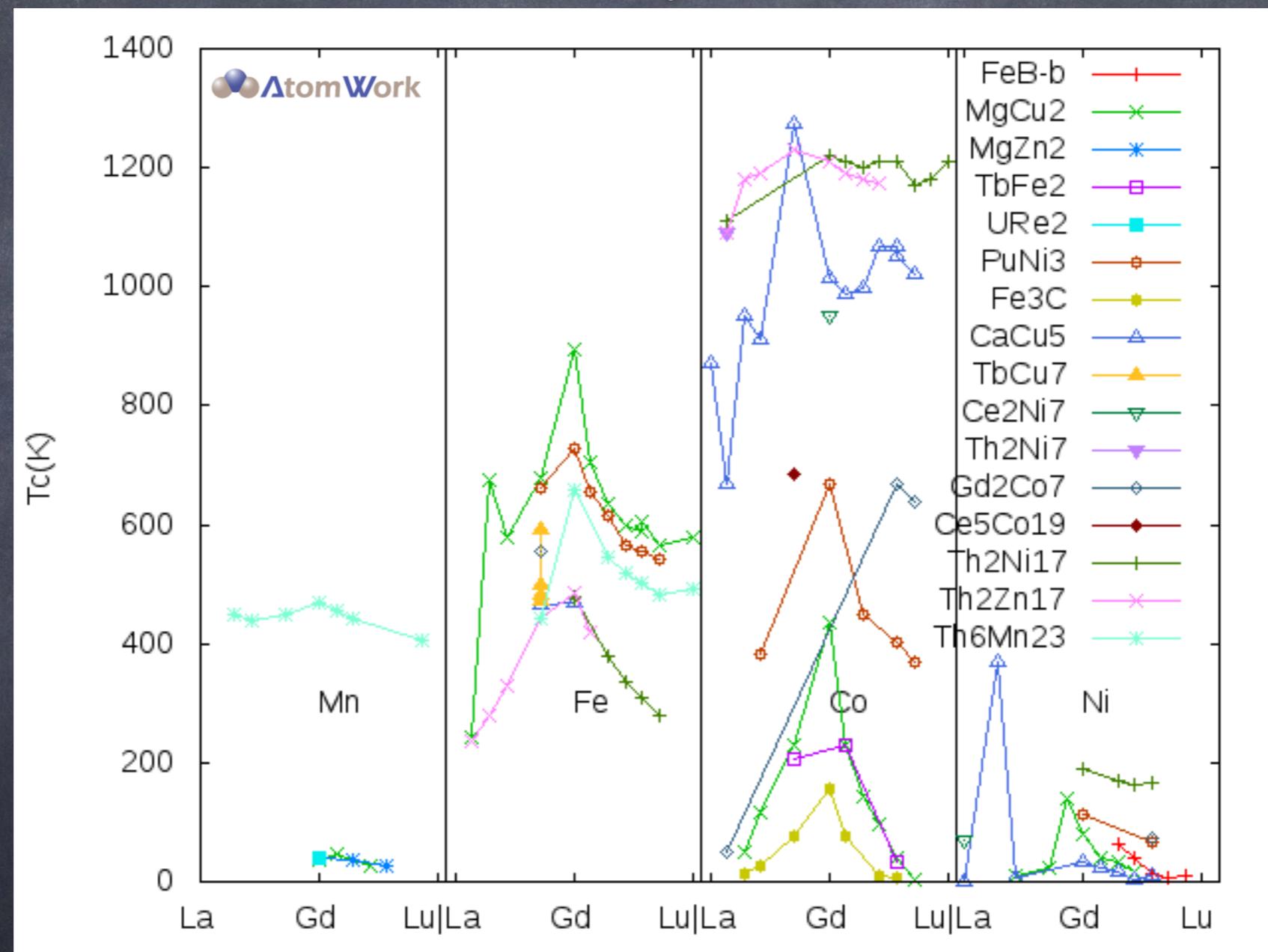
# Curie temperature of 3d-4f binary alloys

d-f bimetal alloys from literatures



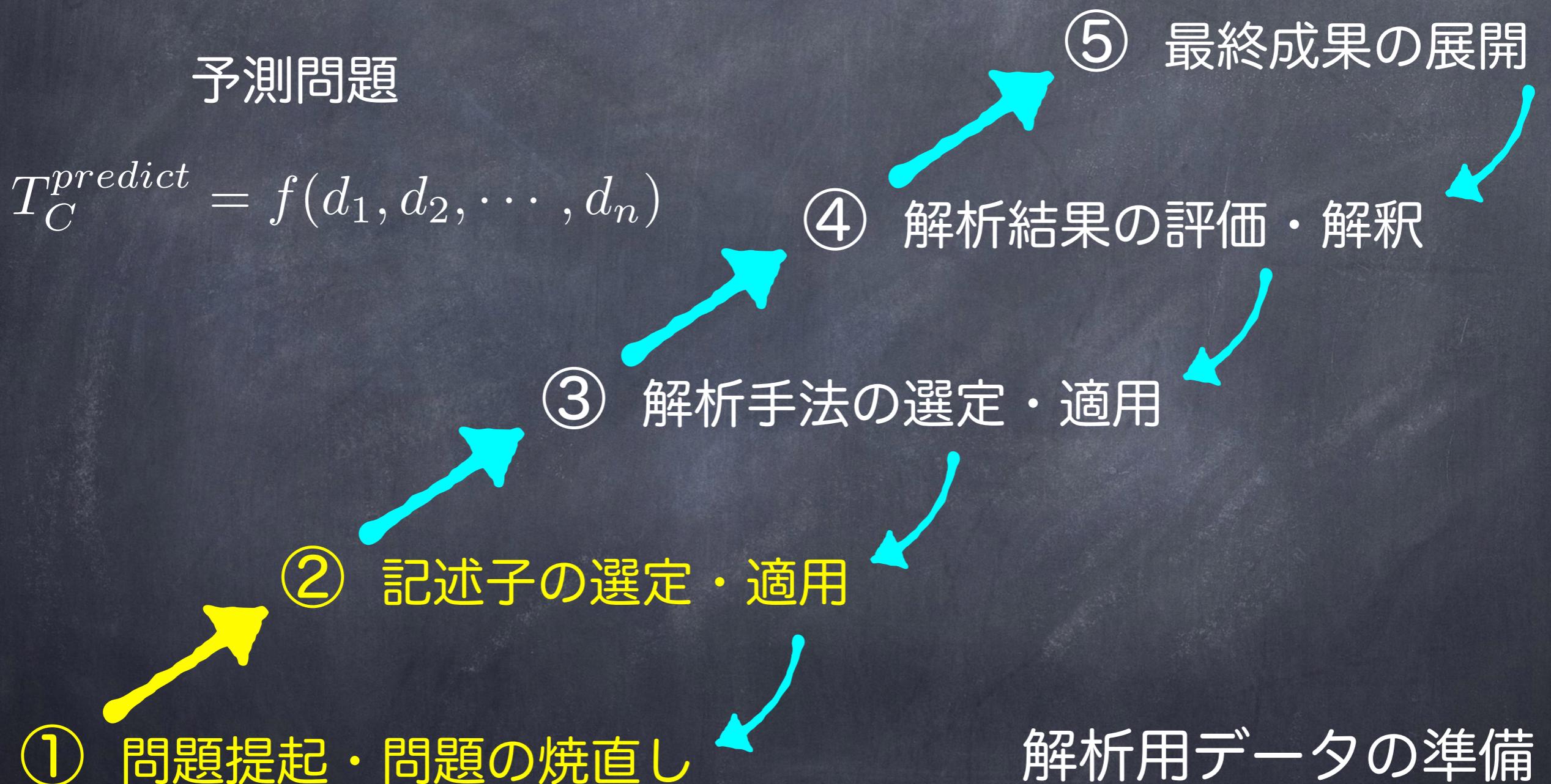
# 問題提起：キュリ温度を予測したい

d-f bimetal alloys from literatures



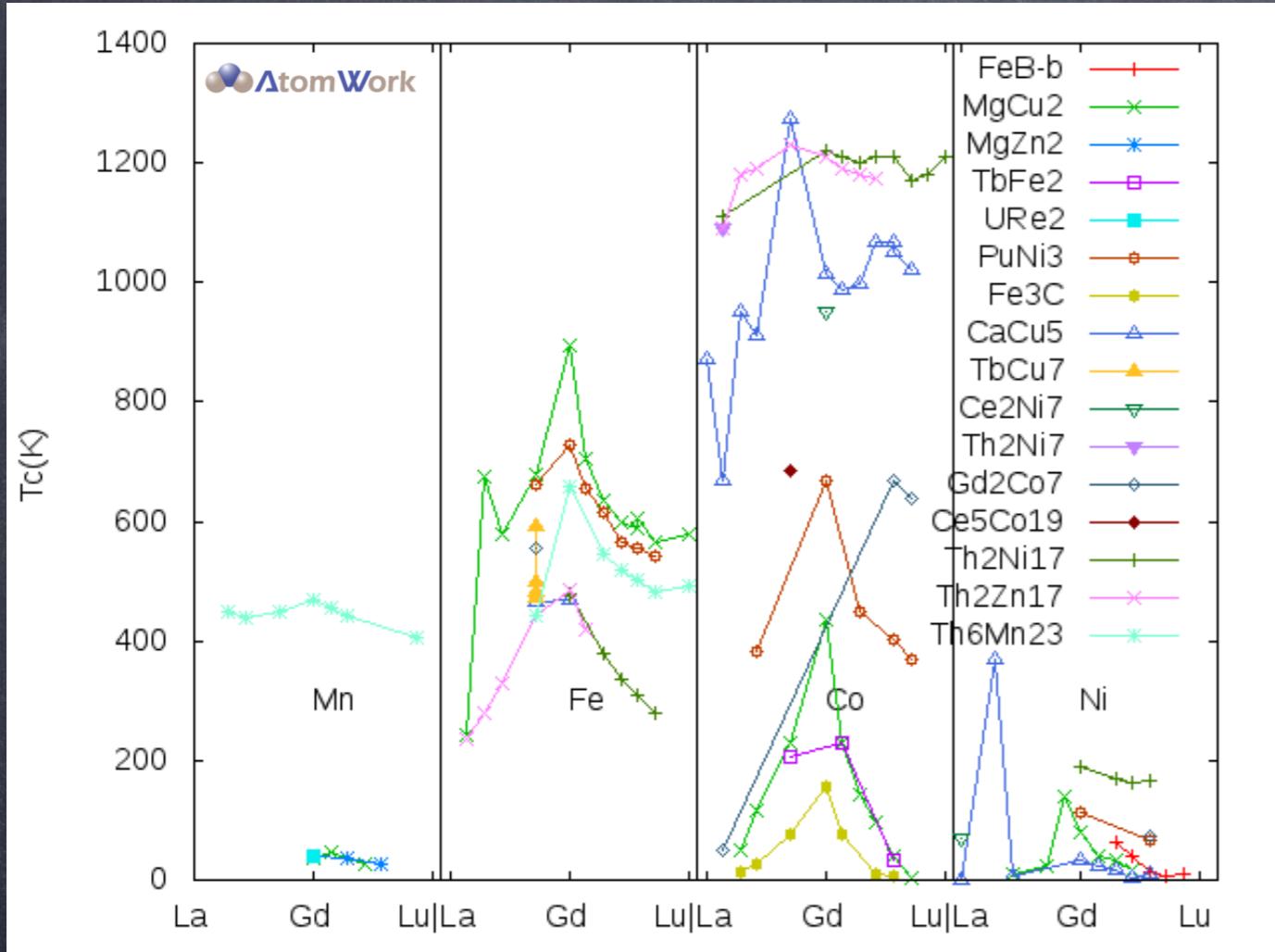
$$T_C^{predict} = f(\square, \diamond, \triangle, \dots)$$

# MIの反復5ステップ



# 2元合金情報の記述

d-f bimetal alloys from literatures



何が特徴なのか？

Material structure

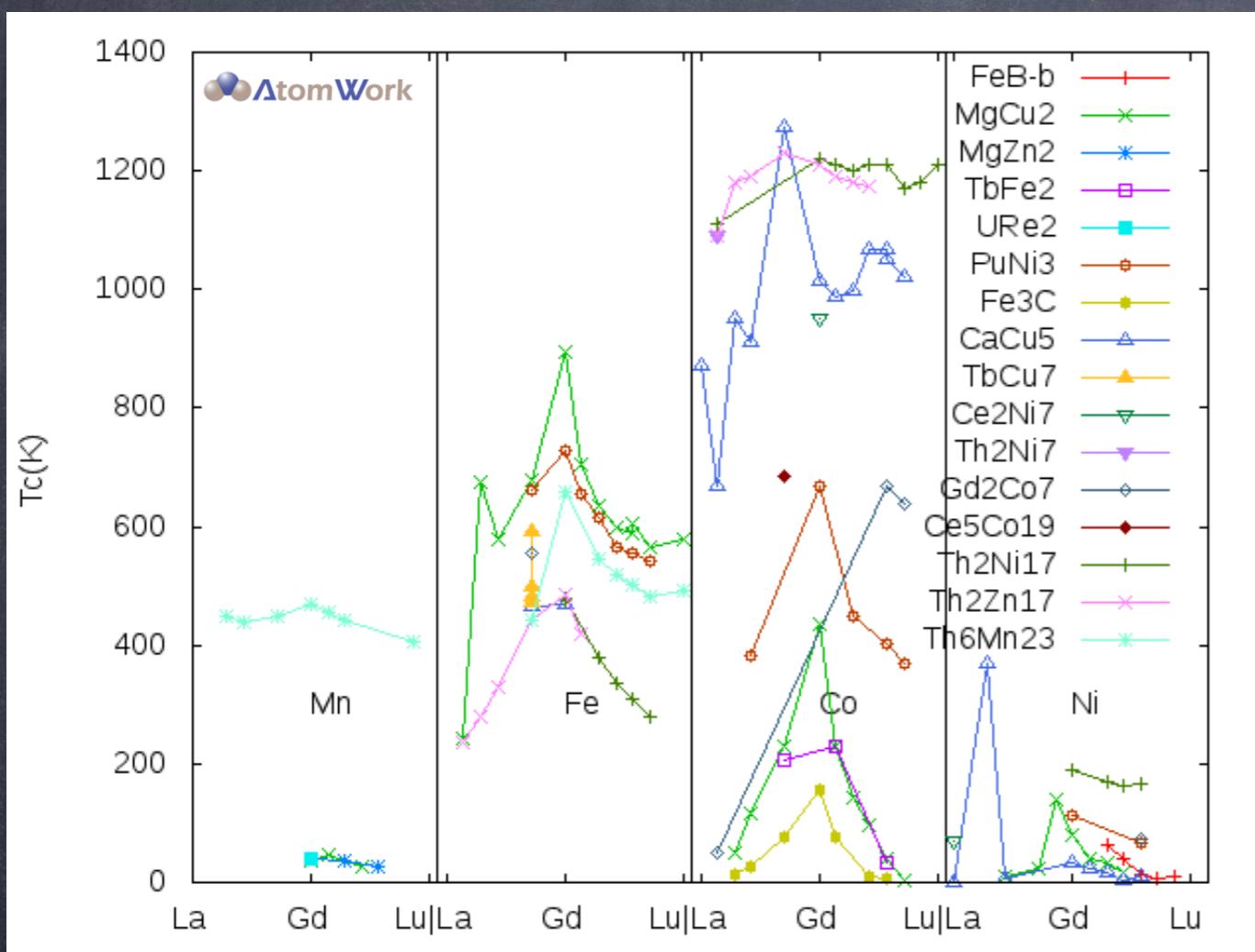
$$\left. \begin{array}{l} a, b, c, \alpha, \beta, \gamma \\ atom_1, x_1, y_1, z_1 \\ atom_2, x_2, y_2, z_2 \\ \vdots \\ atom_i, x_i, y_i, z_i \\ \vdots \\ atom_n, x_n, y_n, z_n \end{array} \right\}$$

記述子

$$T_C^{predict} = f(\square, \diamond, \triangle, \dots)$$

# 2元合金データの記述子の選定

d-f bimetal alloys from literatures



まず、予測できるかをテストする

$$T_C^{predict} = f(d_1, d_2, \dots, d_n)$$

想定する記述子

Structural information:

- R-R Nearest neighbor distance
- R-T Nearest neighbor distance
- T-T Nearest neighbor distance
- Number of R surrounding R
- Number of T surrounding T
- Number of R surrounding T
- Number of T surrounding R
- Concentration of R
- Concentration of T

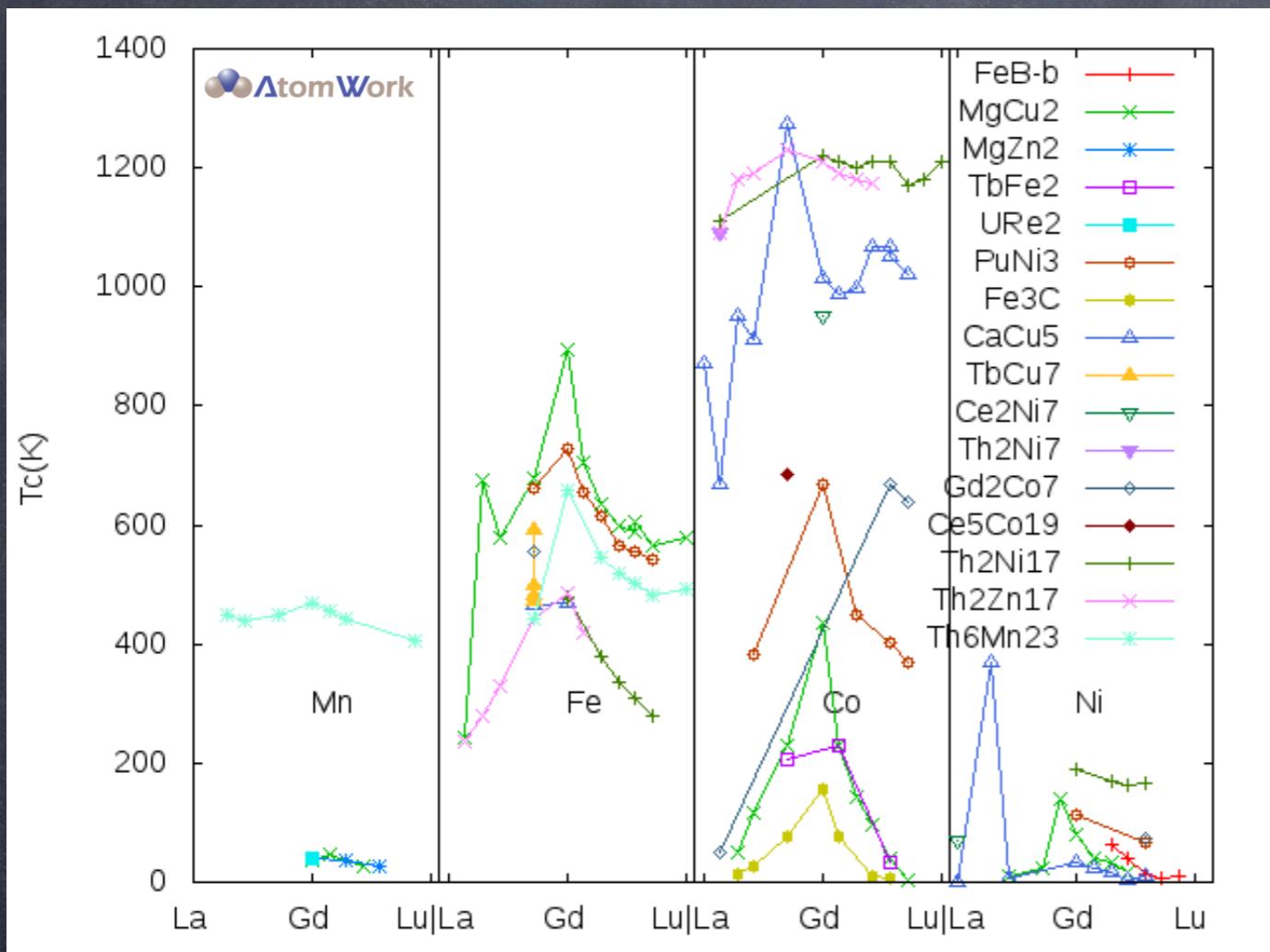
Atomic information:

- Atomic Number of  $Z_R, Z_T$
- Covalent Radius of  $r_R, r_T$
- Electron Negativity of  $\chi_R, \chi_T$
- Ionization Potential  $IP_R, IP_T$
- Orbital angular momentum  $L_R, L_T$
- Spin angular momentum  $S_R, S_T$
- Total angular momentum  $J_R, J_T$
- Landé g-factor  $g_{JR}, J_R g_{JR}, J_R(1-g_{JR})$

Empirical descriptors

# 2元合金データの記述子の適用

d-f bimetal alloys from literatures



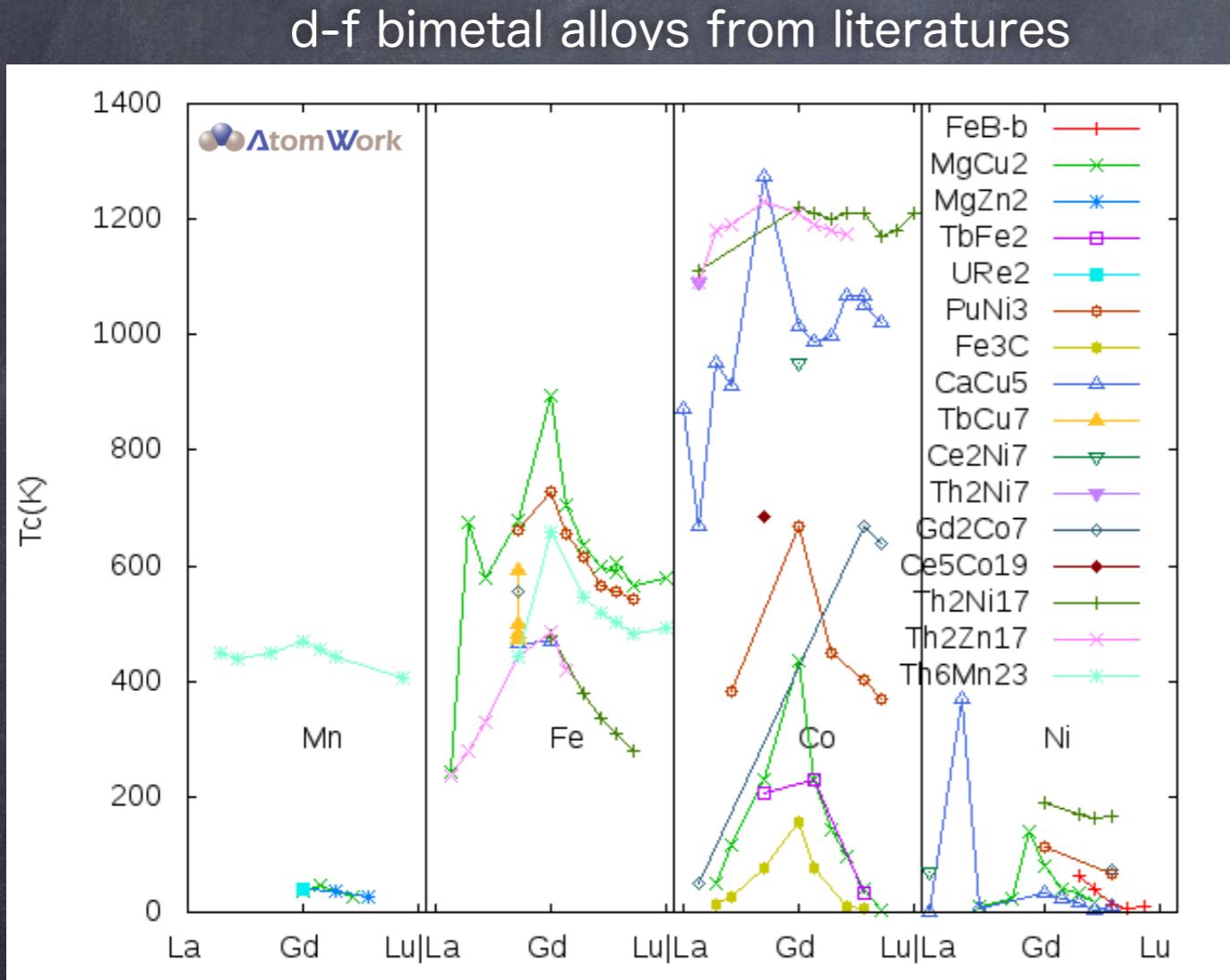
解析用データ

<i>ID</i>	$d_1$	$d_2$	...	$d_n$	$T_C$
合金 1	$d_1^1$	$d_2^1$	...	$d_n^1$	$T_C^1$
合金 2	$d_1^2$	$d_2^2$	...	$d_n^2$	$T_C^2$
...	...	...	...	...	...
合金 $n$	$d_1^n$	$d_2^n$	...	$d_n^n$	$T_C^n$

まず、予測できるかをテストする

$$T_C^{predict} = f(d_1, d_2, \dots, d_n)$$

# 2元合金データの記述子の適用



まず、予測できるかをテストする

$$T_C^{predict} = f(d_1, d_2, \dots, d_n)$$

# 解析用データ

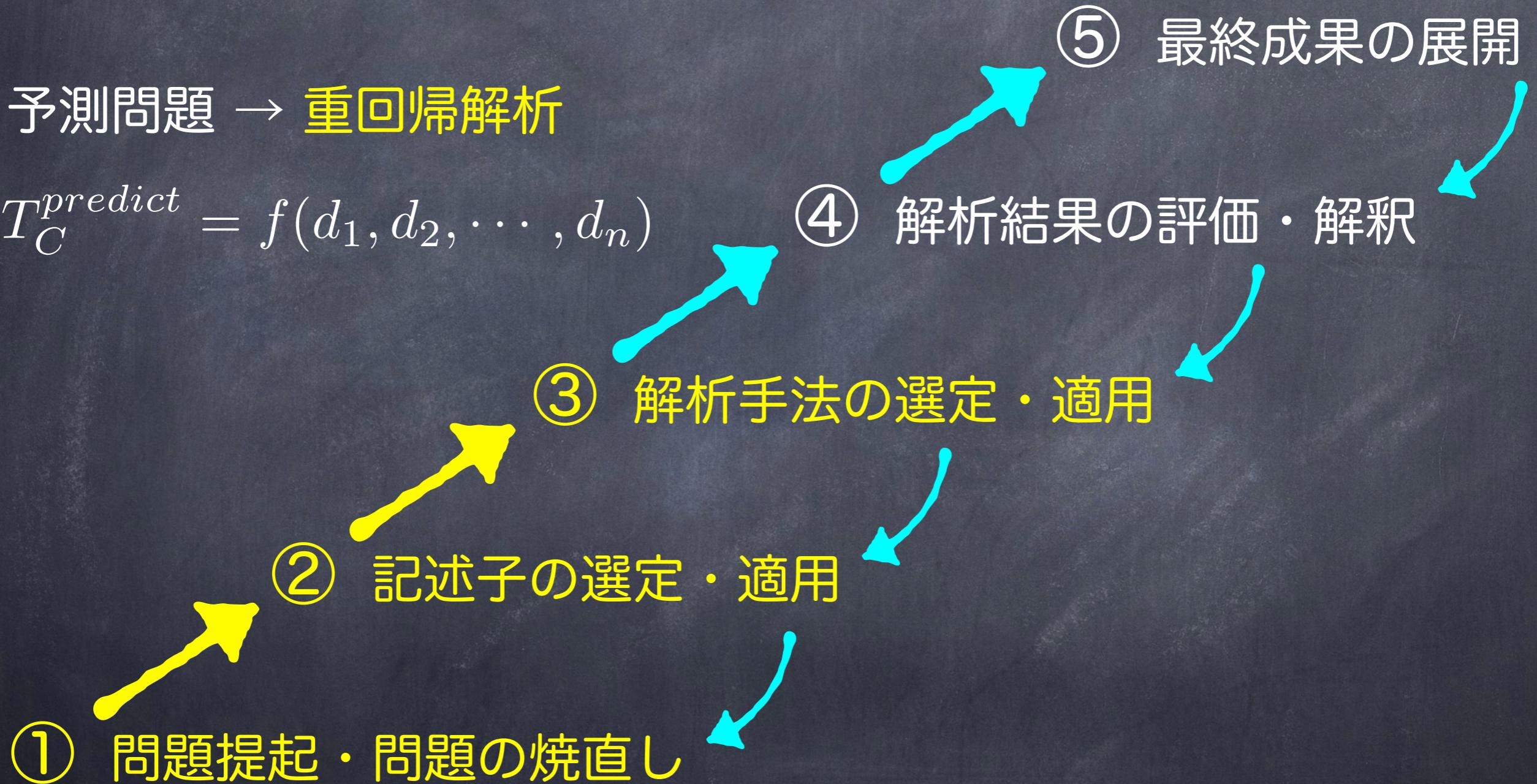
$ID$	$d_1$	$d_2$	$\dots$	$d_n$	$T_C$
合金 1	$d_1^1$	$d_2^1$	$\dots$	$d_n^1$	$T_C^1$
合金 2	$d_1^2$	$d_2^2$	$\dots$	$d_n^2$	$T_C^2$
$\dots$	$\dots$	$\dots$	$\dots$	$\dots$	$\dots$
合金 $n$	$d_1^n$	$d_2^n$	$\dots$	$d_n^n$	$T_C^n$

(背後にある隠れた前提)  
データ対象間の類似度は  
ベクトル演算によって定量評価する

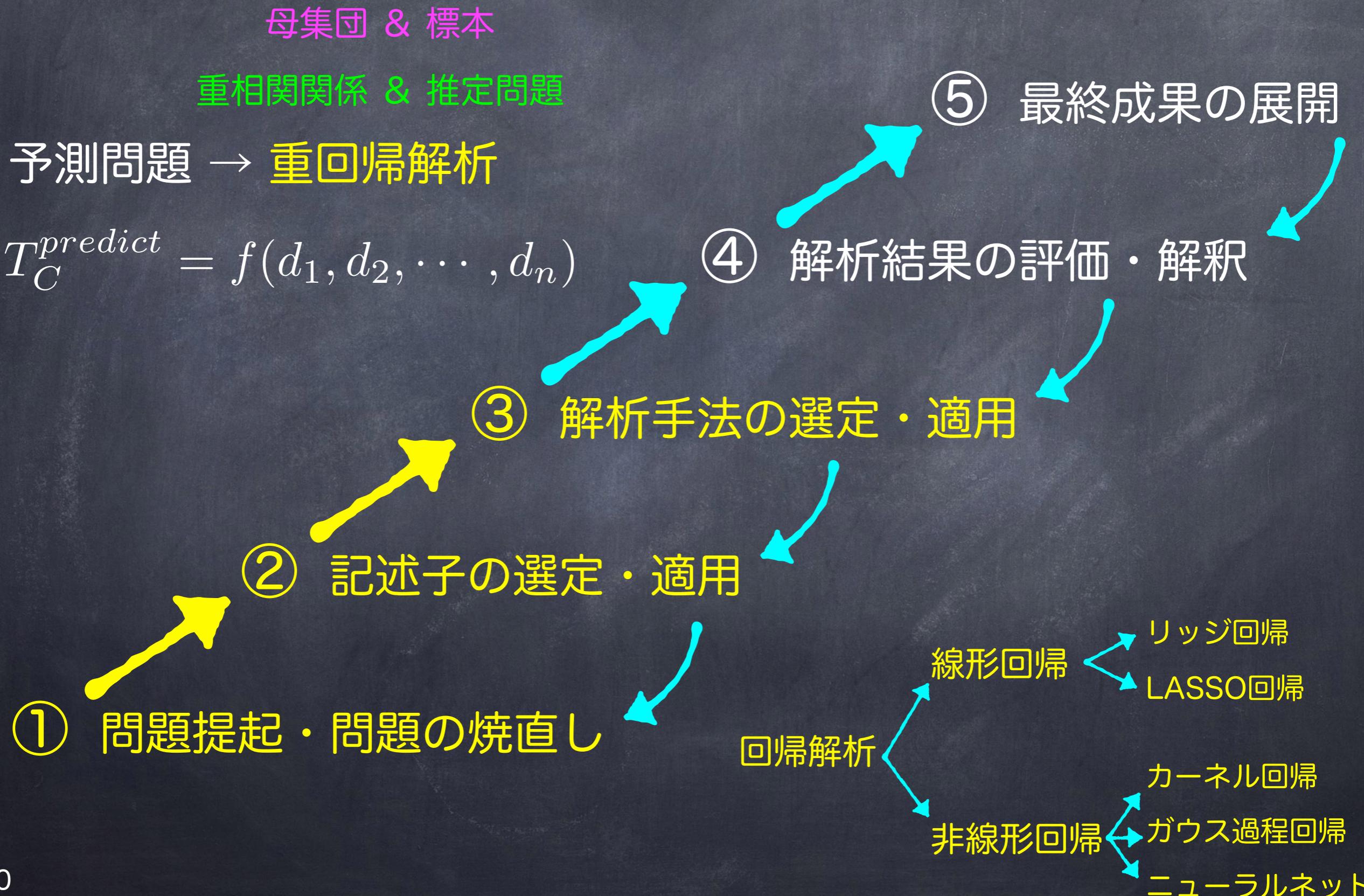
# MIの反復5ステップ

予測問題 → 重回帰解析

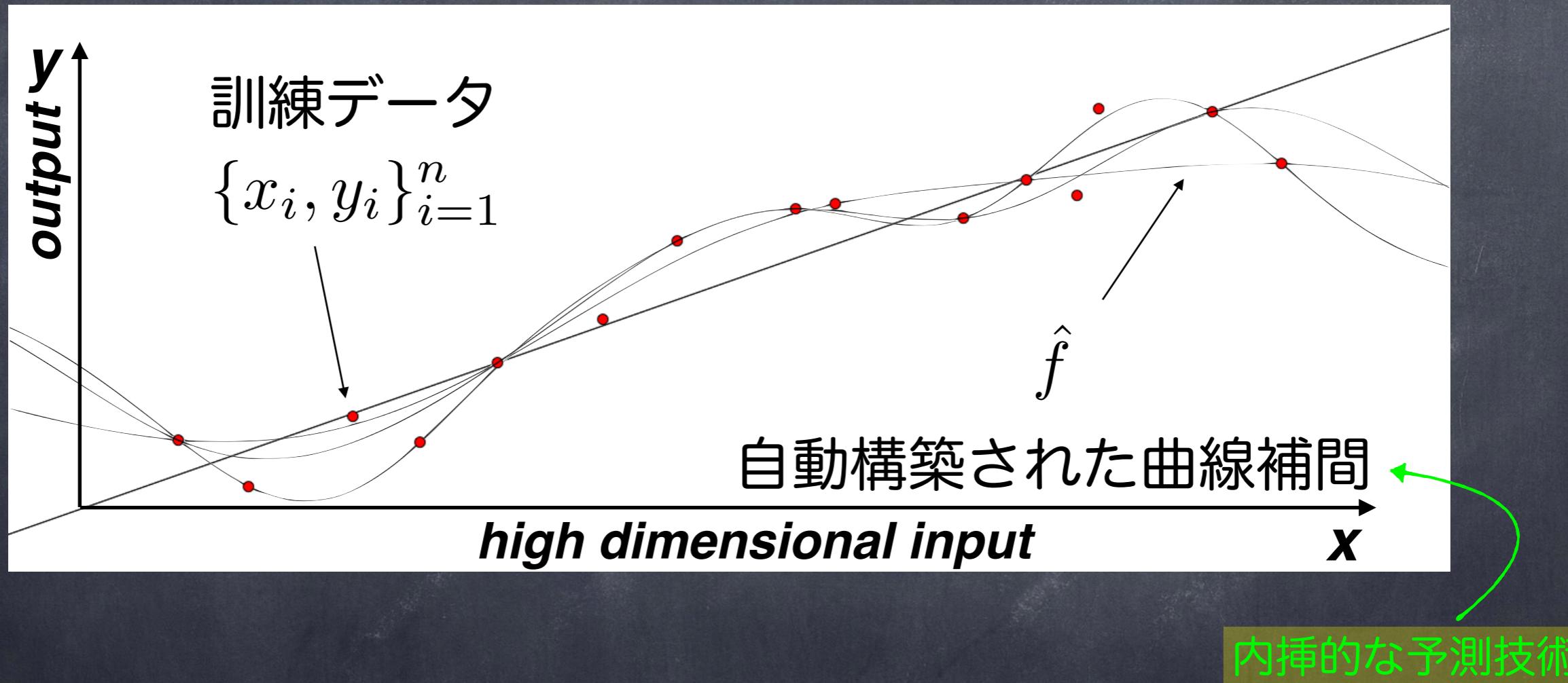
$$T_C^{predict} = f(d_1, d_2, \dots, d_n)$$



# MIの反復5ステップ



# Supervised learning: High dimensional curve fitting



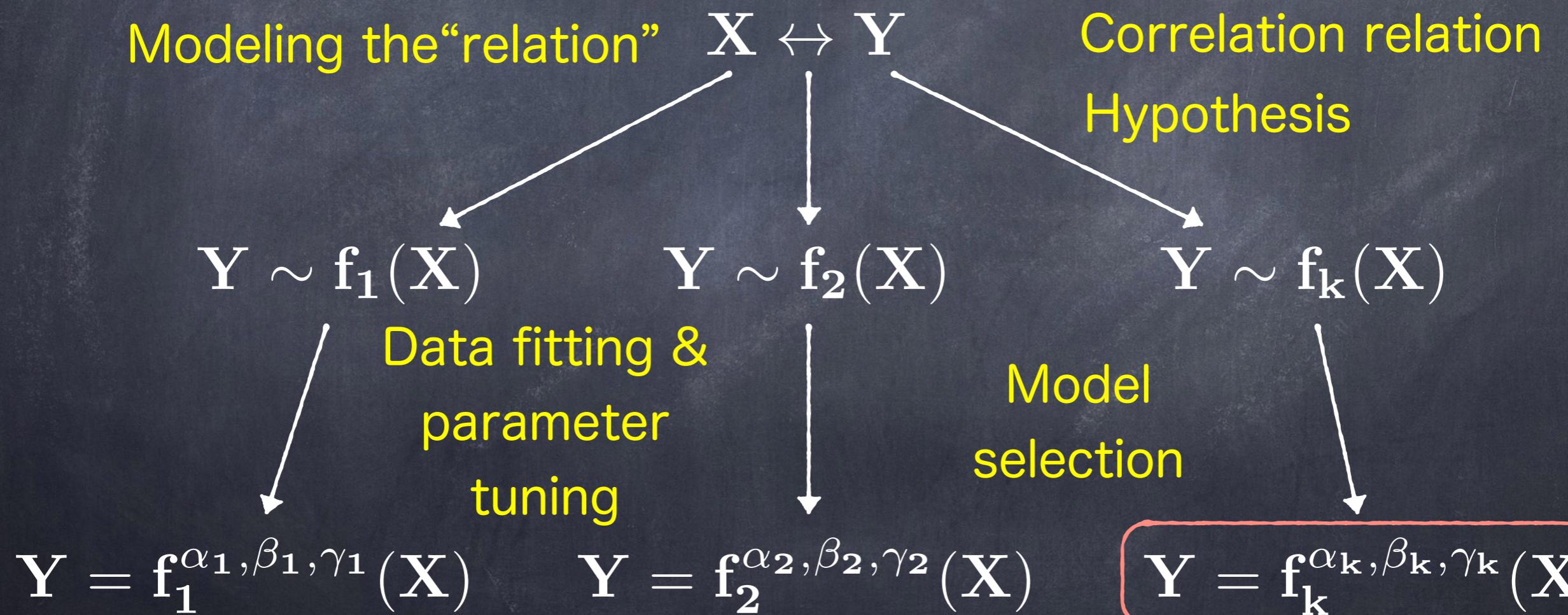
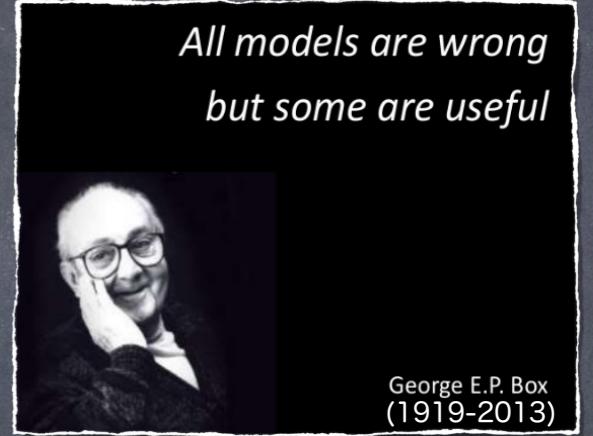
Tools libraries:

scikit-learn, Weka, DeepLearning (Tensorflow, pyTorch, Chainer, …)

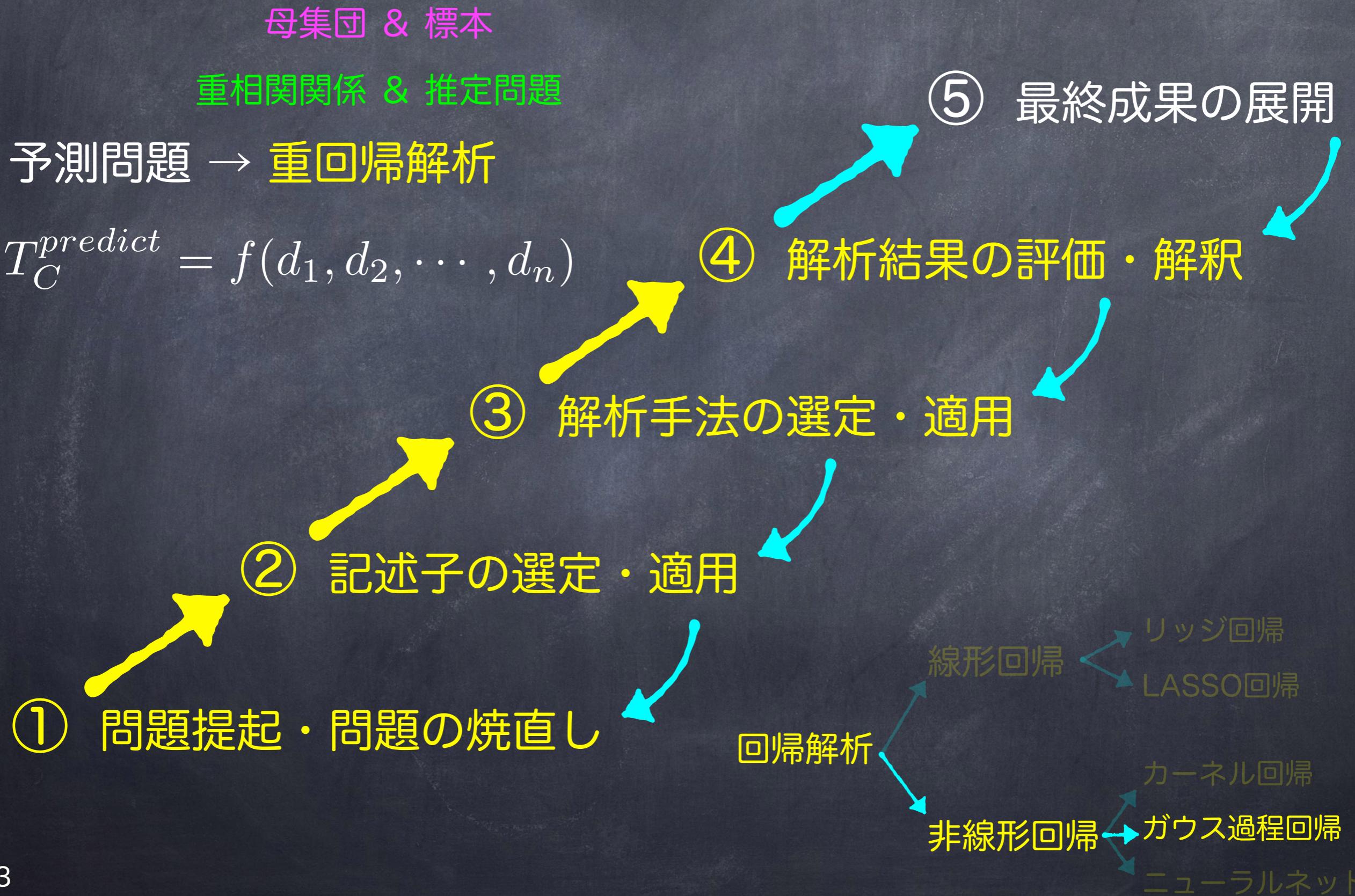
Spark ML, Apple CoreML, Google Cloud AI/Cloud ML, Amazon ML, …

# Regression and Data fitting

Data  $\{X_i, Y_i\}$



# MIの反復5ステップ



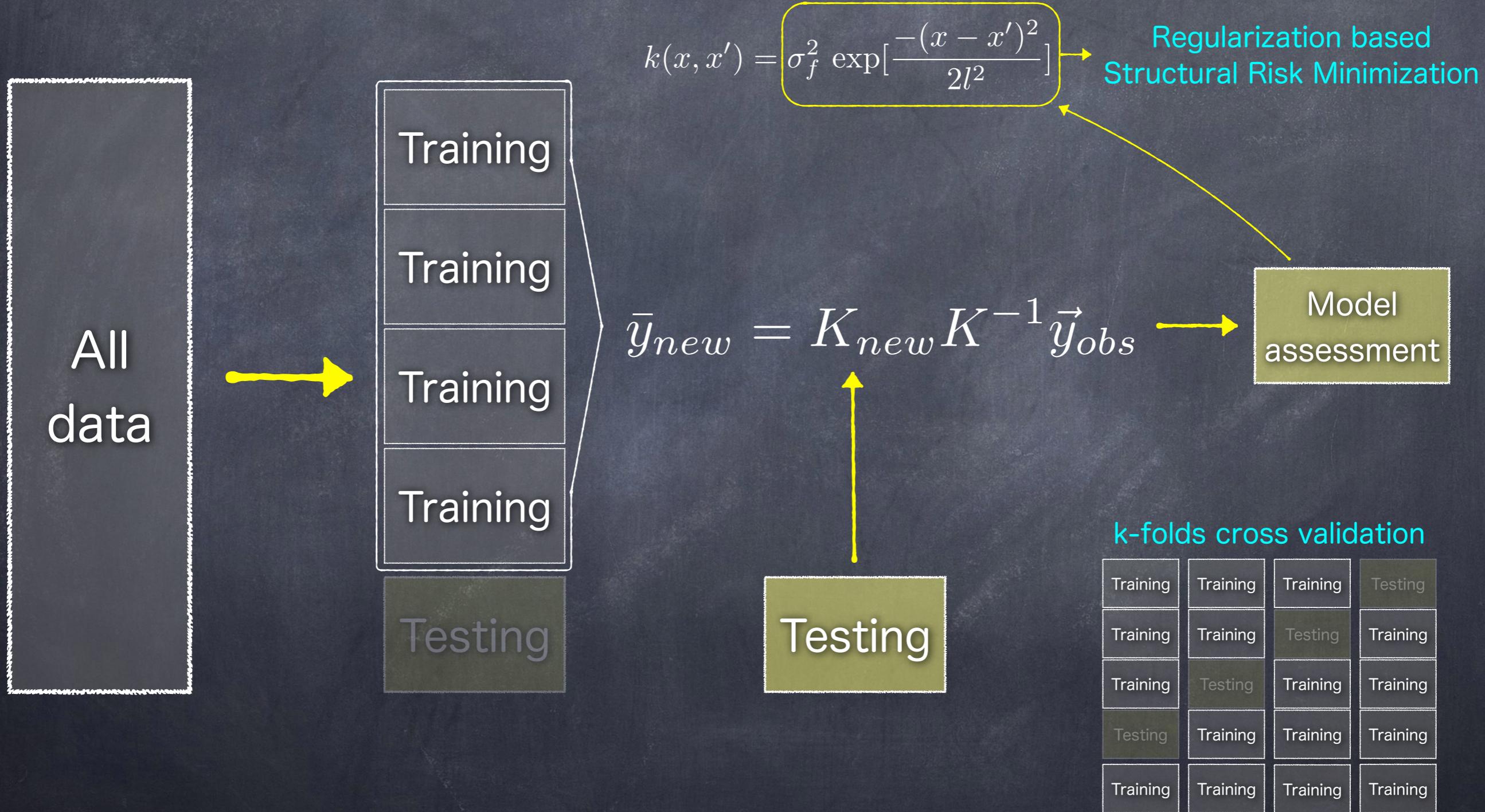
# Model assessment

Statistical estimation  
交差検定法

1. Predictive accuracy - 予測精度
2. Speed - 予測速度
3. Robustness - 頑健性
4. Scalability - 拡張性
5. Interpretability - 解釈容易性

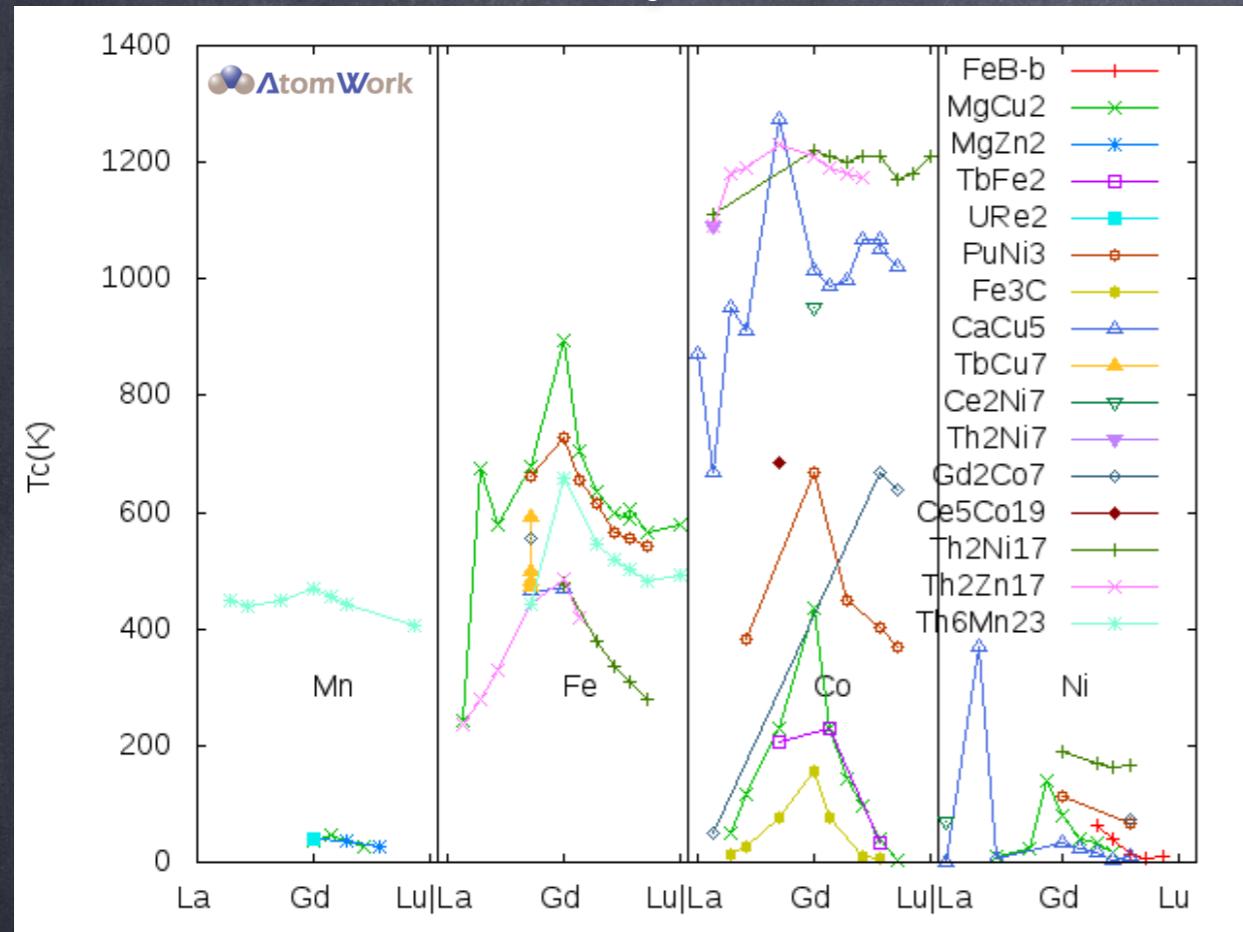
“Use and Abuse of Regression”, George Box,  
活用 亂用 Technometrics, Vol. 8, No. 4, (1966)

# Model assessment by cross validation

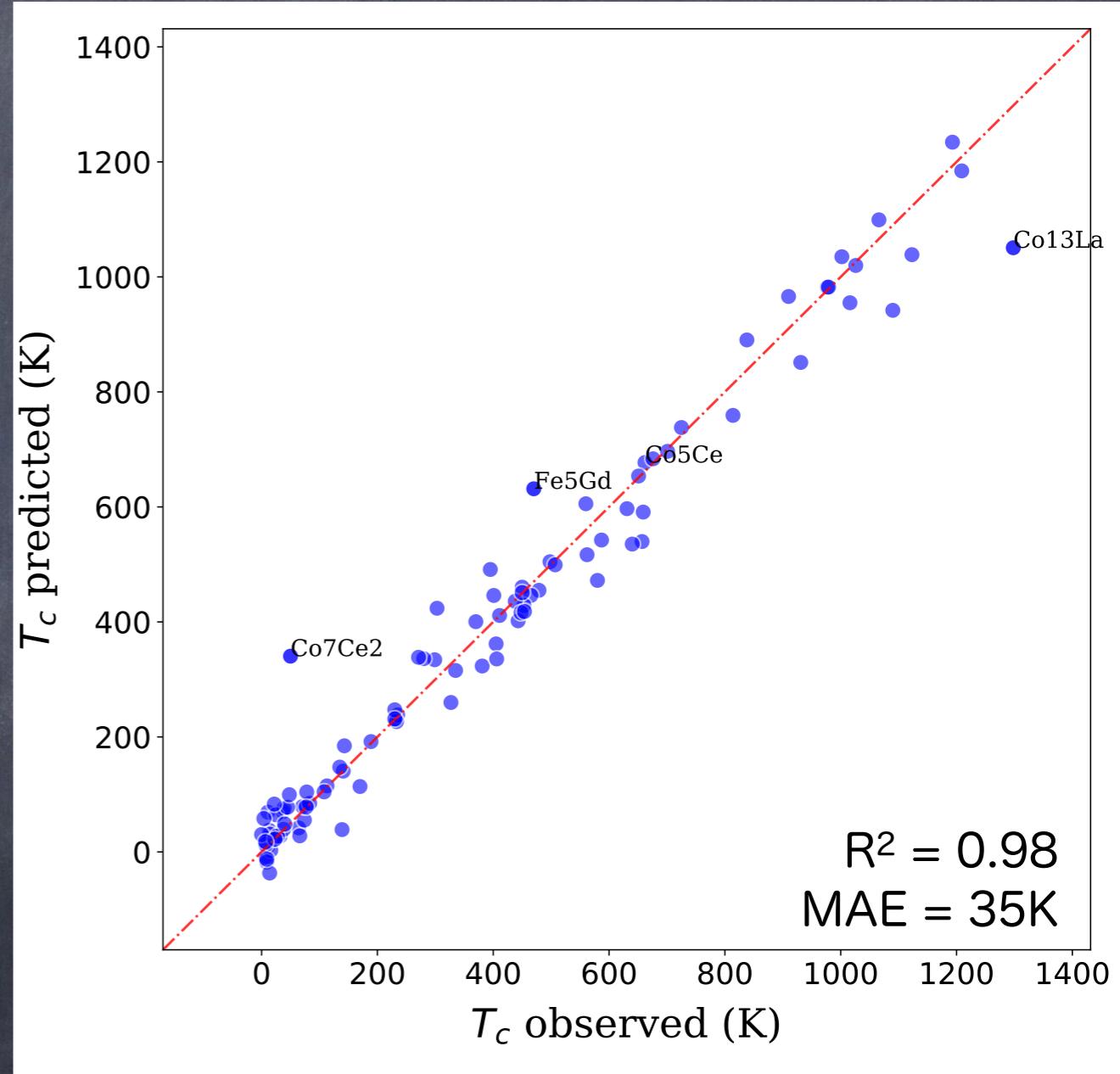


# Prediction of Curie temperature

108 d-f bimetal alloys from literatures



交差検定による  $T_c$  予測モデルの評価

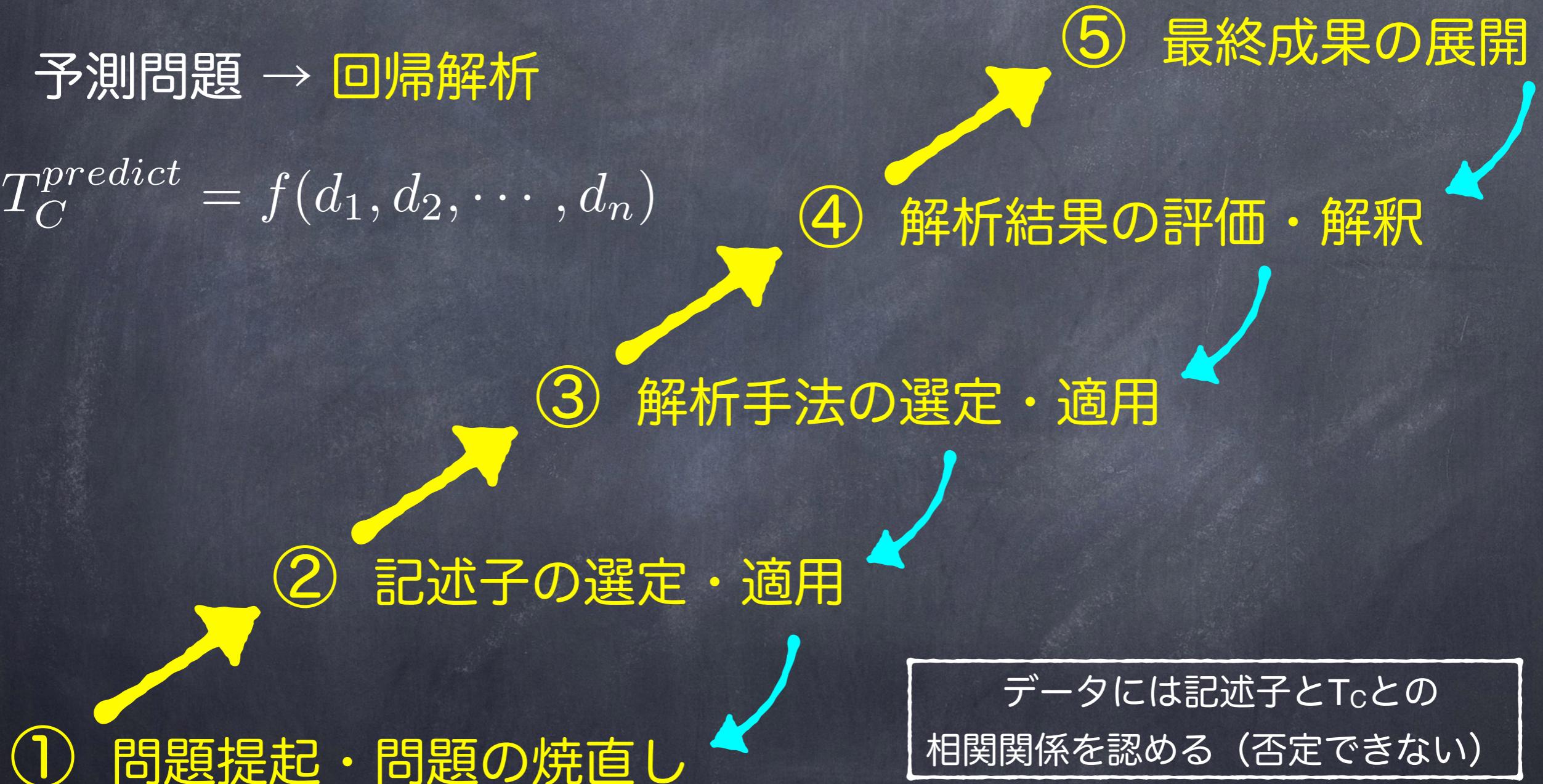


組成が違う新規合金の  $T_c$  を予測できるか？

# MIの反復5ステップ

予測問題 → 回帰解析

$$T_C^{predict} = f(d_1, d_2, \dots, d_n)$$

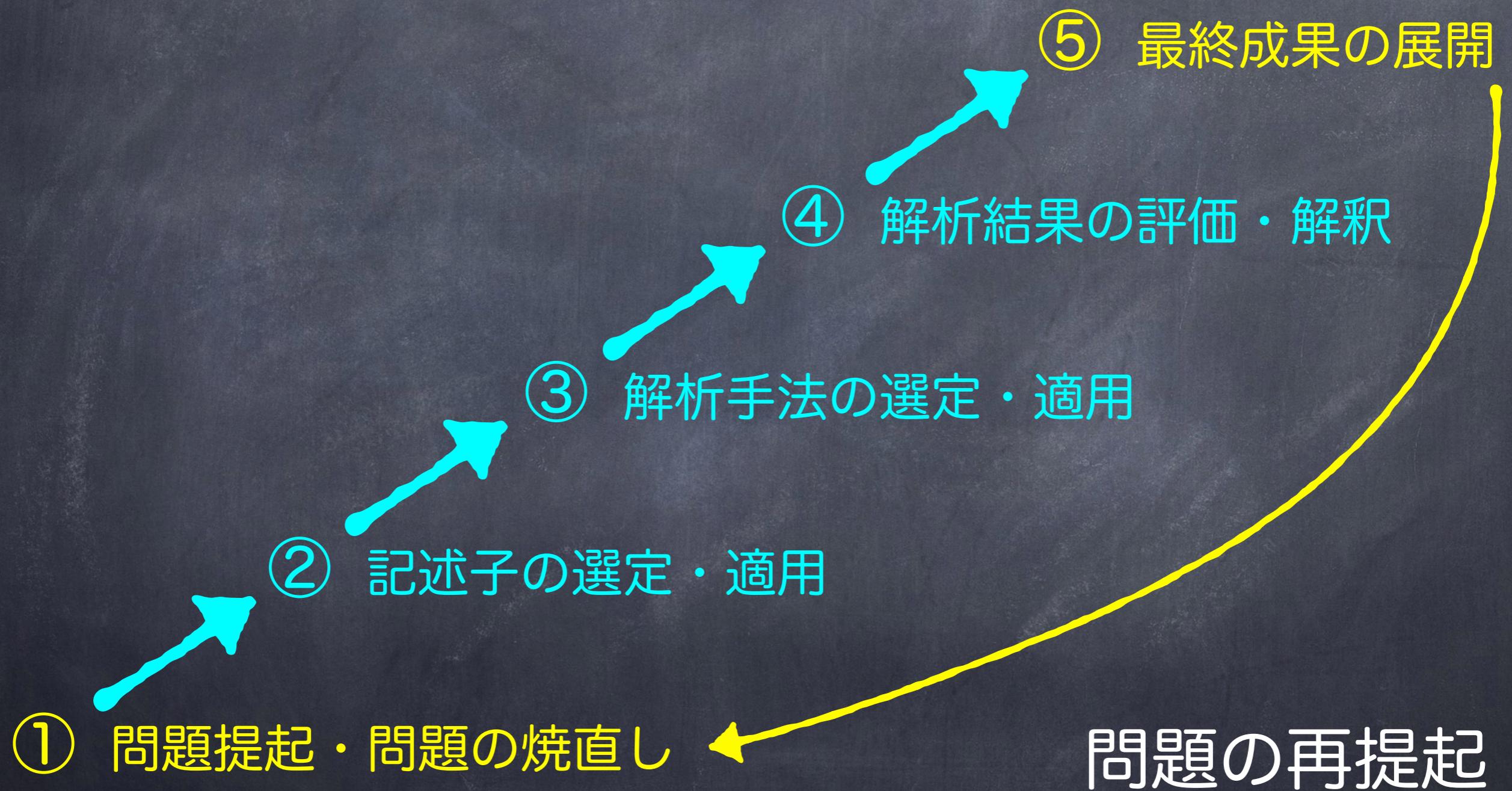


データには記述子と  $T_c$  との  
相関関係を認める (否定できない)

組成が違う新規合金の  $T_c$  を予測できるか？

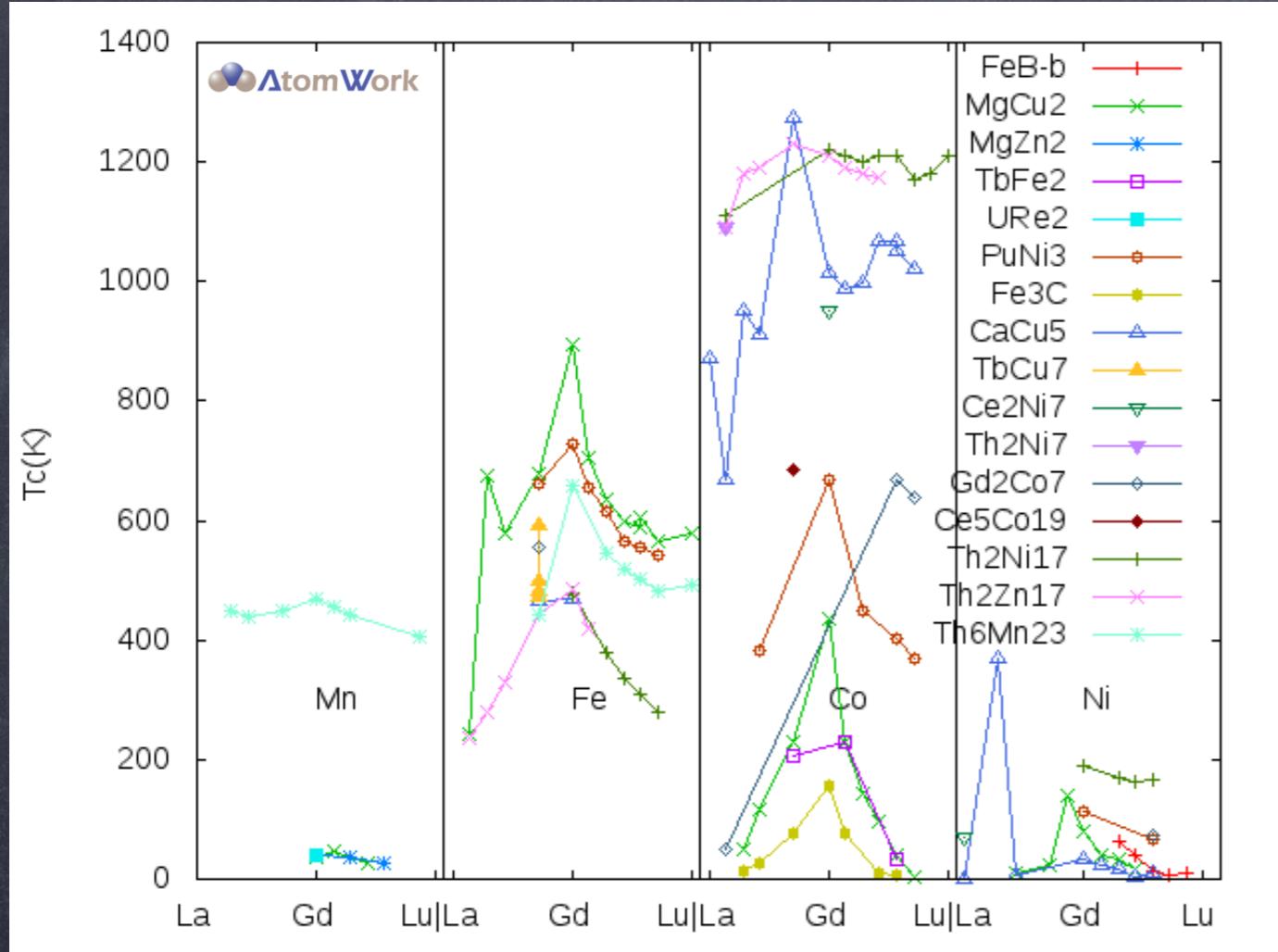
外挿問題にぶつかる！

# MIの反復5ステップ



# 問題の再提起：キュリ温度の決める機構

d-f bimetal alloys from literatures



Material structure

$a, b, c, \alpha, \beta, \gamma$

$atom_1, x_1, y_1, z_1$

$atom_2, x_2, y_2, z_2$

⋮

$atom_i, x_i, y_i, z_i$

⋮

$atom_n, x_n, y_n, z_n$

記述子

「どうなっている？」  
と繰り返して問う

記述子と Tc との相関関係が観測された

「なぜ？」  
と問う

次には

Tc の決める機構を理解するためのヒントが欲しい！

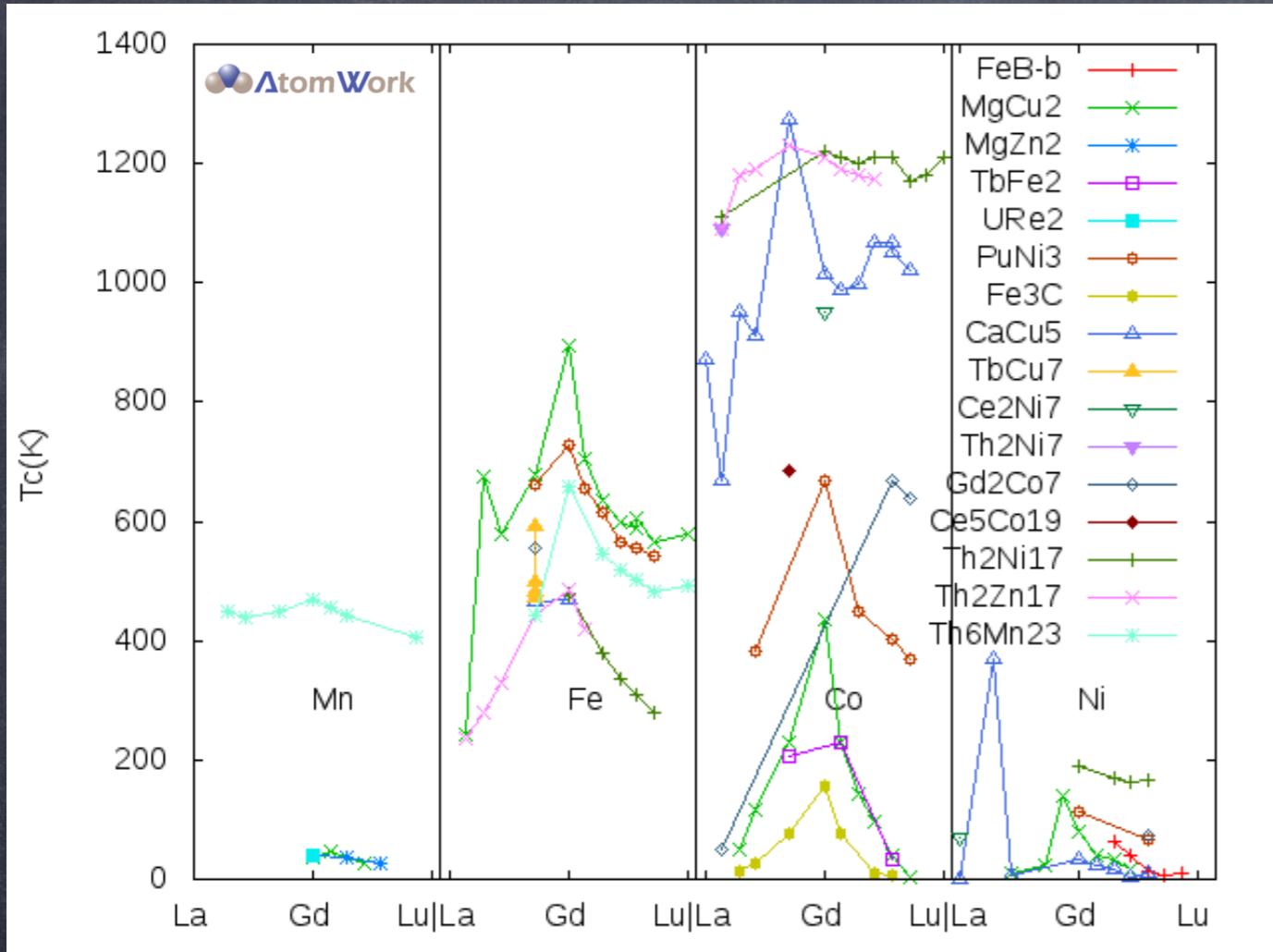


如何にヒントを抽出するか？

$$T_C^{predict} = f(\square, \diamond, \triangle, \dots)$$

# Curie temperature of 3d-4f binary alloys

d-f bimetal alloys from literatures



Material structure

$a, b, c, \alpha, \beta, \gamma$

$atom_1, x_1, y_1, z_1$

$atom_2, x_2, y_2, z_2$

⋮

$atom_i, x_i, y_i, z_i$

⋮

$atom_n, x_n, y_n, z_n$

記述子

「どうなっている？」  
と繰り返して問う

記述子と  $T_c$  との相関関係が観測された

「なぜ？」  
と問う

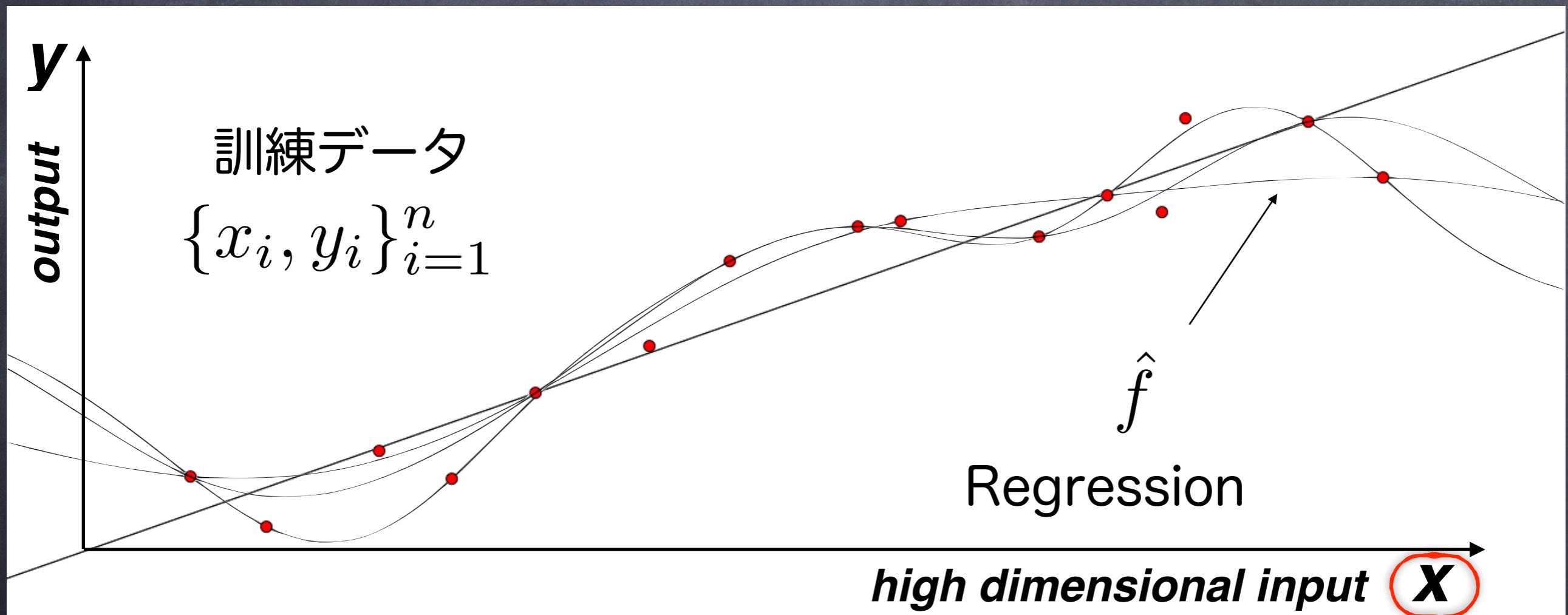
次には

$T_c$  の決める機構を理解するためのヒントが欲しい！

理解可能な構造物性の潜む関係構造を抽出したい！

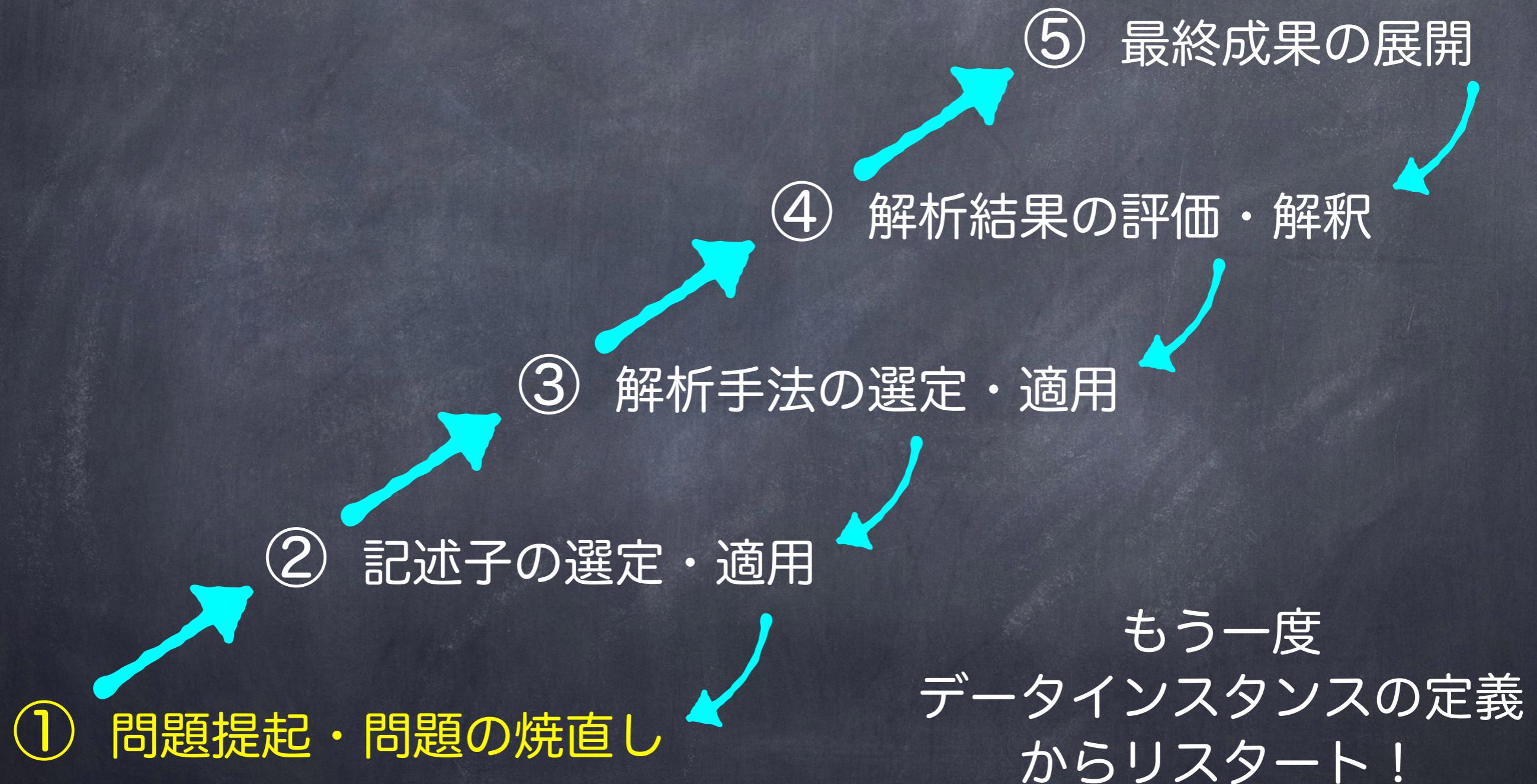
$$T_C^{predict} = f(\square, \diamond, \triangle, \dots)$$

# Supervised learning: High dimensional curve fitting



予測精度に効く記述子はどれか？

# MIの反復5ステップ



# Regression based feature selection

想定する記述子

Structural information:

- R-R Nearest neighbor distance
- R-T Nearest neighbor distance
- T-T Nearest neighbor distance
- Number of R surrounding R
- Number of T surrounding T
- Number of R surrounding T
- Number of T surrounding T
- Concentration of R
- Concentration of T

Atomic information:

- Atomic Number of  $Z_R, Z_T$
- Covalent Radius of  $r_R, r_T$
- Electron Negativity of  $\chi_R, \chi_T$
- Ionization Potential  $IP_R, IP_T$
- Orbital angular momentum  $L_R, L_T$
- Spin angular momentum  $S_R, S_T$
- Total angular momentum  $J_R, J_T$
- Landé g-factor  $g_{JR}, J_R g_{JR}, J_R(1-g_{JR})$

ガウス過程回帰



$T_c^{predict}$

$$T_c^{predict} = f_{gp}(x^1, \dots, x^m)$$

Model data

$$\{\{x^1, \dots, x^m\}, R^2(T_c^{obs}, T_c^{predict})\}$$

$2^{27}$

All combinations



$$f_{gp}(x^1, \dots, x^m)$$



Top N-best  
prediction accuracy

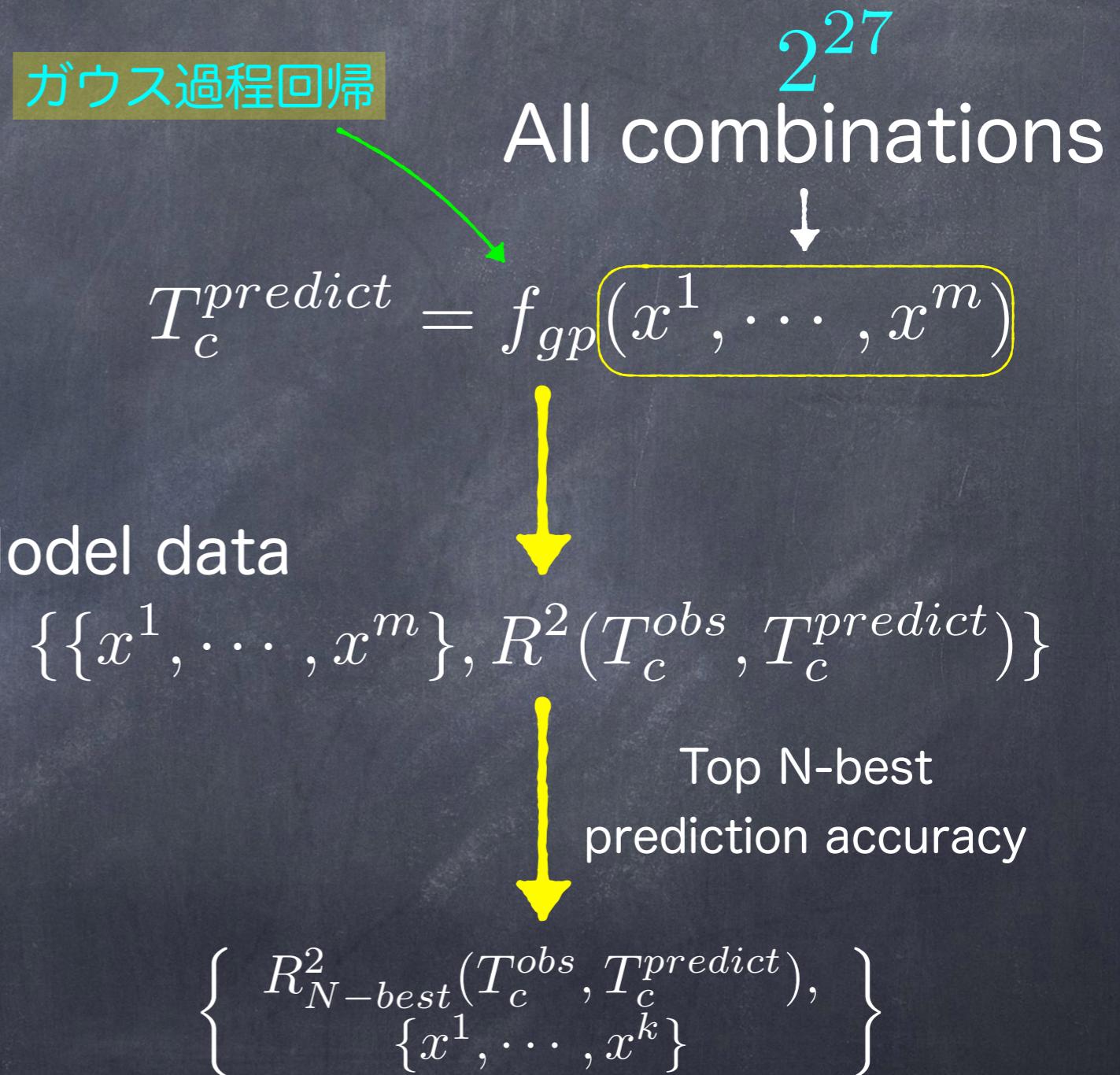
$$\left\{ \begin{array}{l} R^2_{N-best}(T_c^{obs}, T_c^{predict}), \\ \{x^1, \dots, x^k\} \end{array} \right\}$$



予測精度に効く記述子はどれか？

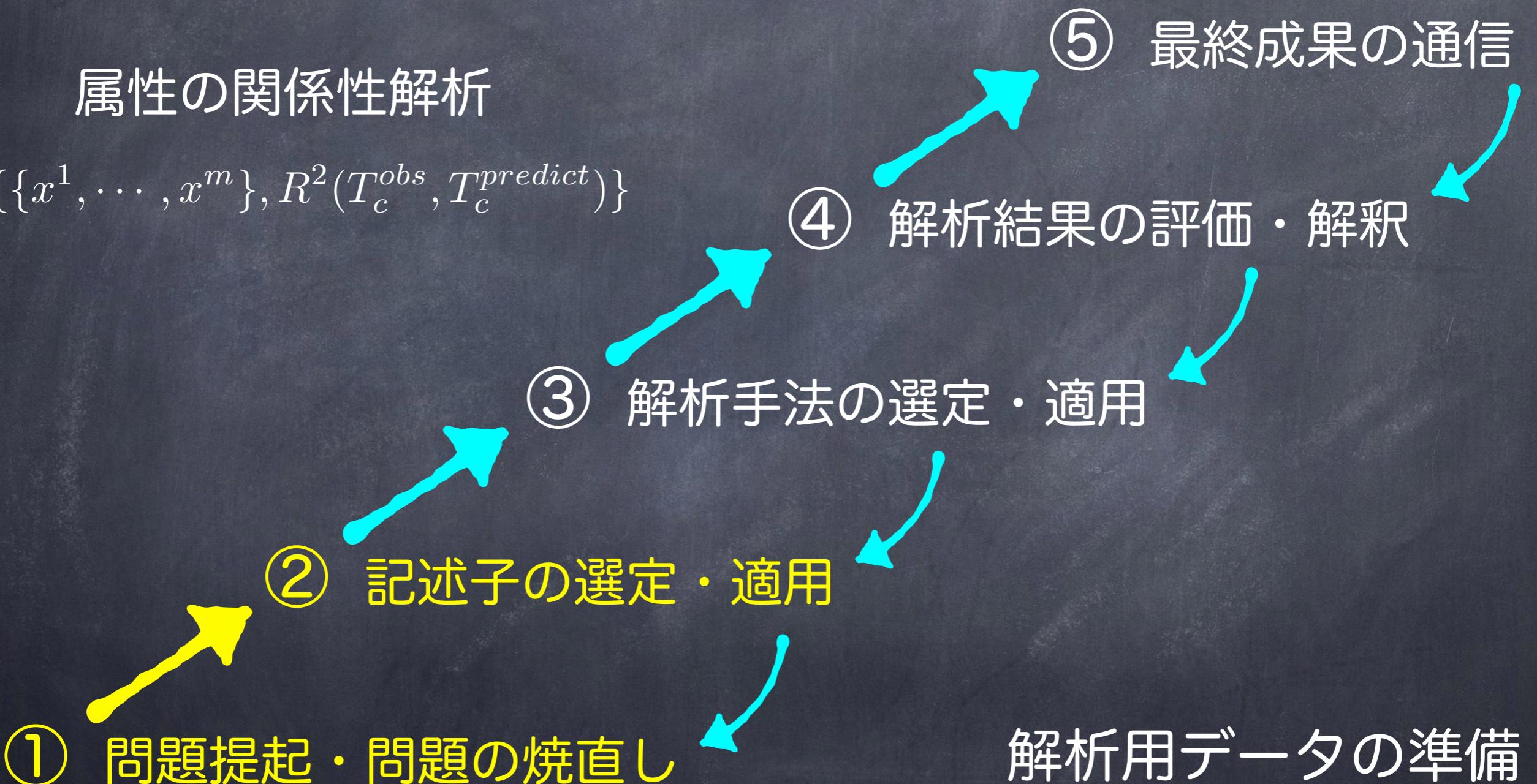
# Regression based feature selection

データインスタンス：  
記述子の組み合わせ  
 $2^{27}$  記述子の全組み合わせ



予測精度に効く記述子はどれか？

# MIの反復5ステップ

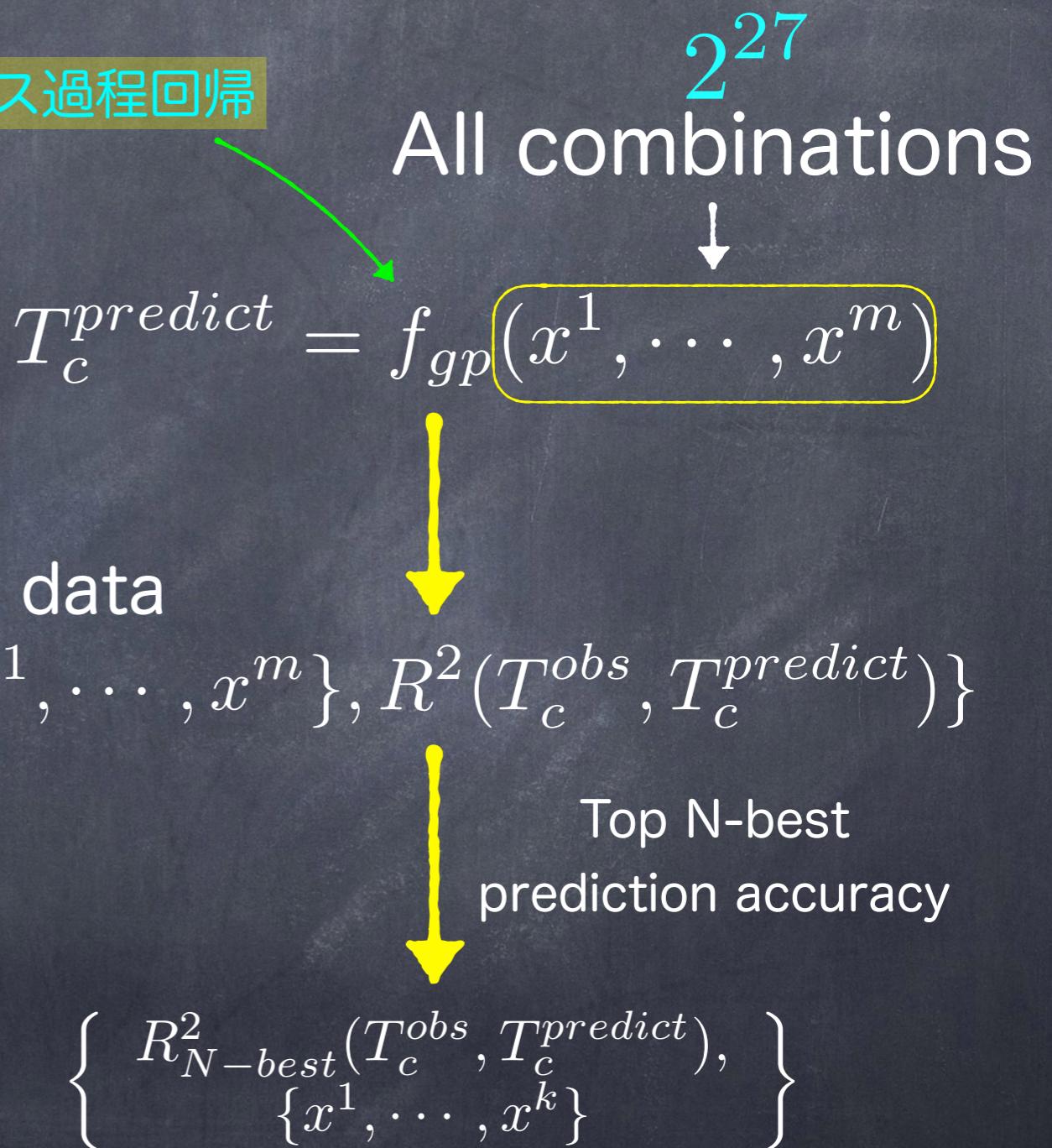


# Regression based feature selection

解析用データ

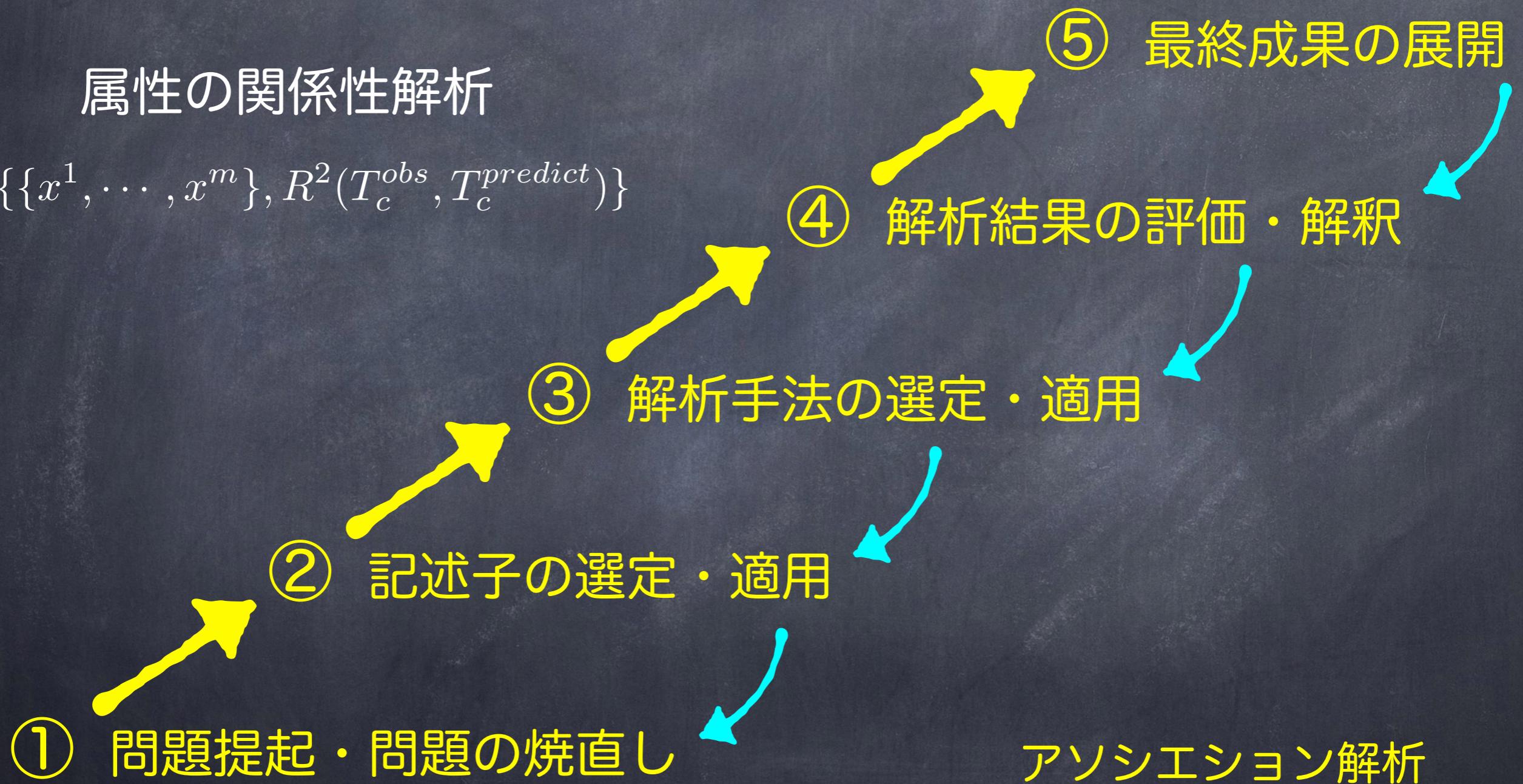
<i>ID</i>	<i>Descriptors</i>	$R^2$
<i>Model 1</i>	$\{d_i, d_j, \dots, d_k\}$	$R^2_1$
<i>Model 2</i>	$\{d_l, d_m, \dots, d_n\}$	$R^2_2$
...	...	...
<i>Model n</i>	$\{d_o, d_p, \dots, d_q\}$	$R^2_n$

ガウス過程回帰

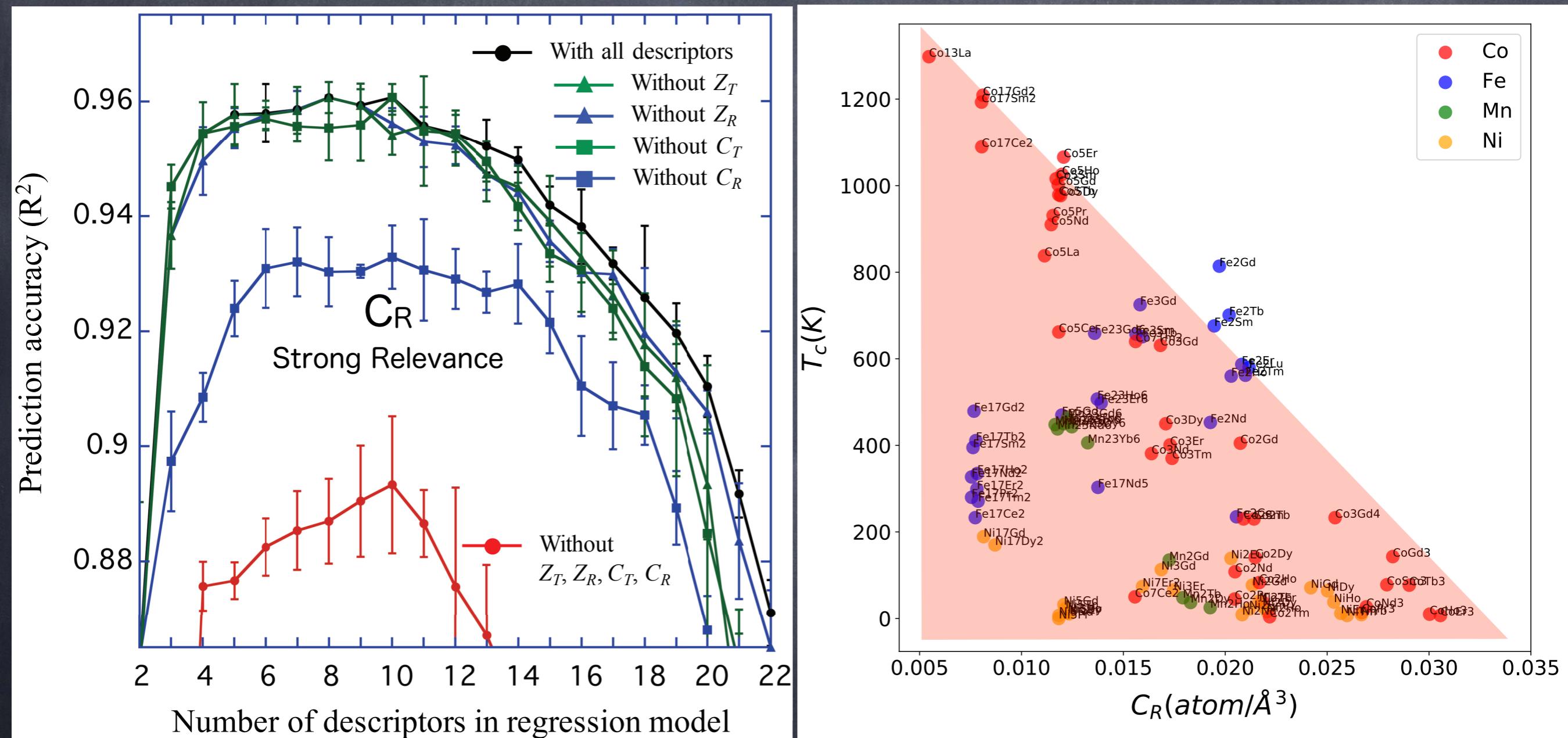


予測精度に効く記述子はどれか？

# MIの反復5ステップ



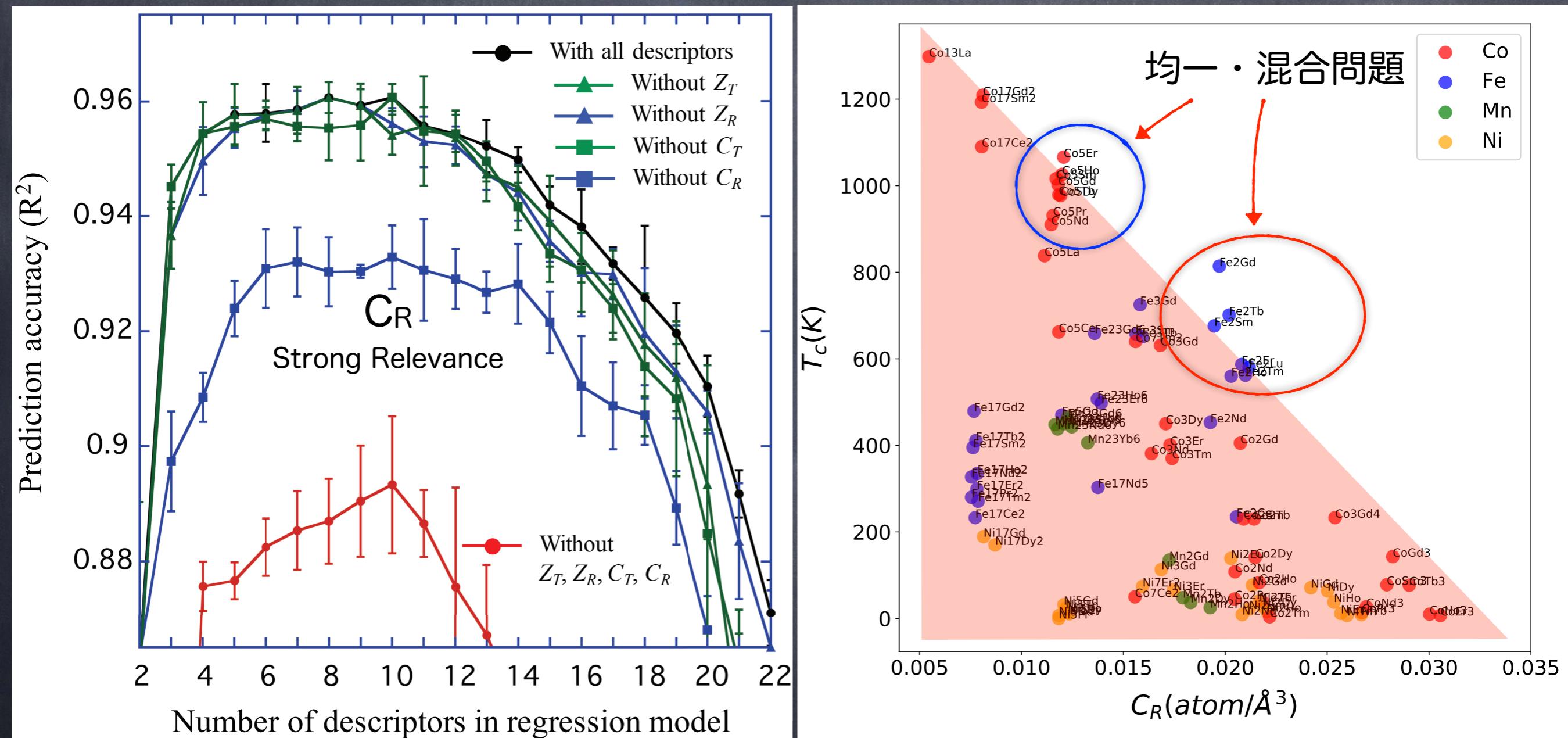
# Regression based descriptor evaluation



$T_c$ の予測精度と活用する記述子との関係性

希土類の濃度が $T_c$ を決定する機構に重要だ！

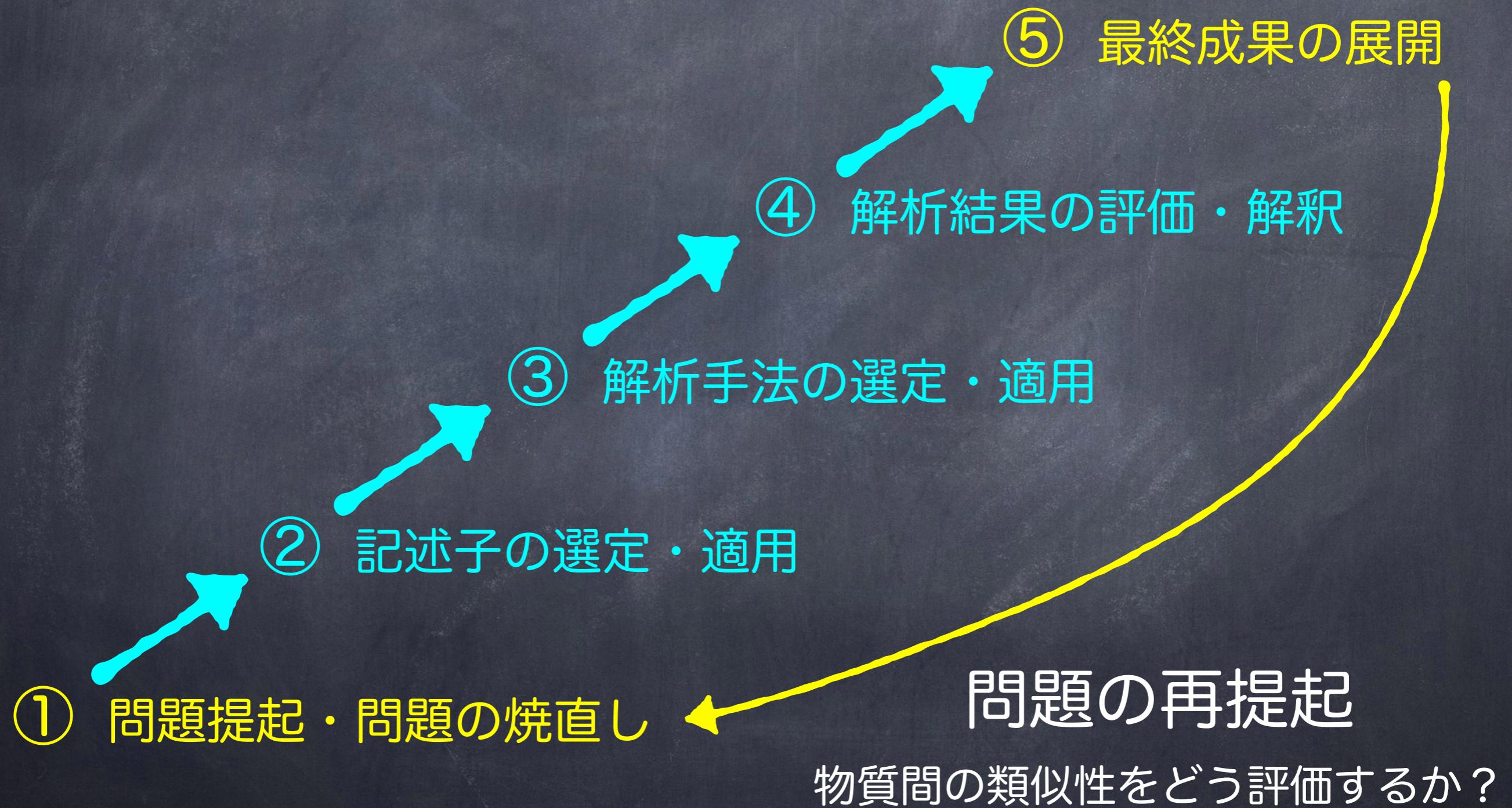
# Regression based descriptor evaluation



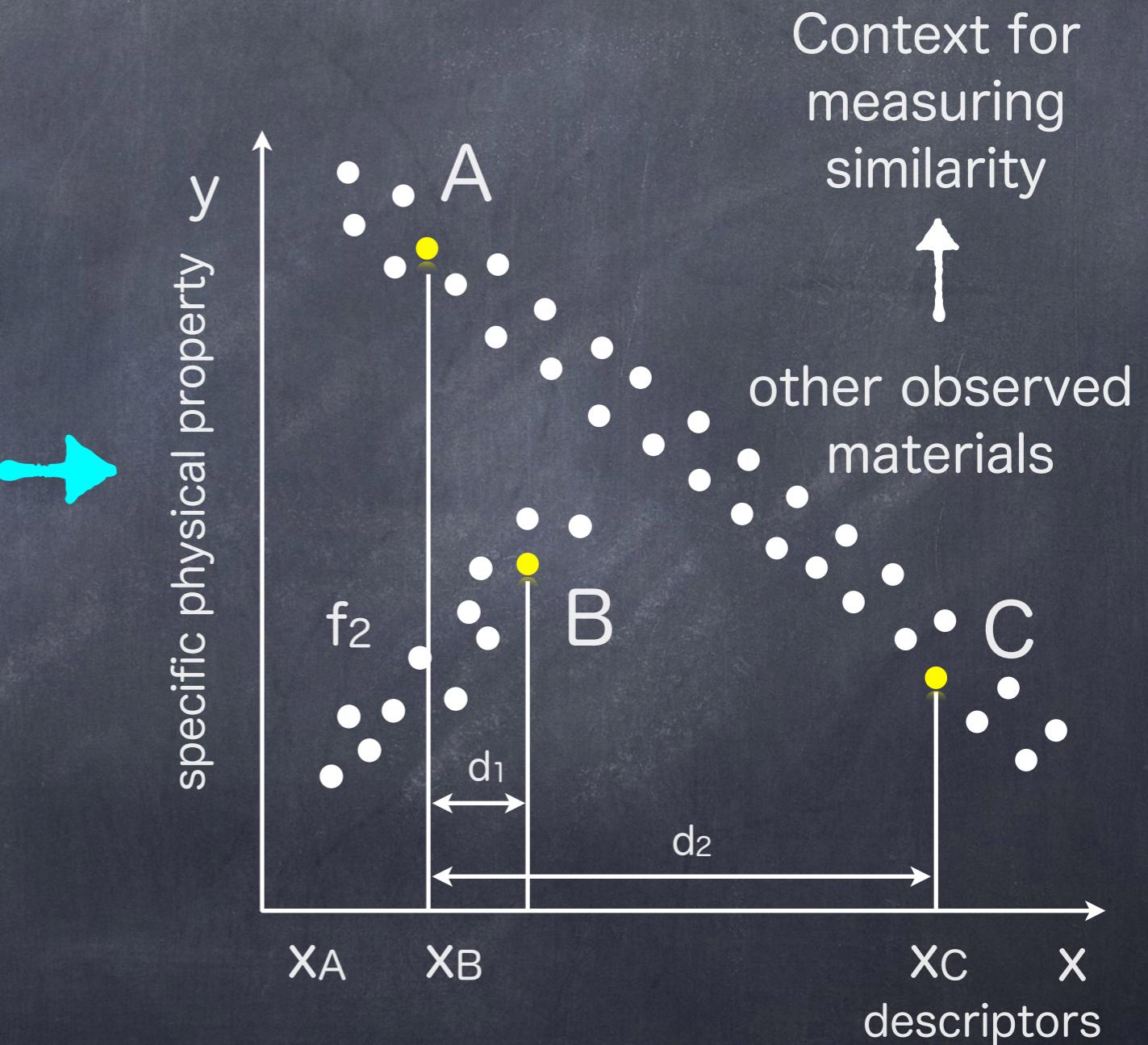
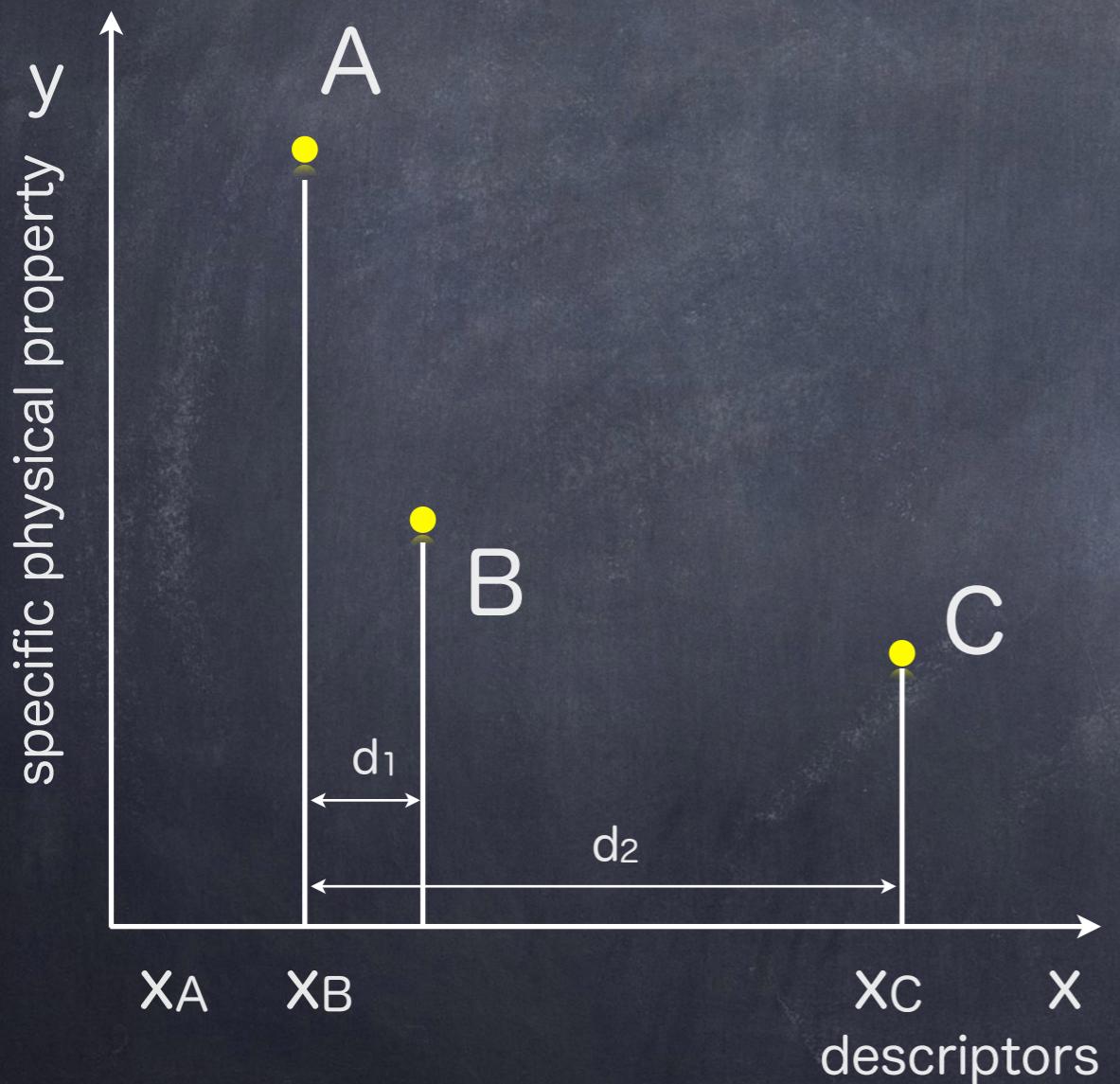
$T_c$ の予測精度と活用する記述子との関係性

希土類の濃度が $T_c$ を決定する機構に重要だ！

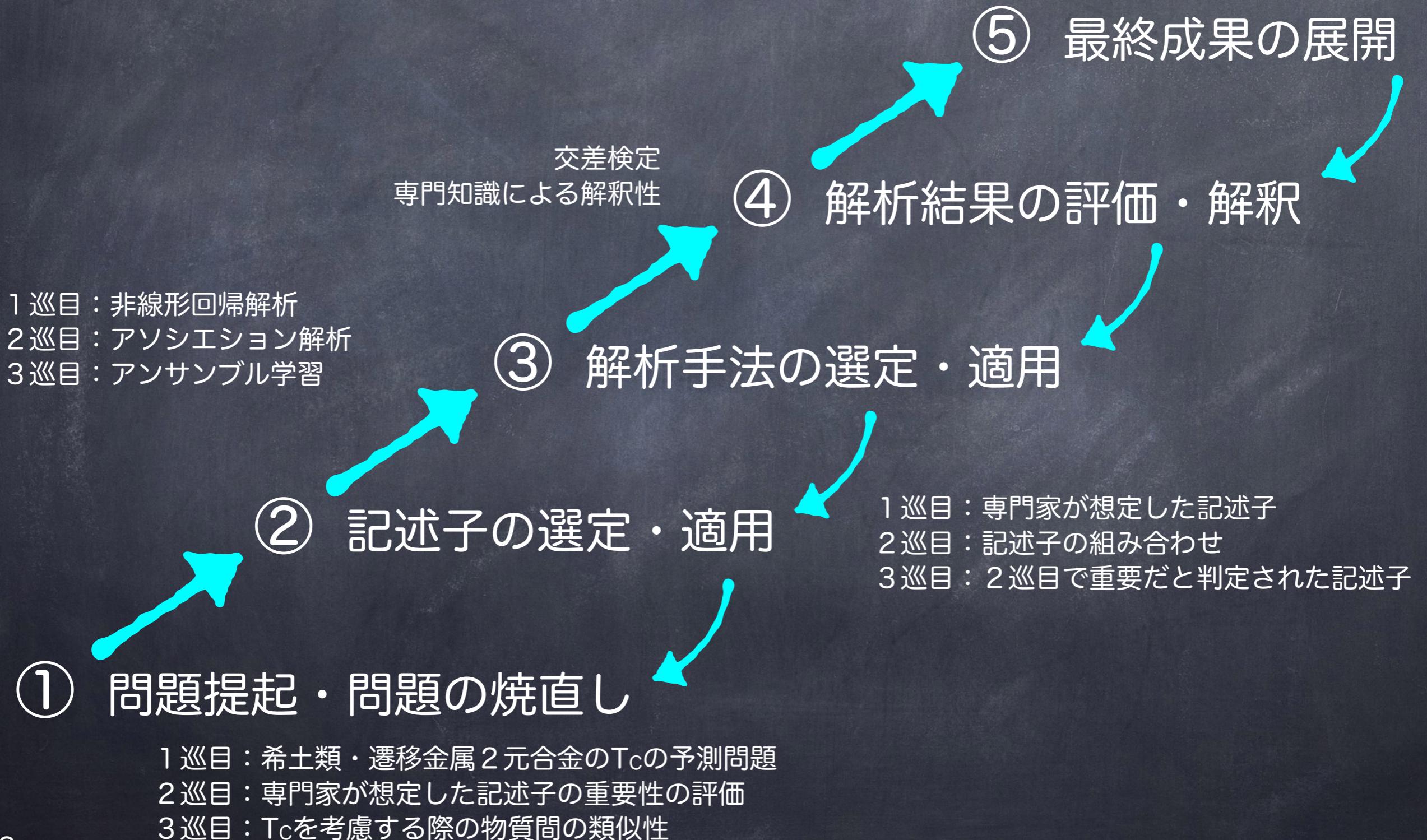
# MIの反復5ステップ



# Development of context-based method for measuring similarity between materials



# Example のまとめ



# **UNCERTAINTY AND DECISION-MAKING FOR MATERIALS DISCOVERY**

Hieu-Chi Dam

Japan Advanced Institute of Science and Technology



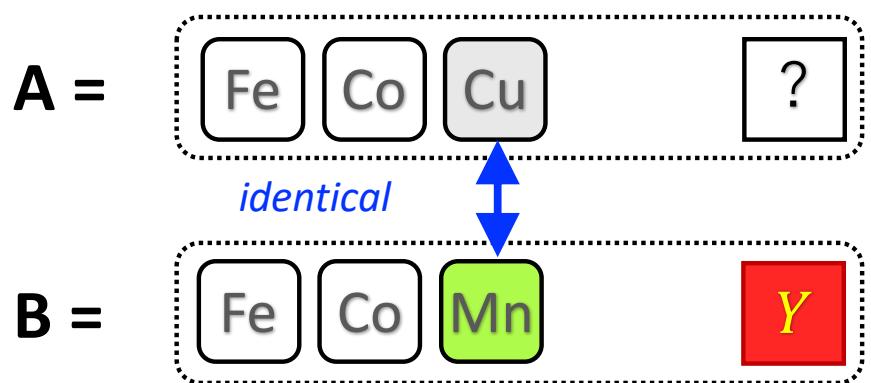
# Material Dataset: An Example with Alloys

Note: descriptions differ from descriptors

Material	Material descriptions	Stability	Magnetization (T)
Material 1	{Cr,Cu,Mo,Tc, structure}	<b>N</b>	0
Material 2	{Co,Cu,Fe,Mn, structure}	<b>Y</b>	0.38
...	...	...	...
Material n	{Ag,Fe,Mn,Zn, structure}	<b>N</b>	0

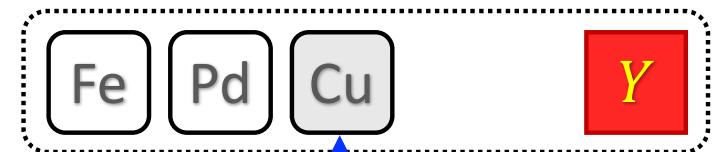
# Similarity: Error of substitution

## Question



## Data

*compositional elements      stability*

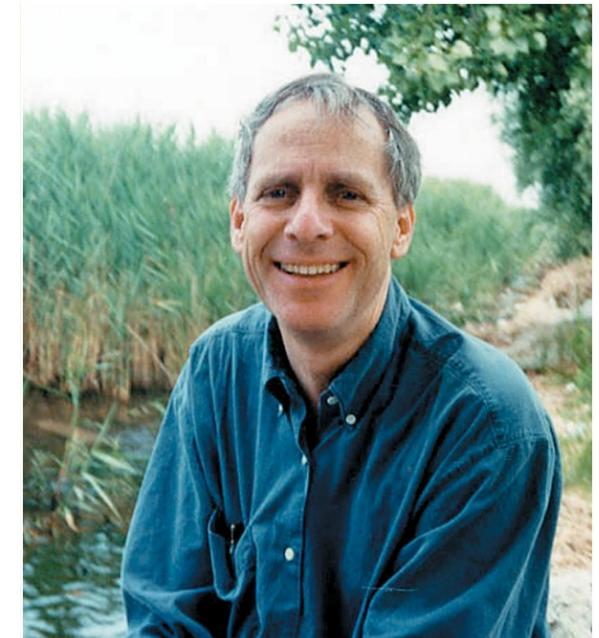


The Key Idea:  
Transforming Data into Evidence of Similarity

# Similarity judgment

Similarity or dissimilarity data appear in different forms:

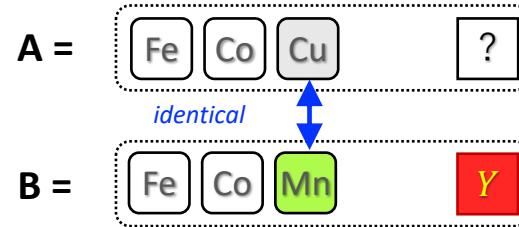
- ✓ Ratings of pairs
- ✓ Sorting of objects
- ✓ Communality between associations
- ➡ ✓ Errors of substitution
- ➡ ✓ Correlations between occurrences



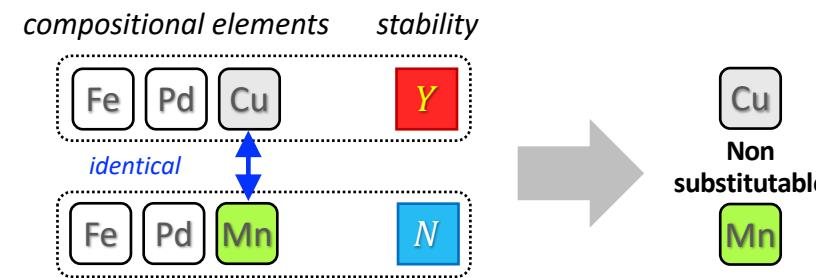
Amos Nathan Tversky  
(1937 – 1996)

# Similarity: Error of substitution

## Question



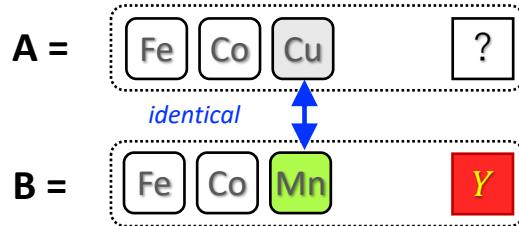
## Data



## Intuitive Evidence

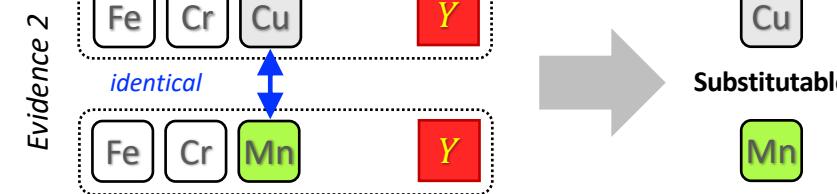
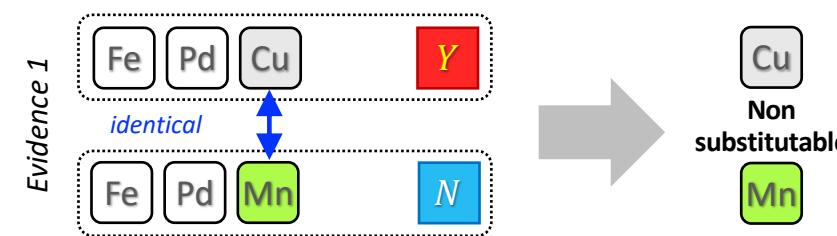
# Similarity: Error of substitution

## Question

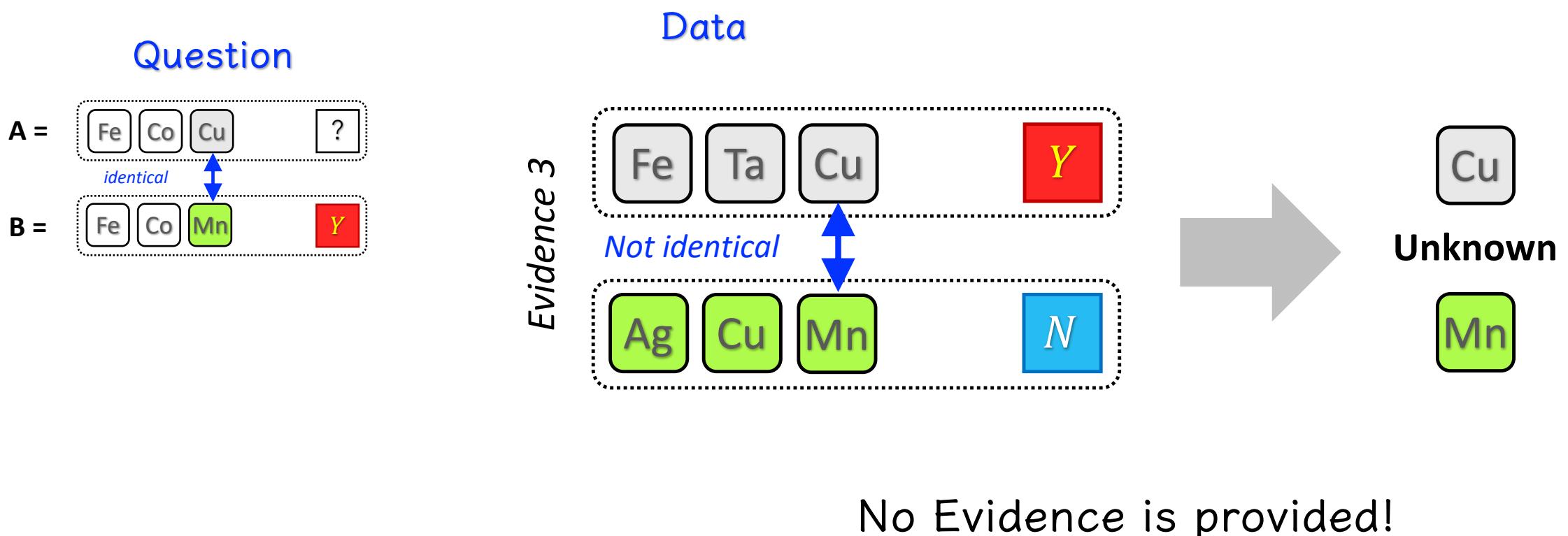


## Data

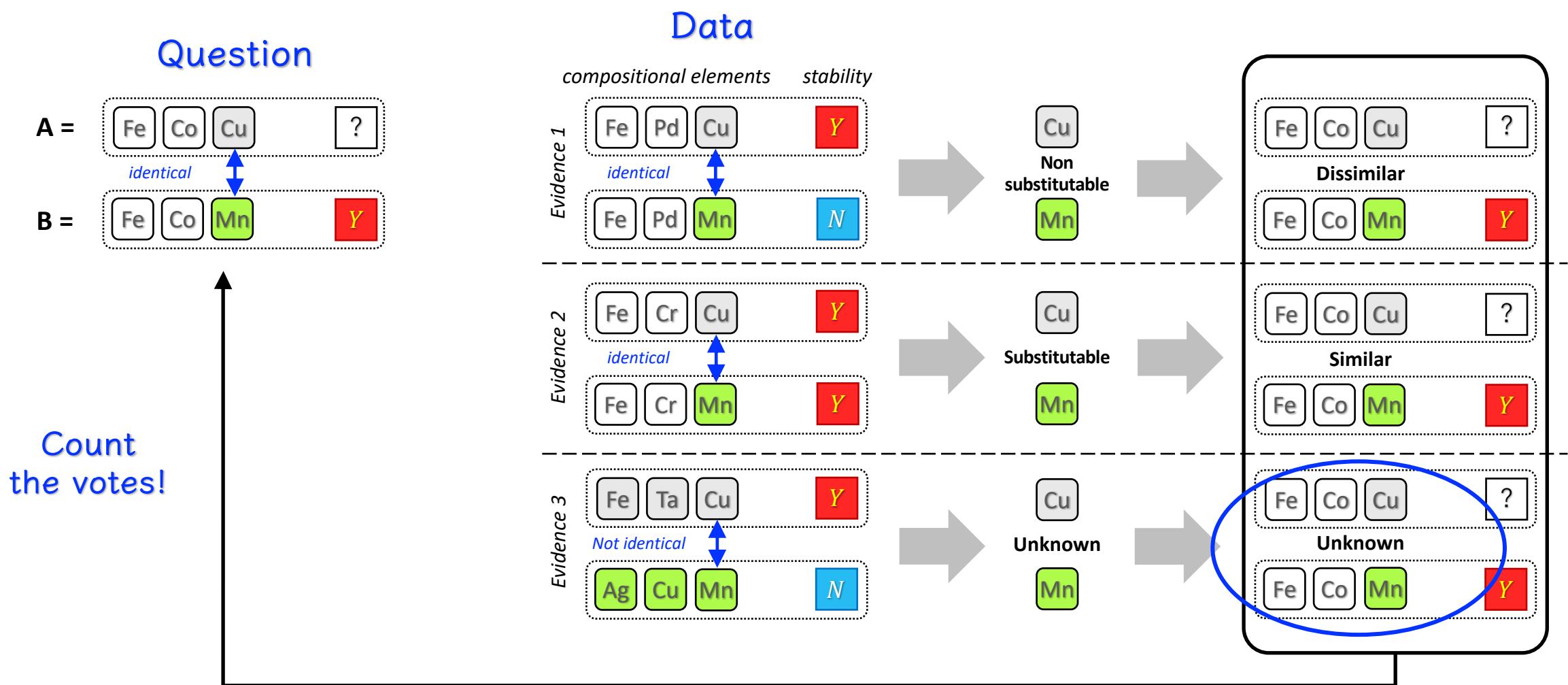
compositional elements      stability



# Similarity: Error of substitution



# Similarity: Error of substitution



# Similarity: Error of substitution

# Question

**A =**

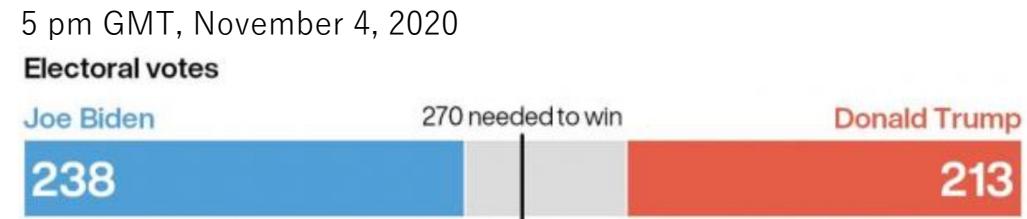
Fe	Co	Cu	?
----	----	----	---

*identical*

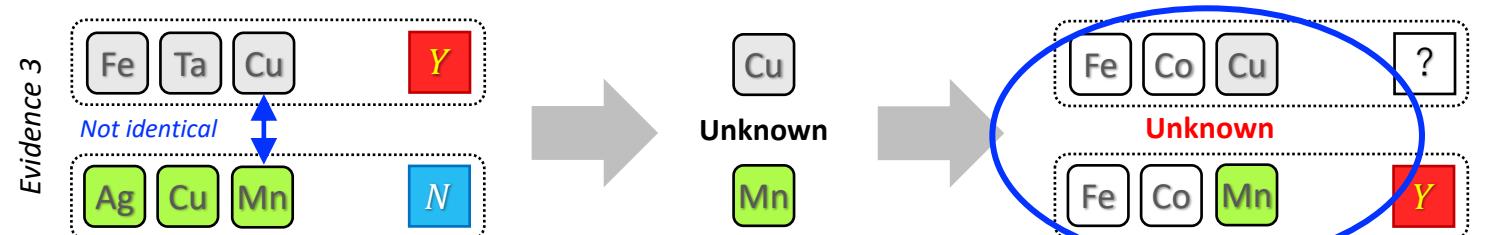
**B =**

Fe	Co	Mn	Y
----	----	----	---

# Uncertainty

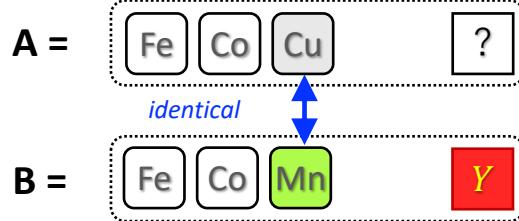


“Unknown” *is not* 50/50



# Similarity: Error of substitution

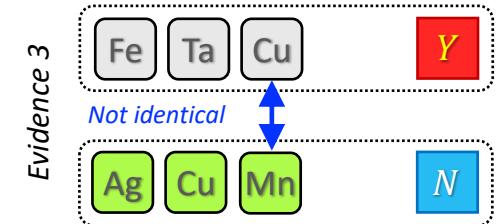
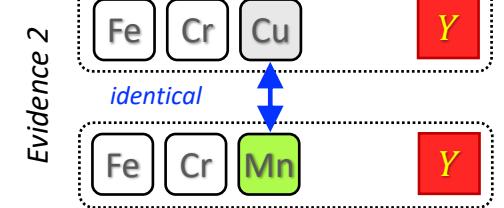
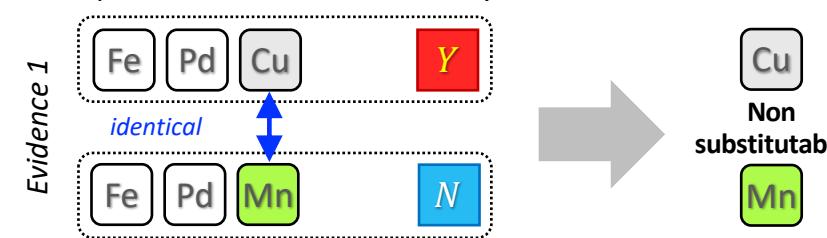
## Question



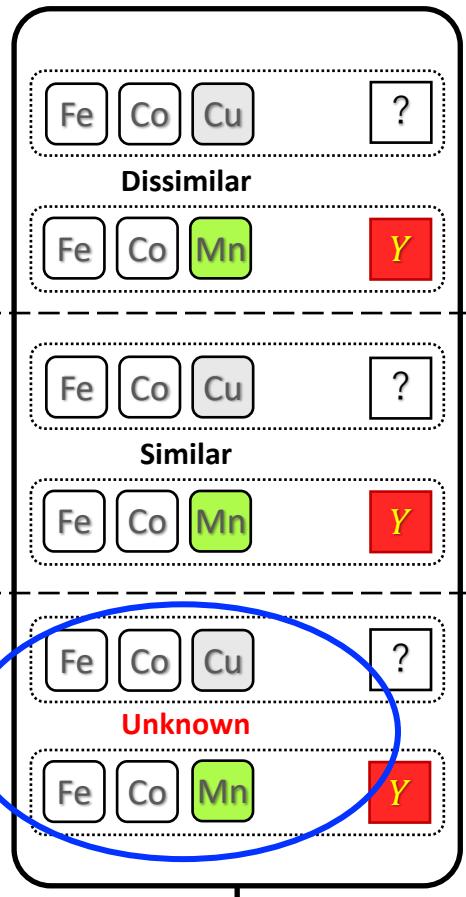
Quantifying the uncertainty?

## Data

compositional elements      stability



Combine evidences



# Dempster–Shafer theory (belief functions)

**Frame of discernments**

$$\Omega_{sim} = \{similar, dissimilar\}$$

Pending judgment of  
Similarity or Dissimilarity

$$2^{\Omega_{sim}} = \{\emptyset, \{similar\}, \{dissimilar\}, \{similar, dissimilar\}\}$$

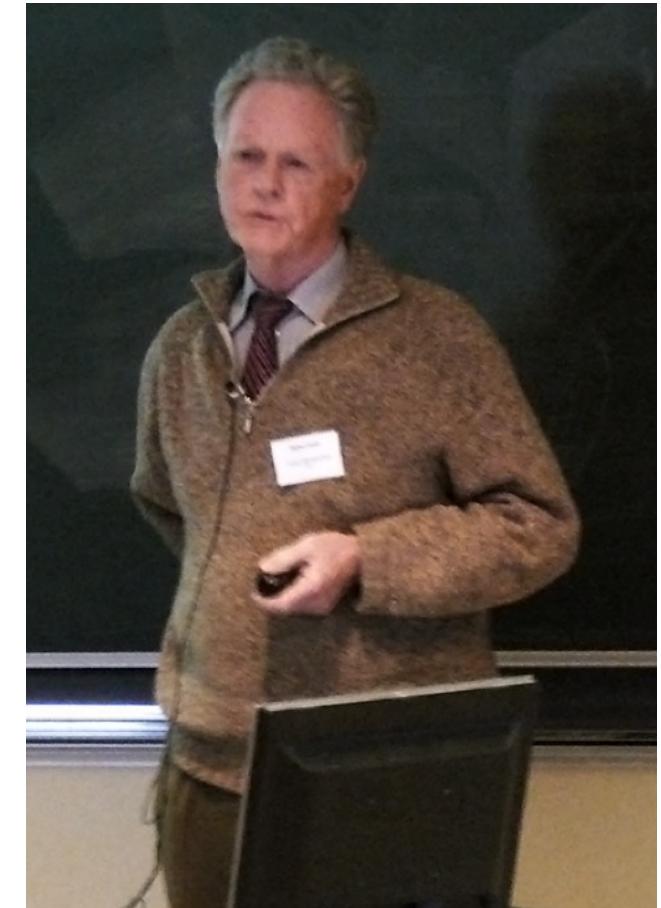
**Mass function**

$$m_{evidence}: 2^{\Omega_{sim}} \rightarrow \mathbb{R}$$

$m_{evidence}(\{similar\})$ : degree of belief that A and B are similar

$m_{evidence}(\{dissimilar\})$ : degree of belief that A and B are dissimilar

$m_{evidence}(\{similar, dissimilar\})$ : degree of belief that the situation is unknown



Prof. Arthur Pentland Dempster

# Challenges in Material Datasets: Size, Bias, and Conflicts

Summary of the 8 datasets of High Entropy Alloys

Datasets	No. alloys	Observation rate	No. HEAs	HEAs rate	No. candidates
$\mathcal{D}_{ASMI16}$	45 binary alloys	13%	45	(100%)	351 binary alloys
$\mathcal{D}_{CALPHAD}$	243 ternary alloys	9%	243	(100%)	2925 ternary alloys
$\mathcal{D}_{AFLLOW}$	117 binary alloys 441 ternary alloys	33% 15%	60 234	(51%) (53%)	351 binary alloys 2925 ternary alloys
$\mathcal{D}_{LTVC}$	117 binary alloys 441 ternary alloys	33% 15%	58 148	(49%) (33%)	351 binary alloys 2925 ternary alloys
$\mathcal{D}_{AFLLOW}^{quaternary}$	1,110 quaternary alloys	6%	754	(68%)	17,550 quaternary alloys
$\mathcal{D}_{AFLLOW}^{quaternary}$	1,110 quaternary alloys	6%	480	(43%)	17,550 quaternary alloys
$\mathcal{D}_{AFLLOW}^{quinary}$	130 quinary alloys	0.16%	129	(99%)	80,730 quinary alloys
$\mathcal{D}_{LTVC}^{quinary}$	130 quinary alloys	0.16%	91	(70%)	80,730 quinary alloys

"HEAs rate": the ratio of No. HEA to No. alloys; "Observation rate": the ratio of No. alloys to No. candidates.

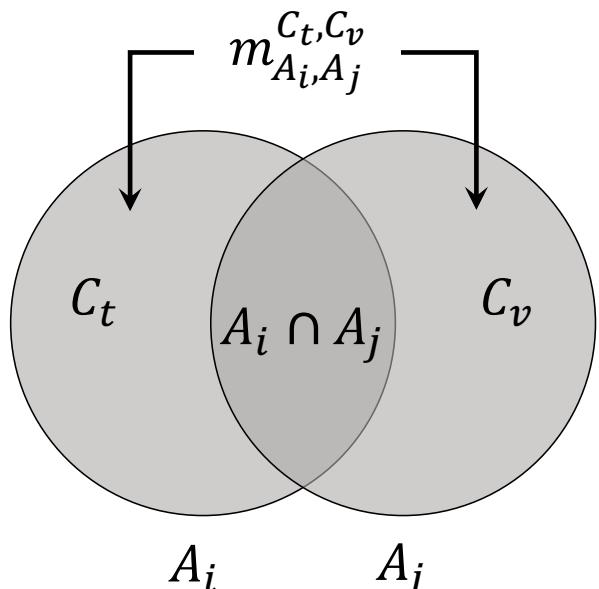
The number of confirmed materials is relatively smaller than the number of candidates

Heavily bias to positive label (HEA phase)

# Modeling Similarity Between Compositions

Quantifying uncertainty  
(Due to Limited or Conflicting Data)

Evidence of similarity



Collecting evidence from data

$$m_{A_i, A_j}^{C_t, C_v}(\{\text{similar}\}) = \begin{cases} \alpha & \text{if } y_{A_i} = y_{A_j}, \\ 0 & \text{otherwise} \end{cases},$$

property of alloy  $A_i$   
property of alloy  $A_j$

confidence of  
similarity  
(substitutable)

$$m_{A_i, A_j}^{C_t, C_v}(\{\text{dissimilar}\}) = \begin{cases} \alpha & \text{if } y_{A_i} \neq y_{A_j}, \\ 0 & \text{otherwise} \end{cases},$$

confidence of  
dissimilarity  
(not substitutable)

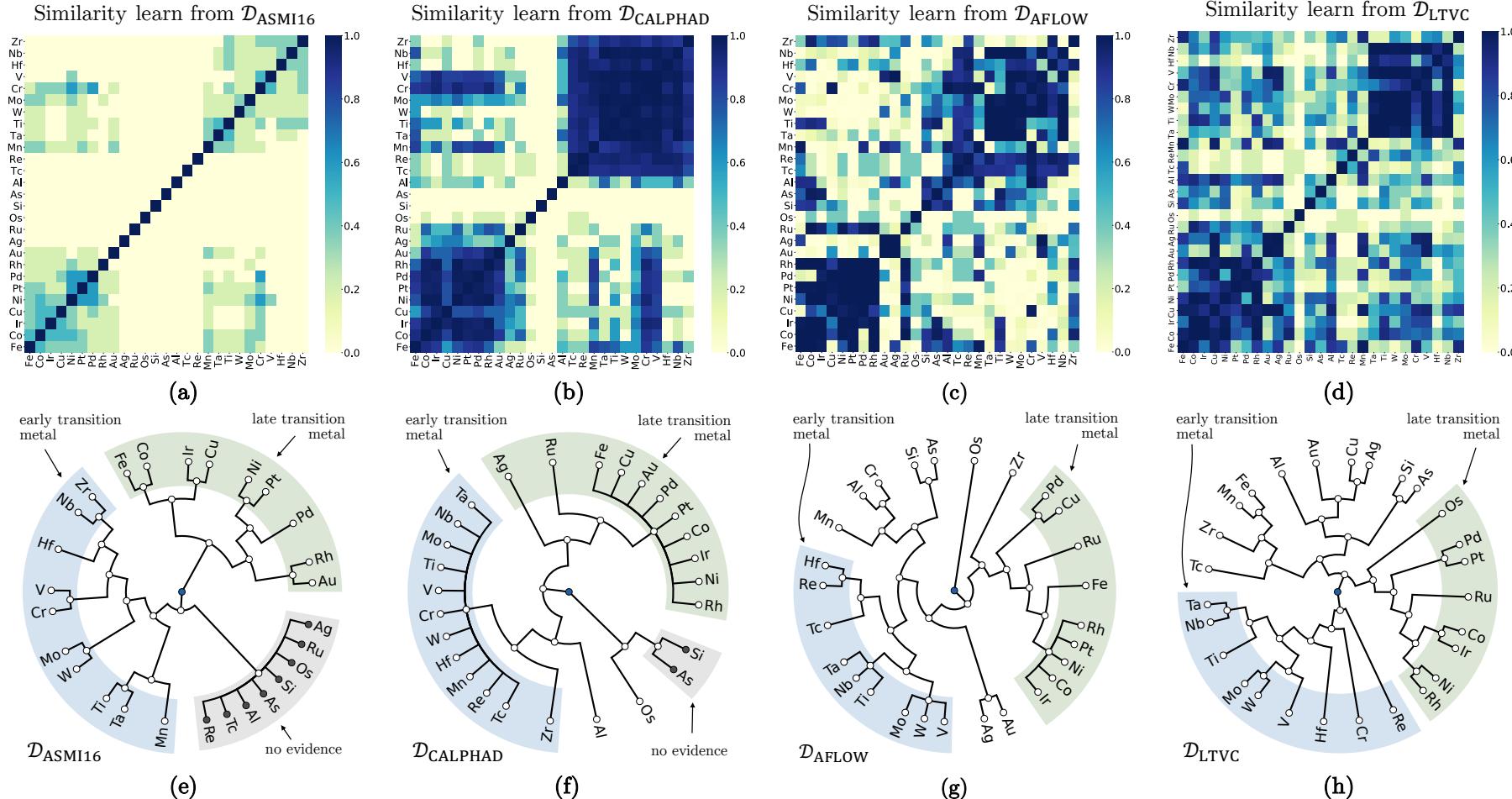
$$m_{A_i, A_j}^{C_t, C_v}(\{\text{similar, dissimilar}\}) = 1 - \alpha$$

degree of  
“unknown”  
(not sure if  
substitutable)

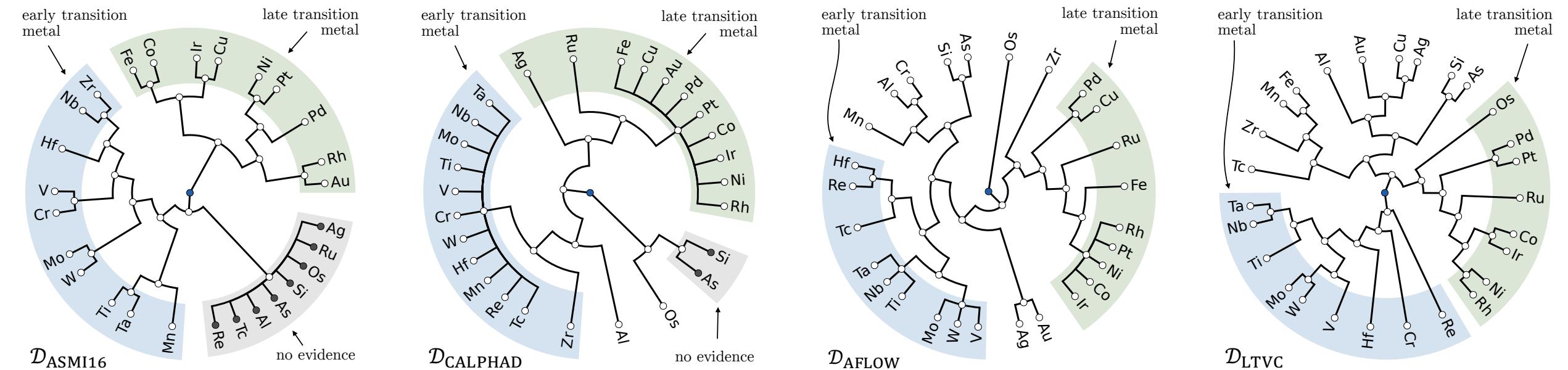
Application:

# Extracting Insights from Material Similarity to Support Decision Making

# Data-Driven Similarities in Forming HEAs



# Data-Driven Similarities in Forming HEAs



# Recommending New Element Combinations

confidence of forming HEA

confidence of not forming HEA

degree of “unknown”

To be substituted by other components to create new candidates

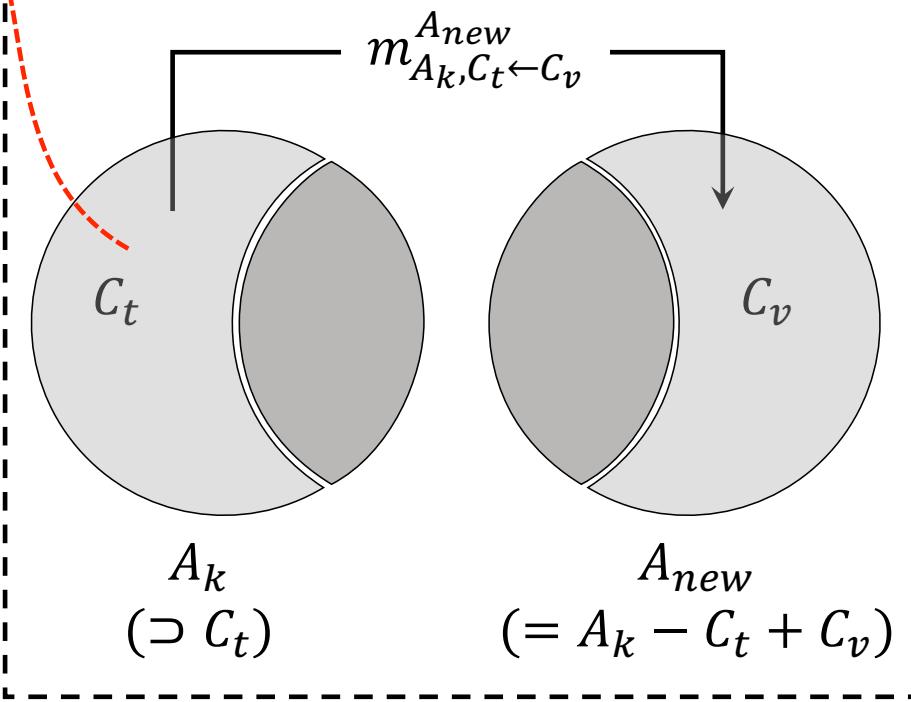
$$m_{A_k, C_t \leftarrow C_v}^{A_{new}}(\{HEA\}) = \begin{cases} s(C_t, C_v) & \text{if } y_{A_k} = HEA, \\ 0 & \text{otherwise} \end{cases}, \quad (6)$$

$$m_{A_k, C_t \leftarrow C_v}^{A_{new}}(\{\neg HEA\}) = \begin{cases} s(C_t, C_v) & \text{if } y_{A_k} = \neg HEA, \\ 0 & \text{otherwise} \end{cases}, \quad (7)$$

$$m_{A_k, C_t \leftarrow C_v}^{A_{new}}(\{HEA, \neg HEA\}) = 1 - s(C_t, C_v), \quad (8)$$

$C_t$  and  $C_v$  can differ in size!

Evidence of element substitution



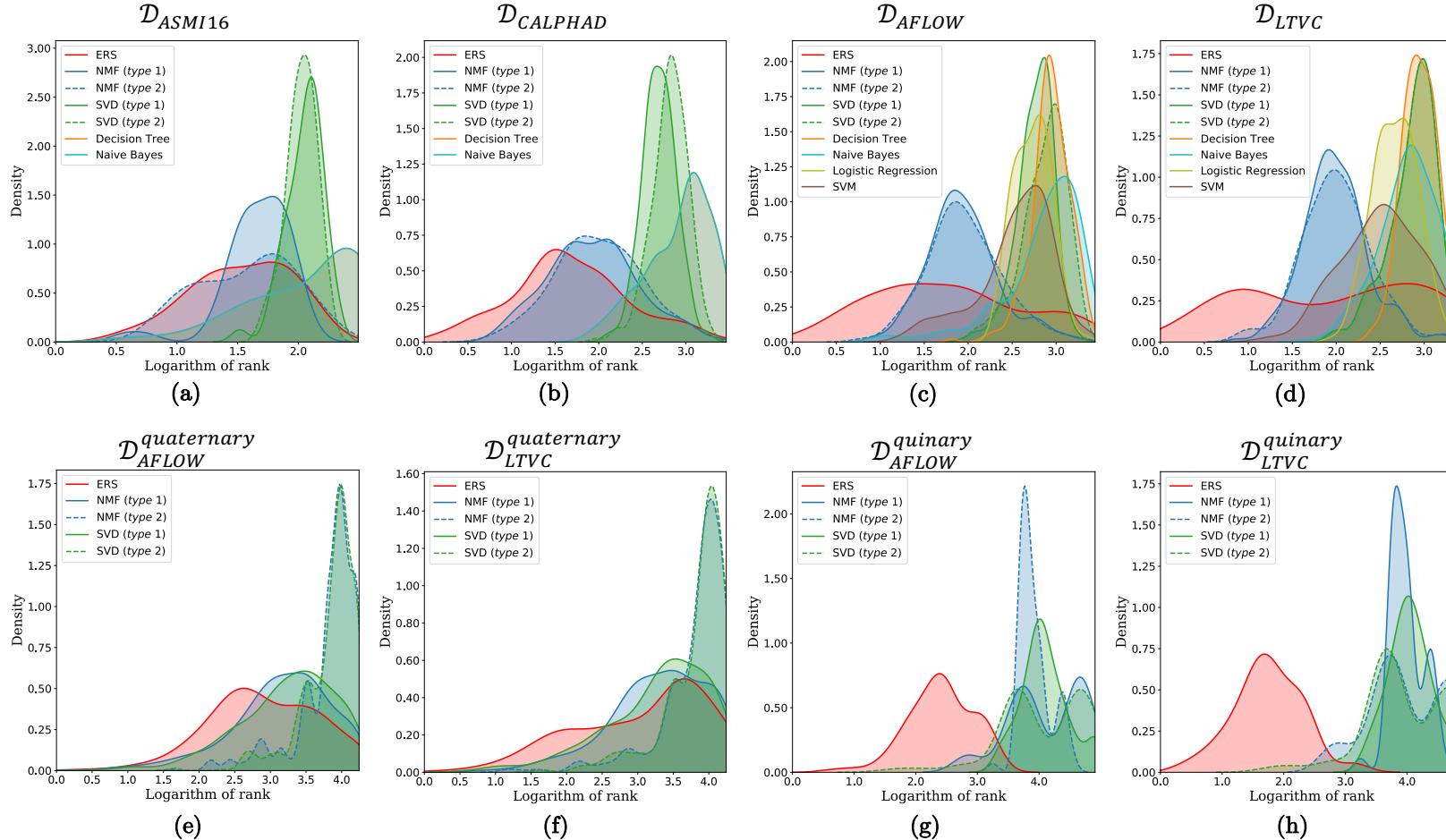
Recommending potential HEAs

# Recommendation capability

$C_t$  and  $C_v$  can differ in size!



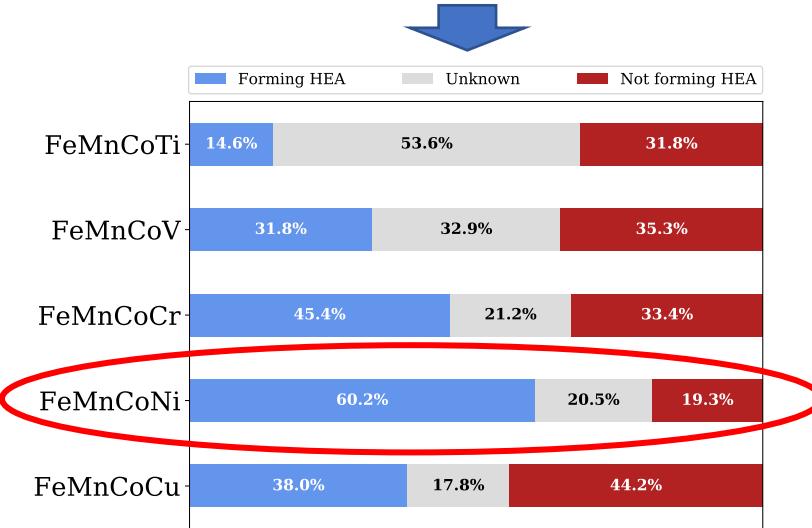
Recommending  
quaternary and  
quinary HEAs by using  
binary, ternary HEAs  
data



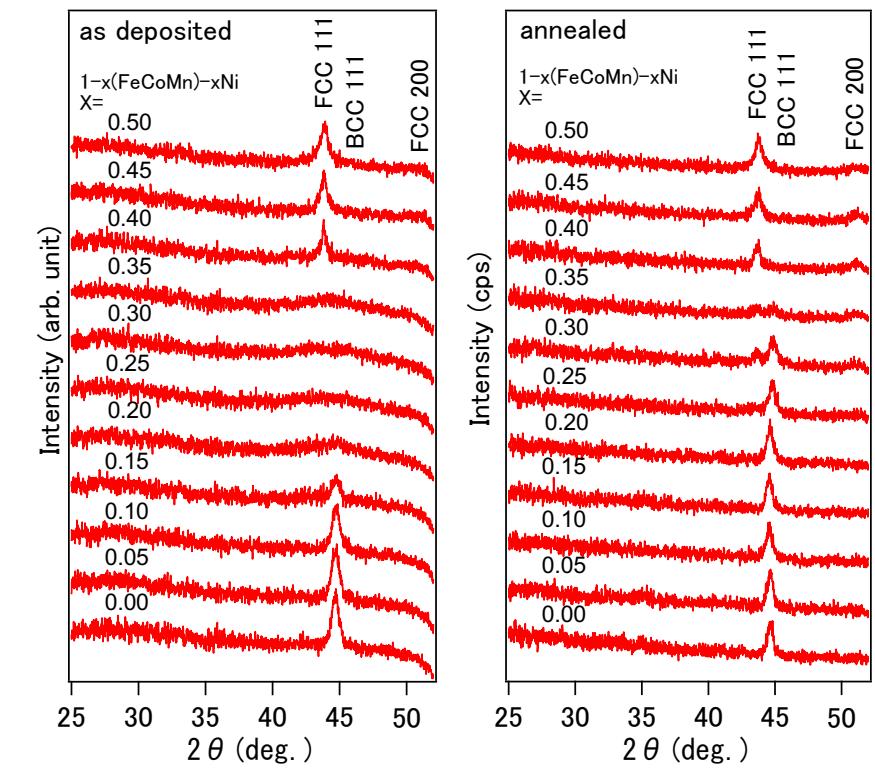
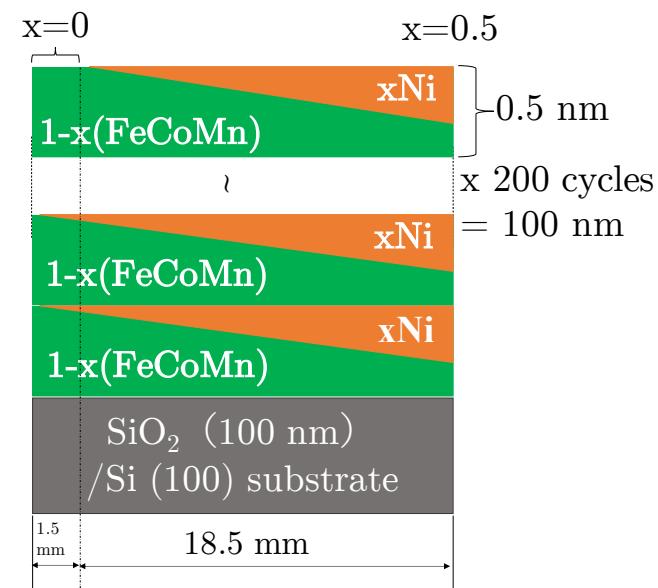
~ 500 times  
higher recall rate

# Experimental validation of FeMnCoNi

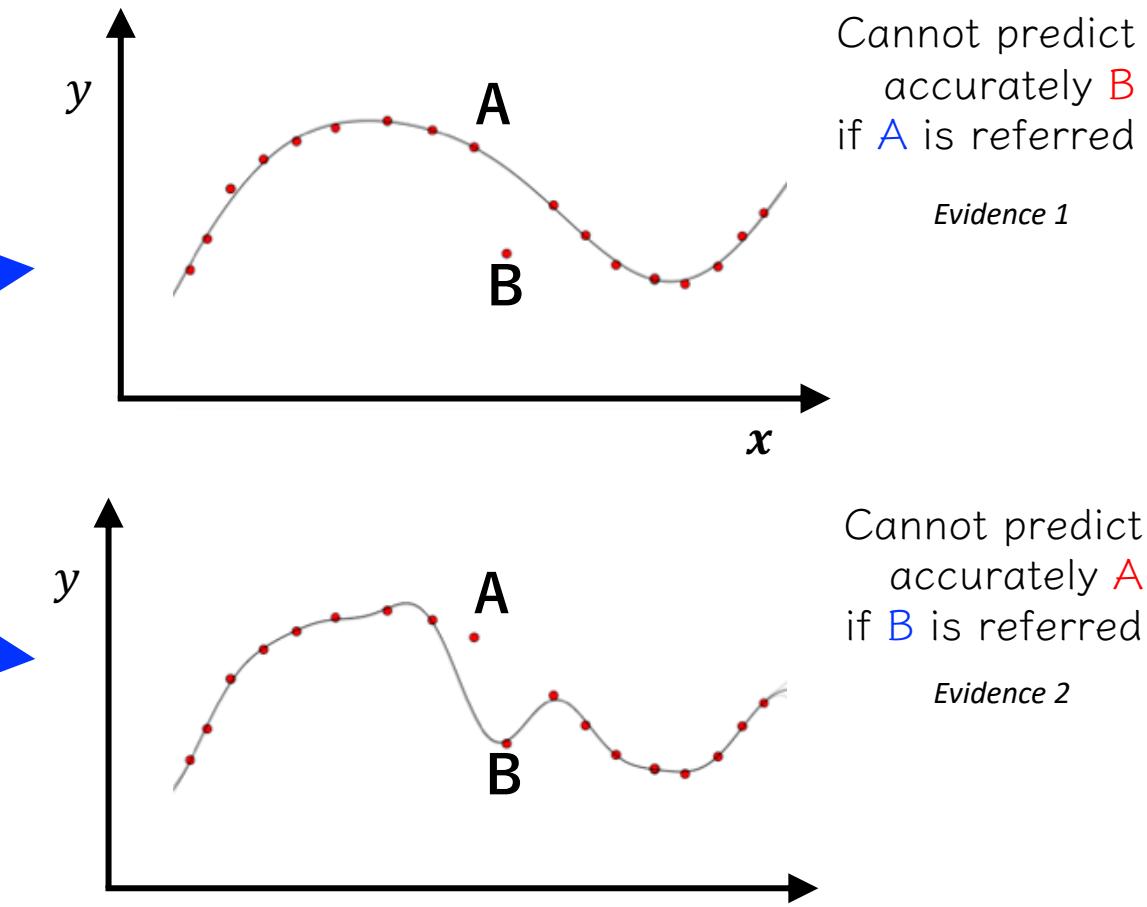
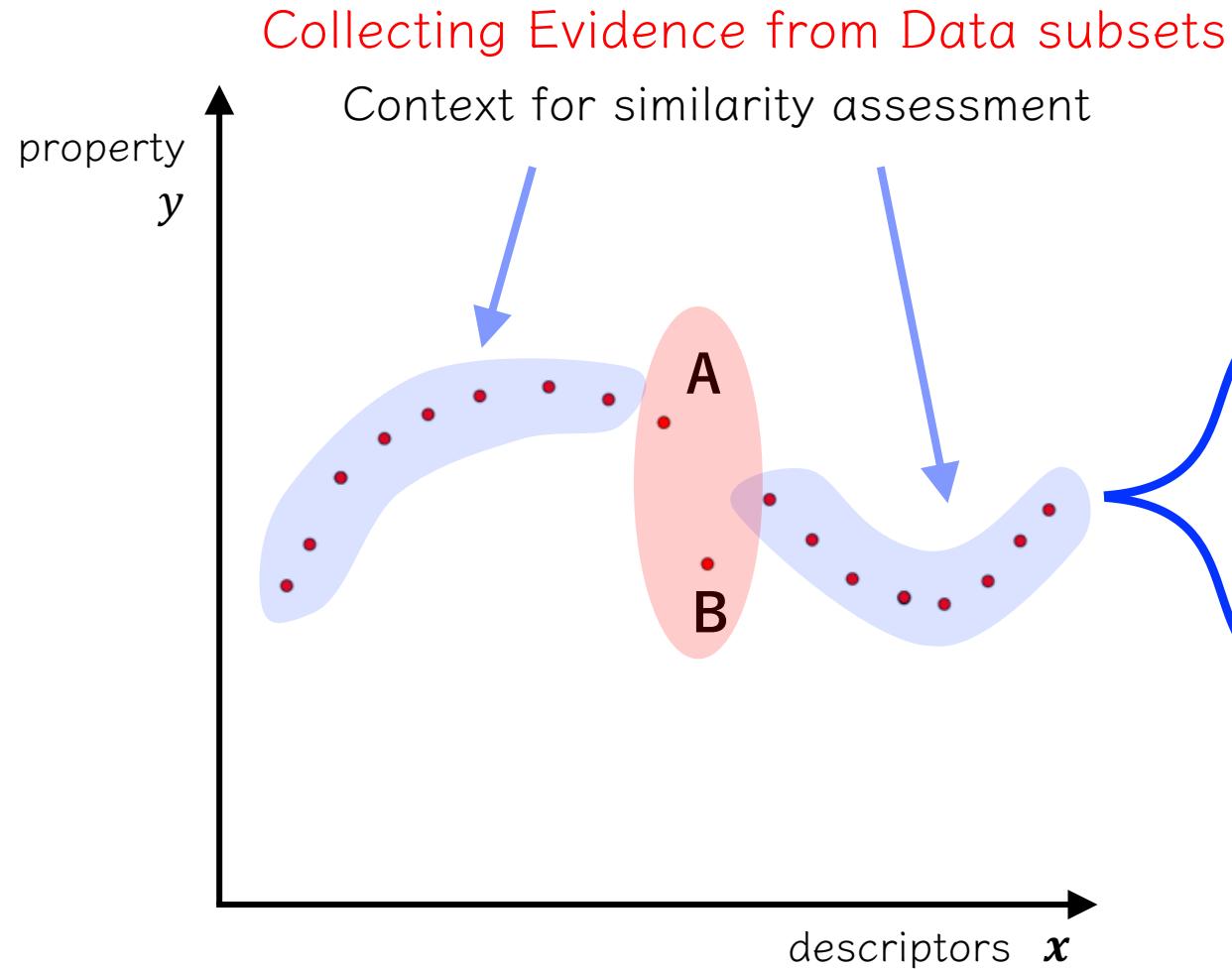
'Unknown' indicates insufficient information from observed data



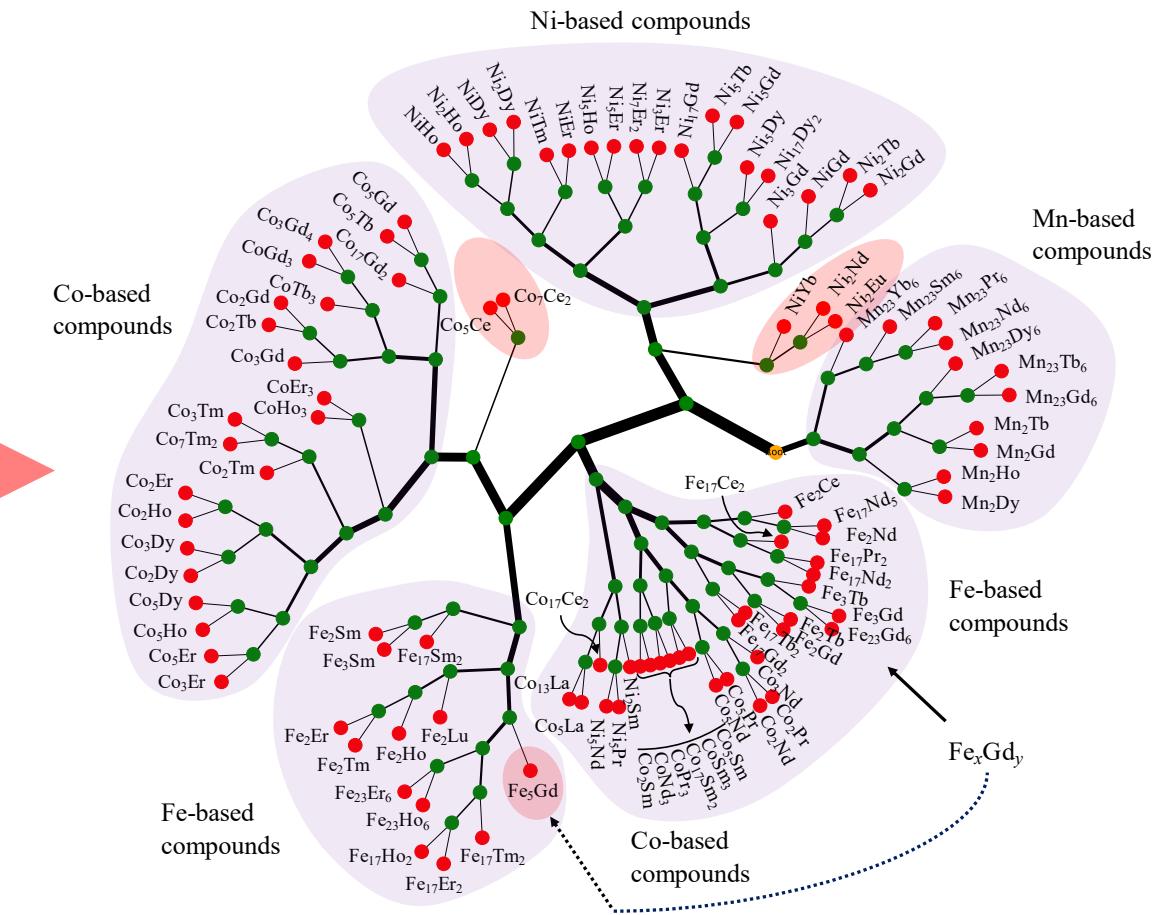
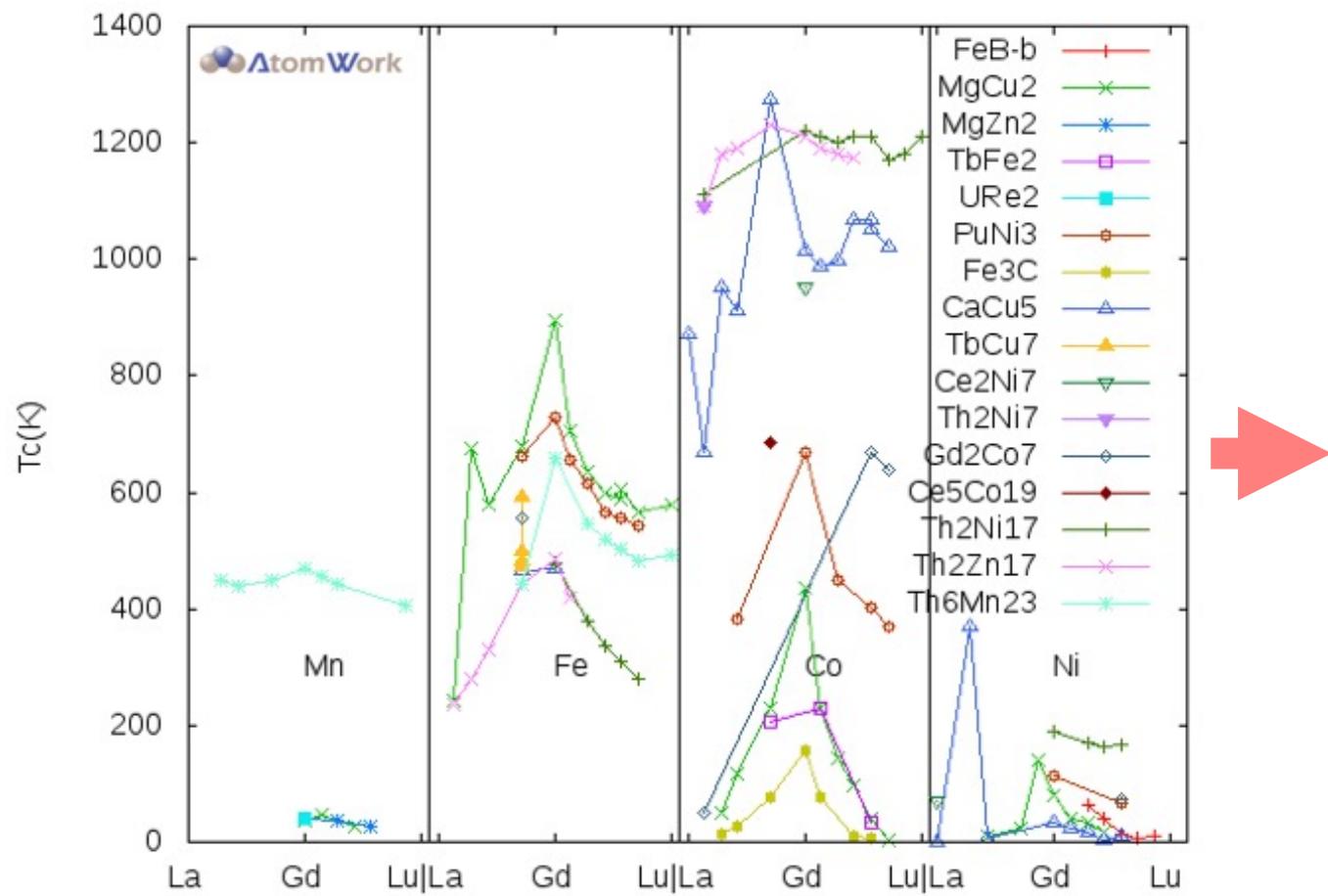
'Unknown' should not be interpreted as a 50/50 probability



# Similarity: Correlation between occurrence



# Similarity between Rare earth – Transition metal alloys



# Gleaning Insights from Material Similarity

