

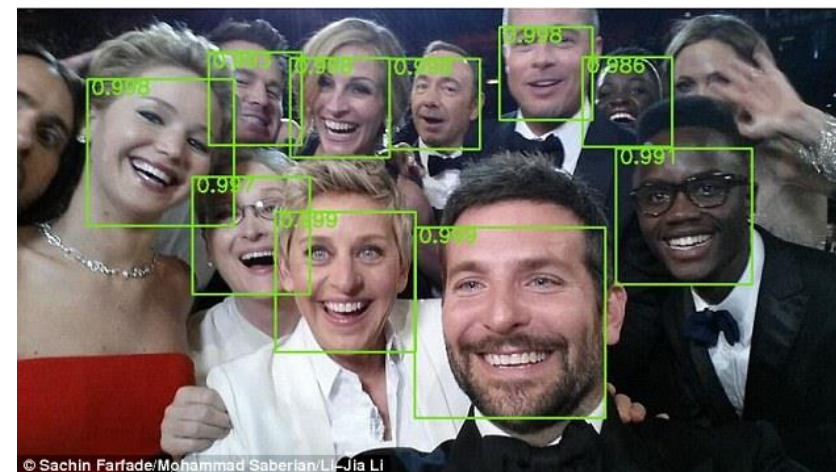
# *Python*ではじめる機械学習

WeekendEngineerコンペティションチーム輪講会 vol.1

データの集合から、  
みえているものからみえていないものを予測する手法



実例としては・・・



## Amazonのレコメンド機能

# 顔認識システム

## なぜ機械学習 (ML) なのか？



「SPAMフィルタ」を例に考えてみる

機械学習を使わない場合

= 人がルールを設計する場合

単語のブラックリストを作成し、  
その単語が  
出てきたらSPAM  
出てこなかったらSPAMではない

問題点

- ・ ルールの設計が必要  
→ 専門家が必要
- ・ タスクが少しでも変わると  
システム全体を  
書き直さなければならない
- ・ 人がルールを  
機械に教え続けるには限界がある

## なぜ機械学習 (ML) なのか？



「画像からの顔認識」を例に考えてみる

~~機械学習を使わない場合~~

= 人がルールを設計する場合

機械学習を使わないとできない！

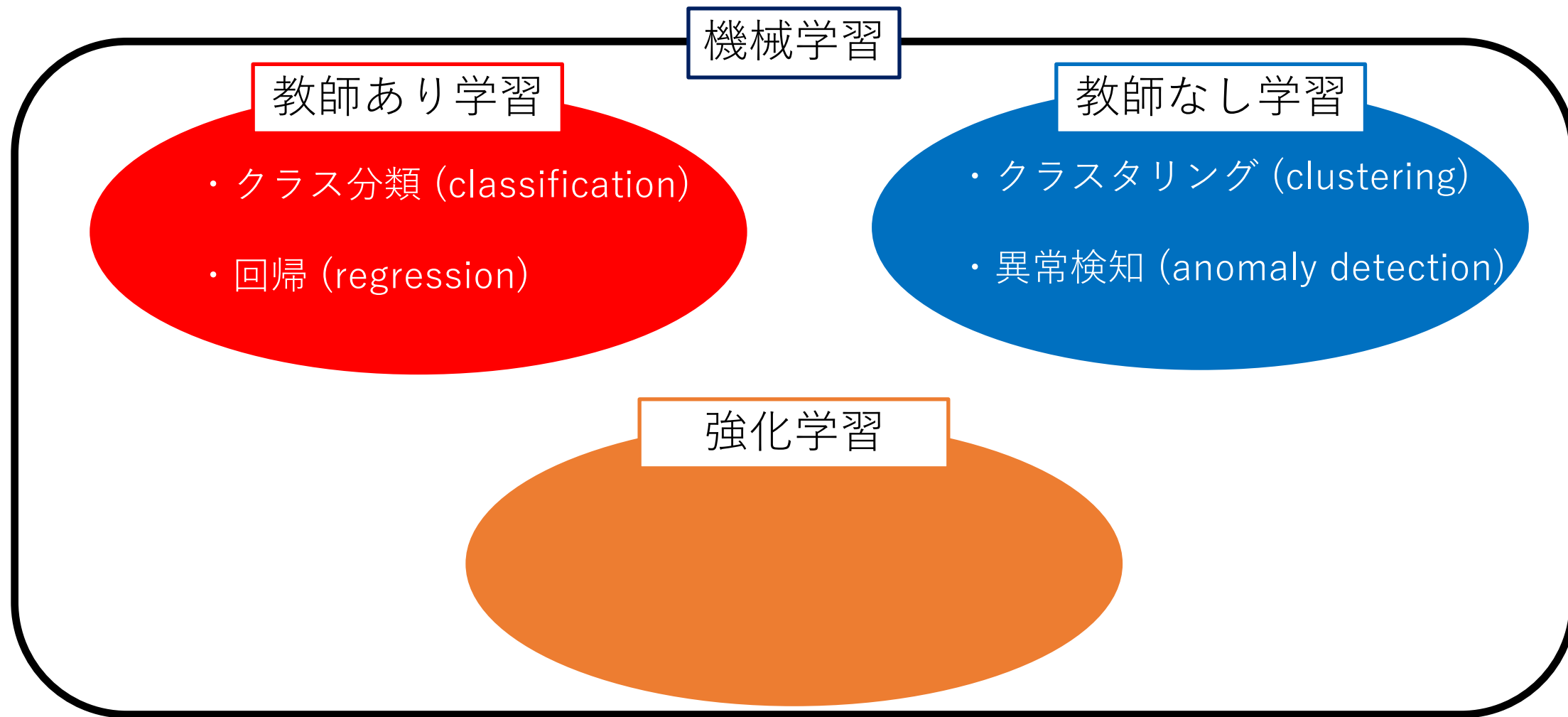
問題点

計算機がピクセルを「知覚」する方法が、  
人間が顔を認識する方法と全く異なる

→ 人がルールを設計出来ない

分類の仕方はイロイロ・・・

## 機械学習 (ML) の分類

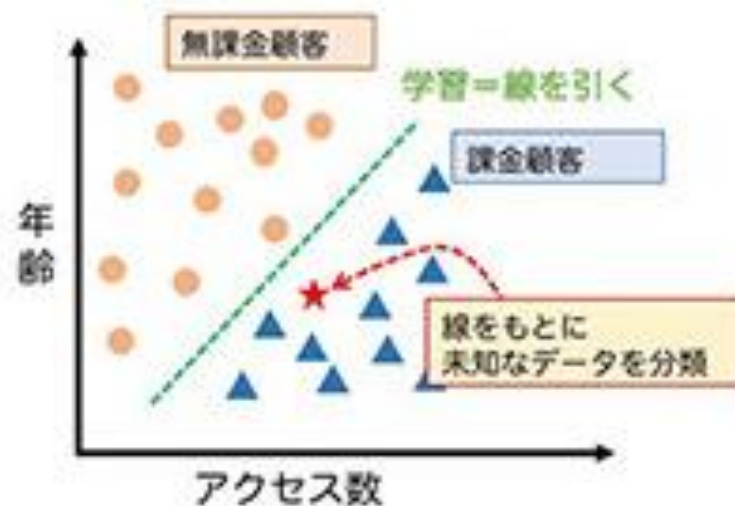


## 教師あり学習 (Supervised Learning) の概要

**学習データと教師データ (正解ラベル) を用いて学習する手法**

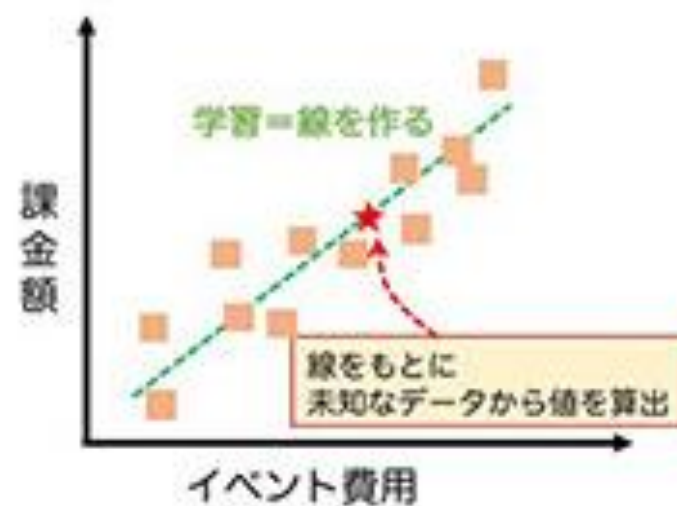
○クラス分類 (classification)

**クラスラベル**を予測する



○回帰 (regression)

**数値**を予測する



## 教師あり学習 (Supervised Learning) の概要

### 学習データと教師データ (正解ラベル) を用いて学習する手法

- 例) ・ KaggleのTitanicコンペティション  
→ 2クラス分類

Titanic号の乗客が事故で  
死んだか(0)生きたか(1)を予測

- ・ MNIST (Mixed National Institute of Standards and Technology database)  
→ 10クラス分類

手書き数字画像「0」～「9」を予測

- ・ iris (アヤメ) の分類  
→ 3クラス分類

花の画像「setosa」、「versicolor」、「virginica」を予測

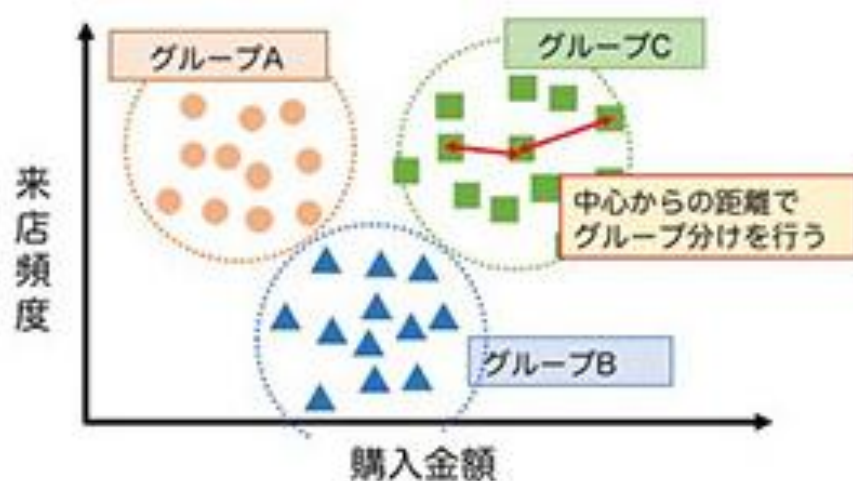
- ・ 売上予測  
→ 回帰

## 教師なし学習 (Unsupervised Learning) の概要

**学習データのみ**を用いて学習する手法

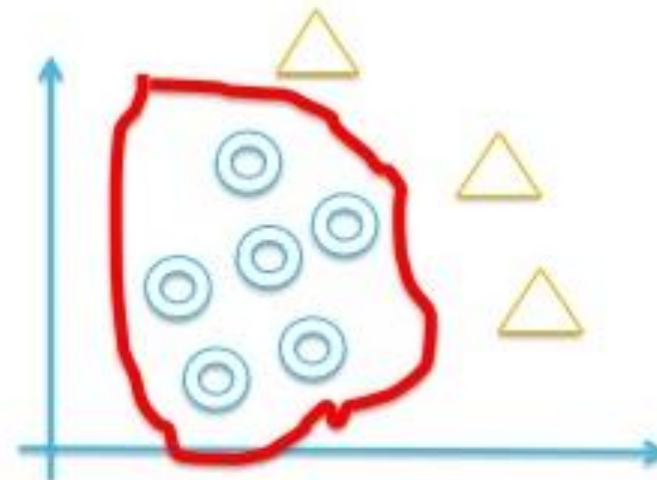
○クラスタリング (clustering)

グループ分けをする



○異常検知 (anomaly detection)

異常値を見つける





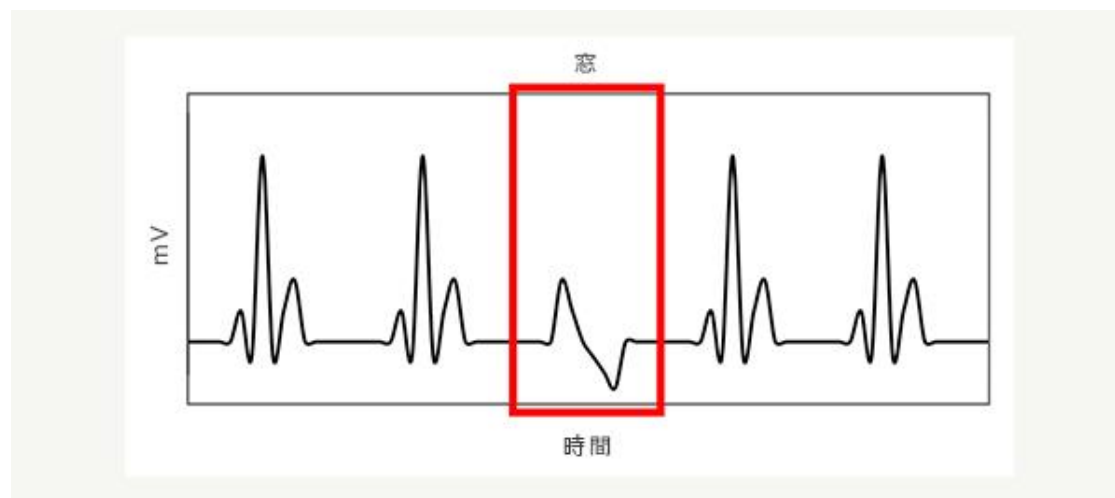
## 教師なし学習 (Unsupervised Learning) の概要

### 学習データのみを用いて学習する手法

例) ・レコメンド

おすすめの商品などを紹介

・心電図の異常検出



## タスクを知り、データを知る

大事

- ・ 扱っているデータを理解する
- ・ 解決しようとしている問題とデータとの関係を理解する

アルゴリズムには、  
それぞれ得意とする  
データの種類や  
問題の設定がある

適当にアルゴリズムを選んでデータを投げればいいわけではない

- ・ 解決しようとしている問題は何か？  
集めたデータでその問題が解決できるのか？
- ・ 機械学習が適しているのか？
- ・ 解決しようとしている問題を解くために十分なデータを集めたか？
- ・ どのような特徴量を抽出しただろうか？  
その特徴量で正しい予測が可能だろうか？
- ・ どのような評価を行えばよいか？

機械学習はあくまで手段であり、  
目的は問題を解決すること

## iris (アヤメ) のクラス分類

iris (アヤメ) の Sepal (がく片)、Petal (花弁) の幅、長さから  
「setosa」、「versicolor」、「virginica」の3品種に分類

setosa



versicolor



virginica



## iris (アヤメ) のクラス分類

①データの読み込み

sklearnのライブラリからインポート

②学習データとテストデータの分割

③データの検査

④学習モデルの構築

⑤モデルの評価

## iris (アヤメ) のクラス分類

詳細は  
5章で

①データの読み込み

②学習データとテストデータの分割

③データの検査

④学習モデルの構築

⑤モデルの評価

データはラベルでソートされているため  
シャッフルしないと出力ラベルが偏る

データセットをシャッフルしてから  
学習データ75%、テストデータ25%に分割

## iris (アヤメ) のクラス分類

①データの読み込み

データがうまく  
分離しているか  
らok

②学習データとテストデータの分割

③データの検査

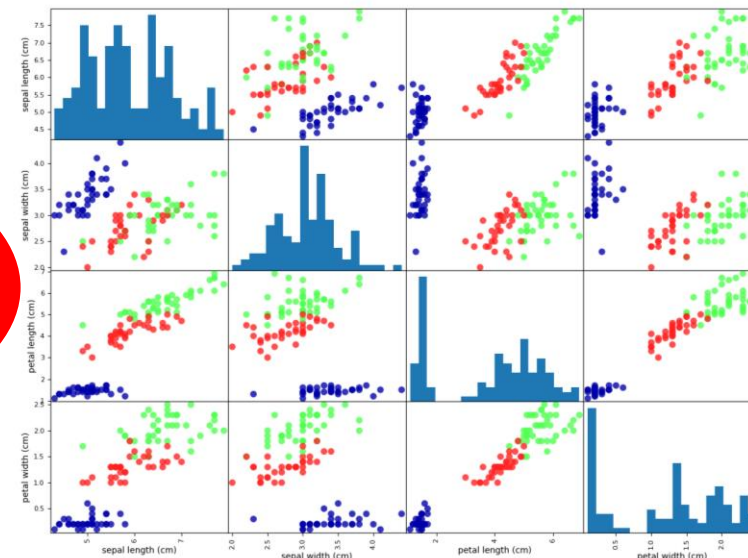
ペアプロットで可視化

④学習モデルの構築

全ての特徴量の組み合わせをプロットしたもの

⑤モデルの評価

全ての特徴量の相関を同時に見ているわけではないので、  
特徴量の数が少ない場合はうまくいく



## iris (アヤメ) のクラス分類

①データの読み込み

②学習データとテストデータの分割

③データの検査

④学習モデルの構築

⑤モデルの評価

学習データに近いデータに合わせる手法  
(詳細は2章で)

k-近傍法 (k-Nearest Neighbors: KNN) を採用

ここでは  $k=1$

## iris (アヤメ) のクラス分類

- ① データの読み込み
- ② 学習データとテストデータの分割
- ③ データの検査
- ④ 学習モデルの構築
- ⑤ モデルの評価



テストデータを用いて評価する

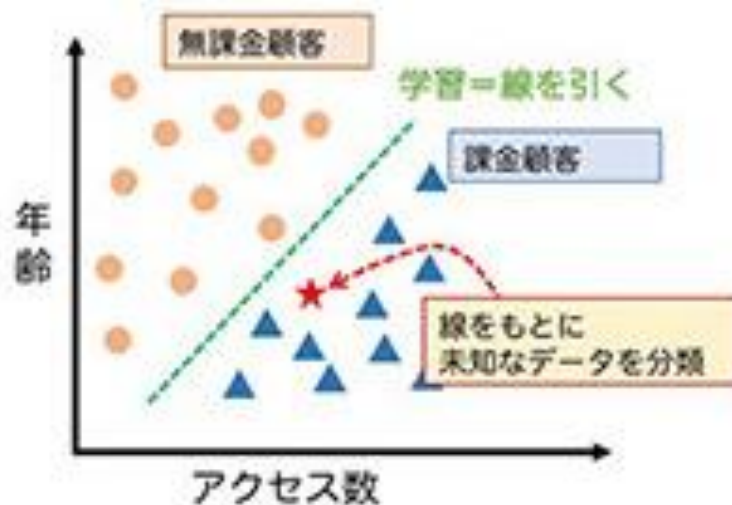


### 教師あり学習 (Supervised Learning) の概要

**学習データと教師データ (正解ラベル) を用いて学習する手法**

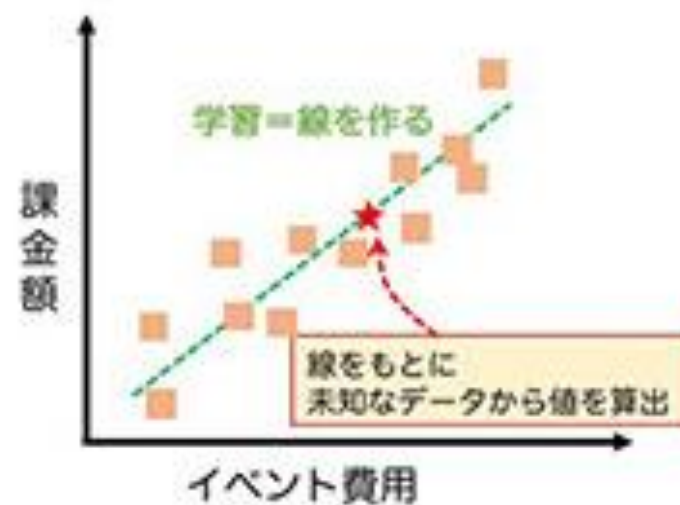
○ クラス分類 (classification)

**クラスラベル**を予測する



○ 回帰 (regression)

**数値**を予測する



### 汎化性能 (generalization performance)

**学習データ以外のデータに対しても正しい予測ができる能力**

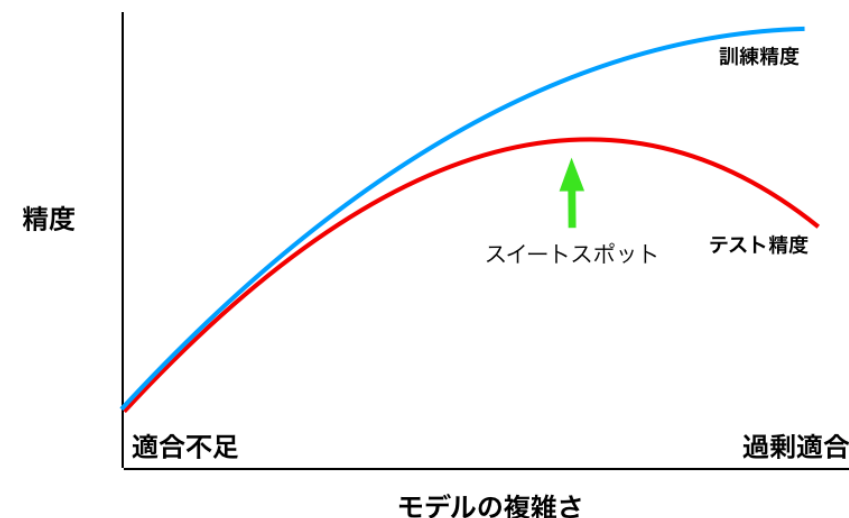
過学習 (over fitting) (過剰適合) ↔ 適合不足 (under fitting)

訓練データにすらうまく学習ができないこと

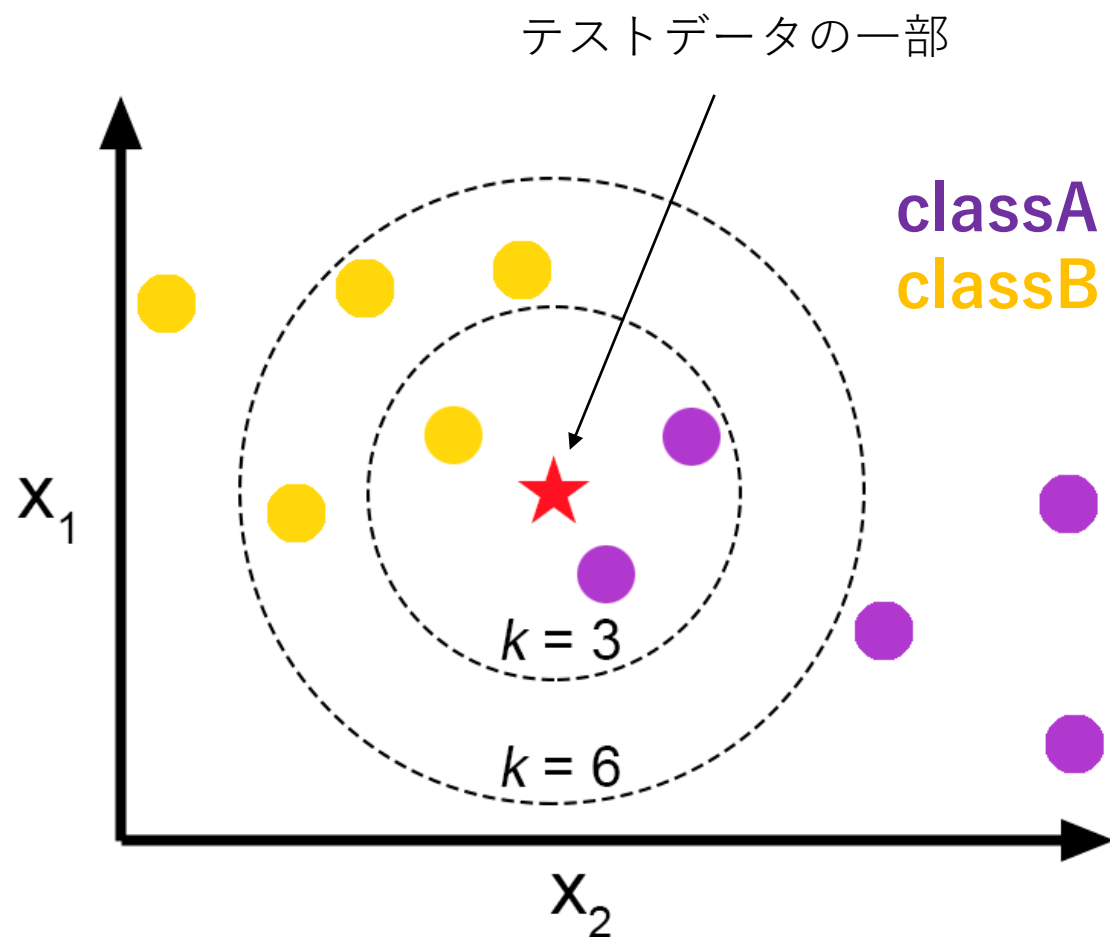
**特定の学習データのみに特化した学習をしてしまうこと**

精度とモデルの複雑さ

最良の汎化性能を示す  
スイートスポットが  
理想のモデルである



### k近傍法(k Nearest Neighbor: kNN)



クラス分類だけでなく、  
回帰にも使える

学習データをベクトル空間上にプロットしておき、  
未知のデータが得られたら、  
そこから距離が近い順に任意のk個を取得し、  
多数決でデータが属するクラスを推定する

### k近傍法(k Nearest Neighbor: kNN)

#### 利点

- ・モデルが理解しやすいため、モデル構築が容易
- ・パラメータ調整による影響が比較的小さい

#### 欠点

- ・多数の特徴量を持つデータセットでは十分な精度がでない
- ・前処理による影響が比較的大きい
- ・データ個数が多いほど処理速度が遅い(計算時間がかかる)



欠点より、  
実際にはほとんど使われていない

## 2章：教師あり学習 (Supervised Learning)

P46,  
47

対応  
ページ

### 線形モデル (linear model) による回帰

$$y = \sum_{i=1}^n (w_i x_i) + b = w_1 x_1 + w_2 x_2 \cdots + w_n x_n + b$$

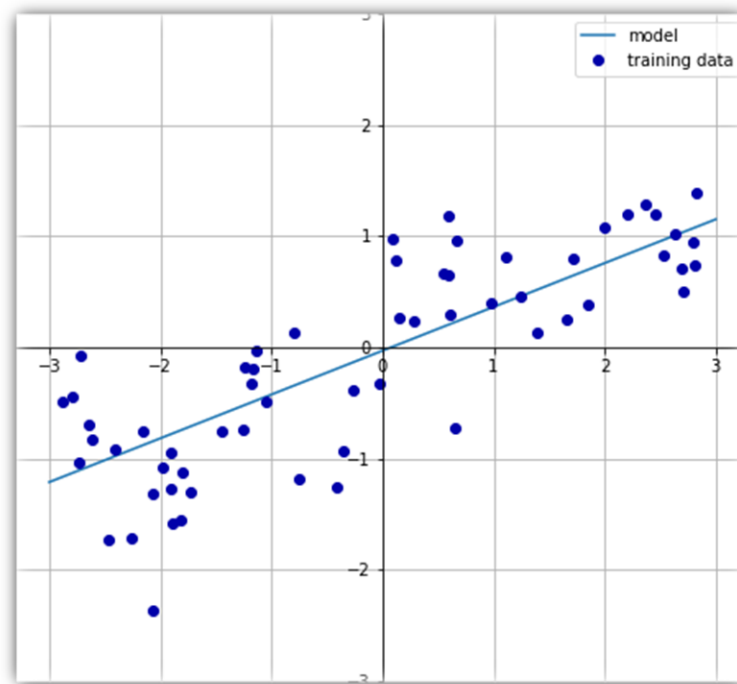
重回帰

$$y = w * x + b$$

単回帰

x : 特徴量  
w : xに対する重み  
b : バイアス  
y : 予測値

低次元空間では、線形モデルを用いるのは制約が強すぎるかもしれない



### 線形回帰 (最小二乗法 (ordinary least squares: OLS))

学習データにおいて、

損失関数 (loss function)

予測値と真値の平均二乗誤差 (Mean Squared Error: MSE)  
を最小にするために**重み  $w$  とバイアス  $b$** を求める

誤差関数

$$MSE(y_i, \hat{y}_i) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$y_i$  : 正解値     $\hat{y}_i$  : 予測値

### リッジ回帰 (Ridge Regression)

基本は線形回帰と同じだが、  
損失関数に、L2正則化項を加える

正則化とは？

過学習を防ぐために  
明示的にモデルを制約すること

過学習を防ぐ

個々の特徴量が出力に与える影響をなるべく小さくしたい  
→ 重み $w$ の大きさが高いものにペナルティを与える

誤差関数

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^n W_i^2$$

$y_i$  : 正解値     $\hat{y}_i$  : 予測値     $\lambda$ (ラムダ) : ハイパーパラメータ

### リッジ回帰 (Ridge Regression)

期待通り

- ・ 線形回帰と比べ、  
学習データに対するスコア → 低い  
テストデータに対するスコア → 高い
- ・ 汎化性能が高い
- ・ 過学習(*over fitting*)の危険が少ない
- ・ 学習データに対するスコアとモデルの簡潔さ(0に近い $w$ の数)はトレードオフの関係



### Lasso

基本は線形回帰と同じだが、  
損失関数に、L1正則化項を加える

いくつかの重みが完全に0になる  
→ 重みが0になった特徴量は無視される

誤差関数

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^n |w_i|$$

$y_i$  : 正解値     $\hat{y}_i$  : 予測値     $\lambda$ (ラムダ) : ハイパーパラメータ

### Lasso

- ・  $\lambda = 1$  の場合、ほとんどの  $w$  が 0 になってしまい、適合不足となった。
- ・  $\lambda$  を小さくしすぎると線形回帰と同じような結果になる
- ・  $\lambda$  を適度に小さくすると、学習に使う特徴量を減らすことができ、シンプルなモデルになる  
→ 特徴量がたくさんあり、そのうち重要なものはわずかしかなことが予測できる場合に有効

### クラス分類のための線形モデル (linear model)

2クラス分類

$$y = \sum_{i=1}^n (w_i x_i) + b = w_1 x_1 + w_2 x_2 \cdots + w_n x_n + b > 0$$

yの値が0より大きければクラスは+1
0より小さければ - 1

- ロジスティック回帰 (logistic regression)
- 線形サポートベクターマシン (Linear Support vector Machine)

### クラス分類のための線形モデル (linear model)

#### 多クラス分類

基本的に2クラス問題にしか適用できない

→各クラスに対して、そのクラスと他のすべてのクラスを分類する  
2クラス分類を用いる

$$y = \sum_{i=1}^n (w_i x_i) + b = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

クラスごとに  
重み $w$ とバイアス $b$ があるからそれを用いて  
クラススコアを計算

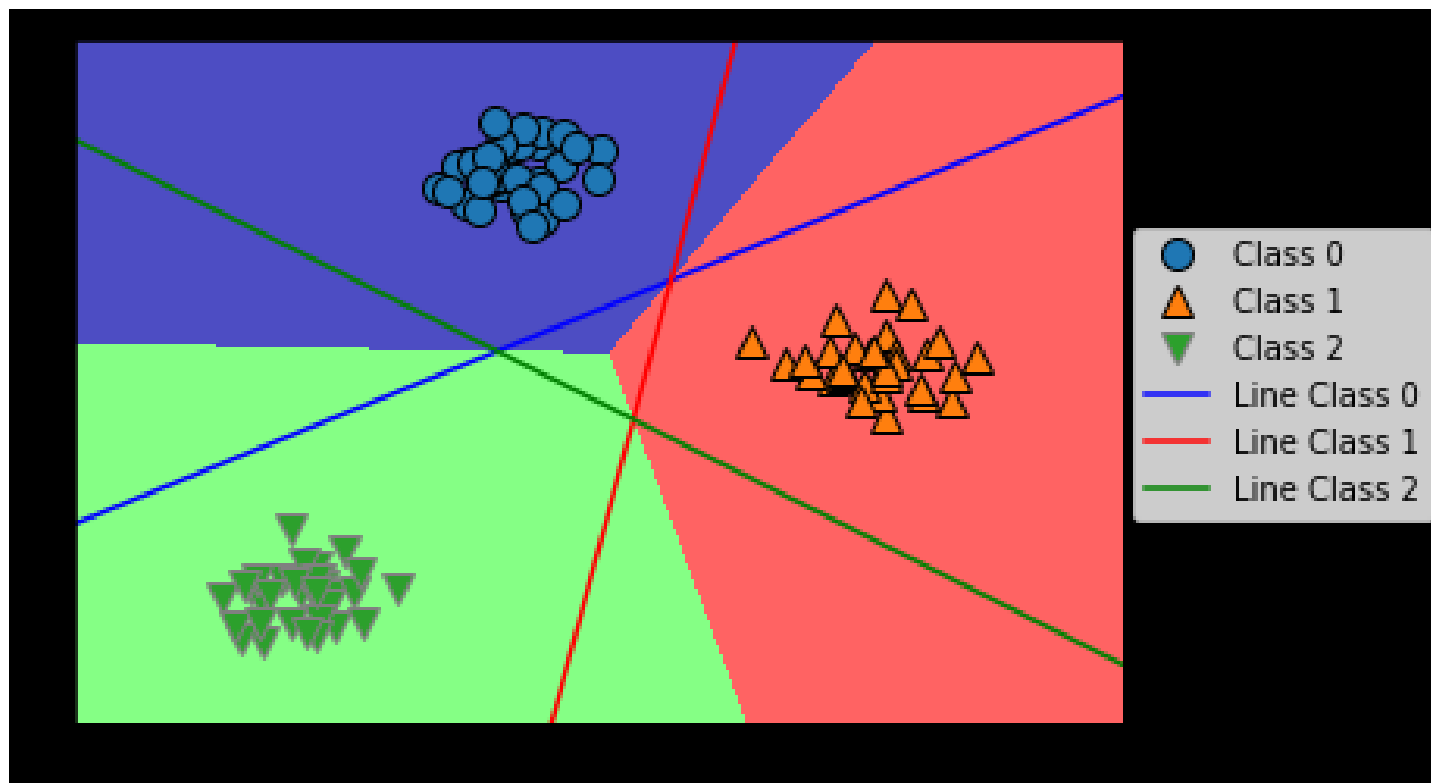
## 2章：教師あり学習 (Supervised Learning)

P63-  
67

対応  
ページ

### クラス分類のための線形モデル (linear model)

多クラス分類



### ナイーブベイズ

クラス分類にしか適用できない  
基本的には線形回帰と同じ

#### 特徴

- ・ 学習が線形回帰よりも高速
- ・ クラス分類でしか使えない
- ・ 線形回帰より精度が低い

scikit-learnで実装されているナイーブベイズ分類器

- ・ GaussianNB  
➡ 任意の連続値データに適用
- ・ BernoulliNB  
➡ 2値データに適用  
例) 0, 1のみで構成された特徴量
- ・ MultinomialNB  
➡ カウントデータ(整数値)  
例) 文中に出てくる単語の出現数

## ここまでのまとめ

とりあえず最初に試すべきアルゴリズム

### 線形モデル

- 処理速度が速い(計算時間がかからない)
- パラメータ調整によって過学習などの問題を解決できる
- 大きいデータや高次元のデータに適する
- パラメータ調節による影響が大きい場合があるためパラメータ調節に時間(コスト)がかかる

### kNN

- モデルが理解しやすく、モデル構築が容易
- パラメータ調整による影響が比較的小さい
- 多数の特徴量を持つデータセットでは十分な精度がでない
- データ個数が多いほど処理速度が遅い(計算時間がかかる)
- 前処理による影響が比較的大きい

データがうまく分離されていたら試すべき

### ナイーブベイズ

- 線形モデルよりさらに処理速度が速い(計算時間がかからない)
- 線形モデルよりさらに高次元、大きいデータに適する
- 線形モデルより精度が劣ることが多い