

機械学習 第6回 回帰

立命館大学 情報理工学部

福森 隆寛

Beyond Borders

講義スケジュール

(第1～4回、第14回) (第5～13回、第15回)

□ 担当教員：村上 陽平先生・福森 隆寛

1	機械学習とは、機械学習の分類
2	機械学習の基本的な手順
3	識別（１）
4	識別（２）
5	識別（３）
6	回帰
7	サポートベクトルマシン
8	ニューラルネットワーク

9	深層学習
10	アンサンブル学習
11	モデル推定
12	パターンマイニング
13	系列データの識別
14	強化学習
15	半教師あり学習

□ 担当教員：叶 昕辰先生（第16回の講義を担当）

今回の講義内容

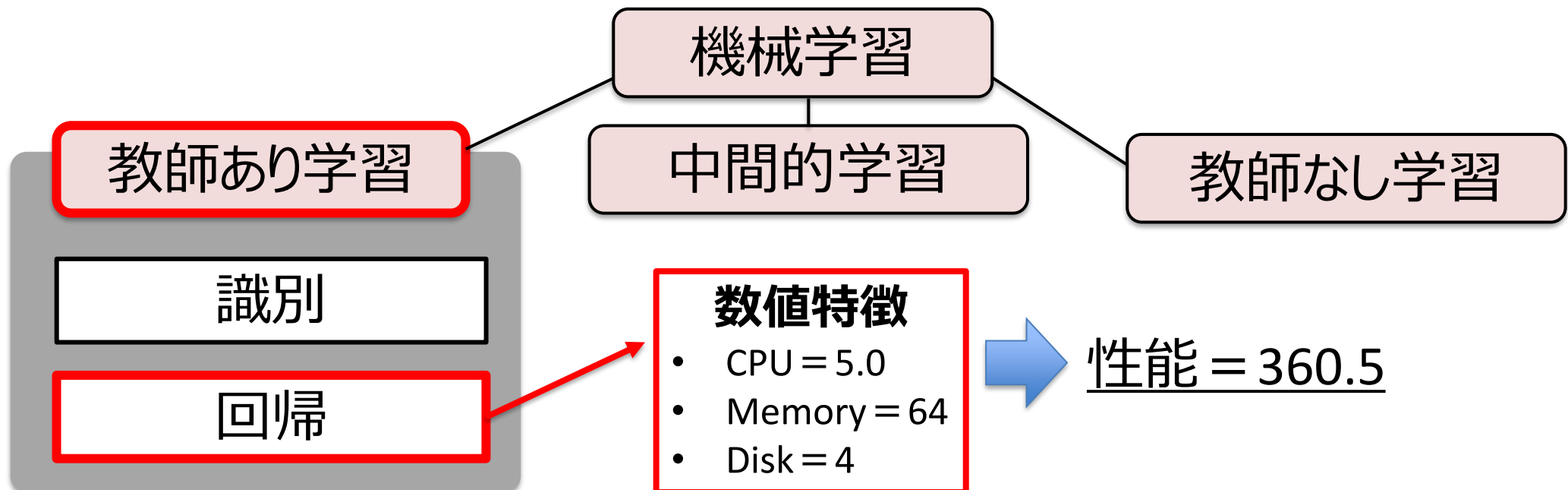
- 取り扱う問題の定義
- 線形回帰
- 回帰モデルの評価
せ い そ く か
- 正則化
- バイアスと分散のトレードオフ
か い き ぎ
- 回帰木
- 演習問題

取り扱う問題の定義：教師あり・回帰

□ **数値**データからなる特徴ベクトルを入力して、**数値**を出力する関数を作る

※ 教師あり学習の回帰問題での学習データは、以下のペアで構成される

入力データの特徴ベクトル $\leftarrow \{\underline{x_i}, \underline{y_i}\}, \quad i = 1, 2, \dots, \underline{N} \longrightarrow$ 学習データの総数
(数値データ) 数値形式の正解情報 \rightarrow 「ターゲット」と呼ぶ



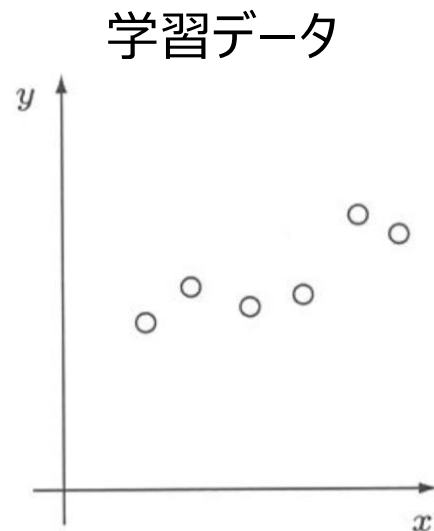
取り扱う問題の定義：教師あり・回帰

- 識別と回帰の境界は、それほど^{めいかく}明確ではない
 - 識別：数値特徴を入力としてクラスを出力
 - 回帰：数値特徴を入力として数値を出力

- クラスによって異なる値をとるクラス変数を^{どうにゆう}導入し
入力からクラス変数の値を予測する問題を考えると
識別問題を回帰問題として考えることもできる

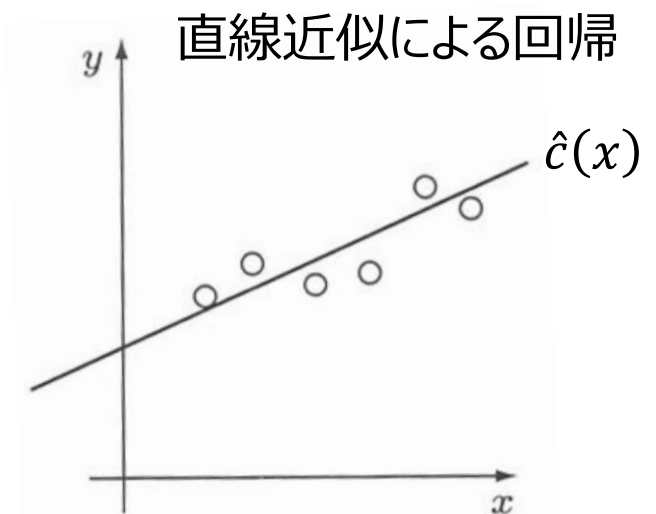
線形回帰

- 最も単純な入力も出力もスカラーである場合の回帰問題を考える
- 学習データから入力 x を出力 y に写像する関数 $\hat{c}(x)$ を推定



入力 x が大きくなると
出力 y も大きい値になる
傾向がみられる

この傾向を直線で表して
入力 x と出力 y を関係づける



最小二乗法と同様の方法で
なるべく誤差の少ない直線を求める

線形回帰

□ 最小二乗法から回帰式を求める

■ 回帰式を $\hat{c}(x) = w_1x + w_0$ とする

■ 誤差の二乗和は

$$E(\mathbf{w}) = \sum_{i=1}^N \{y_i - \hat{c}(x_i)\}^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

- \mathbf{X} : 1列目の全要素が1、2列目 i 行の要素が x_i のパターン行列
- \mathbf{w} : 重みベクトル ※ $\mathbf{w} = (w_0, w_1)^T$

■ \mathbf{w} で微分したものを0とすると、線形回帰式の重みは下式で計算できる

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

■ 入力 x が d 次元の場合も、同様の方法で計算できる

回帰モデルの評価

□ 回帰式が未知データに対して正しく出力値を予測するかを評価

□ 評価指標

■ 相関係数 R : 正解と予測が、どの程度似ているのか^に

■ 決定係数 R^2

- 「正解との離れ具合^{はな}」と「平均との離れ具合」の比を1から引く
- \tilde{y} : y_i の平均値

$$R^2 = 1 - \frac{\sum_{i=1}^N \{y_i - \hat{c}(x_i)\}^2}{\sum_{i=1}^N (y_i - \tilde{y})^2}$$

^{へんけい}
式変形により相関係数の二乗と一致するので R^2 と表記する

演習問題6-1（10分間）

□ 右表のような身長と体重のデータが与えられた

□ 身長 x から体重 $\hat{c}(x)$ を予測する線形回帰式が

$$\hat{c}(x) = 0.625x - 48.604$$

であるときの決定係数と相関係数を計算せよ

身長と体重データ

番号	身長 [cm]	体重 [kg]
1	147.9	41.7
2	163.5	60.2
3	159.8	47.0
4	155.1	53.2
5	163.3	48.3
6	158.7	55.2
7	172.0	58.5
8	161.2	49.0
9	153.9	46.7
10	161.6	52.5

正則化

□ 望ましい線形回帰式

- 汎化能力という点では、入力が少し変化したときに、出力も少し変化する回帰式が良い

- 重みが大きいと、入力が少し変化するだけで出力が大きく変化
- そのような回帰式は、たまたま学習データの近くを通っても、未知データに対する出力は信用できない
しんよう

- 線形回帰の重みは、値が0となる次元を多くすれば良い

- 回帰式の係数 w に関して、「大きな値の重みが、なるべく少なくなる」あるいは「0となる重みが多くなる」ような方法が必要

- このような工夫が**正則化**

- 誤差関数の式に正則化項を追加する
ついか

正則化 : Ridge回帰

□ Ridge回帰

- パラメータ \mathbf{w} の二乗を正則化項とする
- パラメータの値が小さくなるように正則化させる

$$E(\mathbf{w}) = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \underbrace{\lambda \mathbf{w}^T \mathbf{w}}_{\text{正則化項}}$$

λ : 正則化項の重み (重みが大きければ、性能よりも正則化の結果を重視)

- 最小二乗法でパラメータを求めたときと同様に、 \mathbf{w} で微分した値が0となる \mathbf{w} の値を求めると...

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad \text{※ } \mathbf{I} : \text{単位行列}$$

正則化：Lasso回帰

□ Lasso回帰

- パラメータ \mathbf{w} の絶対値を正則化項とする
- 値を0とするパラメータが多くなるように正則化される

$$E(\mathbf{w}) = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \sum_{j=1}^d |w_j|$$

λ ：正則化項の重み（大きければ、値を0とする重みが多くなる）

w_0 ：回帰式の切片は汎化能力に影響なし（通常は正則化の対象としない）
せつぺん えいきょう

- Lasso回帰の解は、解析的に求められない
 - 原点で微分不可能な絶対値を含むため
 - 正則化項の上限を微分可能な2次関数で押さえ、その関数のパラメータを誤差が小さくなるように逐次更新する方法が提案
ちくじ ていあん

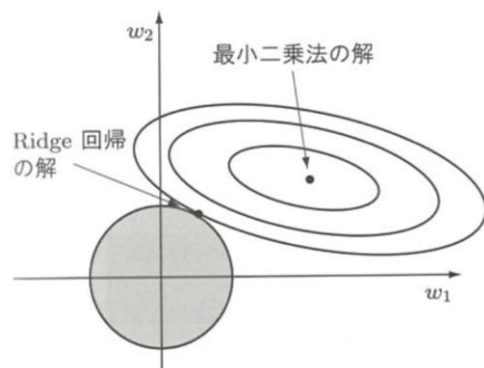
正則化：正則化の振る舞い

□ Ridge回帰

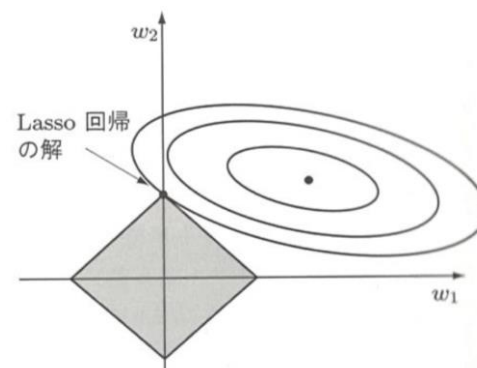
- パラメータの存在する範囲を円（ d 次元では超球）の中に限定して、それぞれの重みが大きな値をとれないようにする
 - 重み：誤差関数の等位線との接点（＝円周上の点）

□ Lasso回帰

- パラメータの和が一定という条件なので、それぞれの軸で角をもつ領域に値が制限
 - 角で誤差関数の等位線と接する（多くのパラメータが0になる）



Ridge回帰における
正則化



Lasso回帰における
正則化

バイアスと分散のトレードオフ

- 回帰式を高次方程式に置き換えて適用できる
 - 特徴ベクトル \mathbf{x} に対して、基底関数ベクトル $\boldsymbol{\phi}(\mathbf{x})$ を考える

$$\boldsymbol{\phi}(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_b(\mathbf{x}))^T$$

- 例：1次元ベクトル x に対して $\boldsymbol{\phi}(x) = (1, x, x^2, \dots, x^b)^T$ となる
- 以下のように回帰式を定義すれば、
係数が線形という条件のもとで最小二乗法が適用可能

$$\hat{c}(\mathbf{x}) = \sum_{j=0}^b w_j \phi_j(\mathbf{x})$$

- 複雑な関数を用いることで、
真のモデルに近い形を表現できるのか？

バイアスと分散のトレードオフ

□ バイアスと分散はトレードオフの関係

- バイアス：真のモデルとの距離
- 分散：学習結果の散らばり具合

□ 単純なモデル → バイアス：大、分散：小

- こべつ 個別のデータに対する誤差が大きくなりやすいが、学習データが少し変動しても結果として得られるパラメータは大きく変動しない
へんどう

□ 複雑なモデル → バイアス：小、分散：大

- 個別のデータに対する誤差を小さくしやすいが、学習データの値が少し変動すると、結果が大きく異なることがある

バイアスと分散のトレードオフ

□ 回帰問題におけるバイアスと分散

■ 線形回帰式の場合

- 求まった超平面は、学習データ内の点をほとんど通らないので、バイアスが大きい

■ 「学習データの個数 - 1」次の高次回帰式の場合

- 求まった回帰式は、全学習データを通る（学習データと一致する関数が求まる）ので、バイアスが小さい
- データが少し動いただけで、この高次式は大きく変動するので、結果の分散は大きい

□ 機械学習では、バイアスー分散のトレードオフを常に意識しなければならない

■ 正則化：緩いバイアスで分散を減らすのに有効

回帰木

□ 回帰木

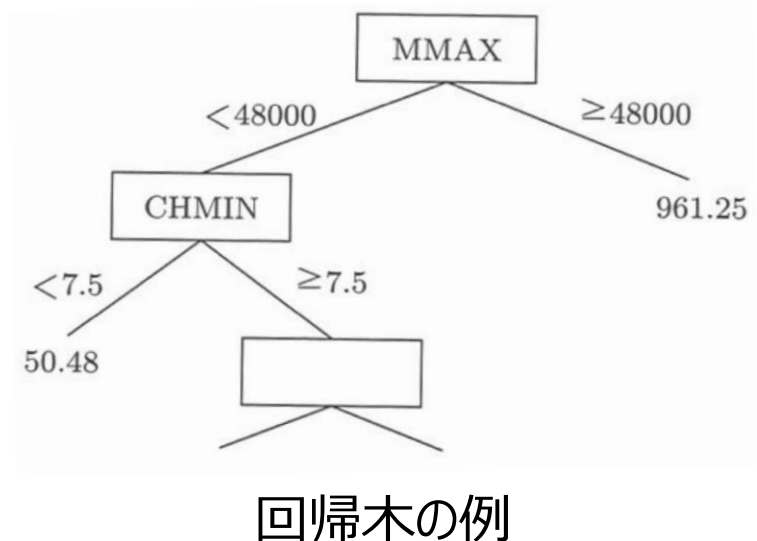
- 識別における決定木の考え方を回帰問題に適用する方法

□ 決定木による識別問題の学習

- 特徴の値によって学習データを同じクラスの集合になるように分割する

□ 回帰木による回帰問題の学習

- 出力値の近いデータが集まるように、特徴の値によって学習データを分割
- 特徴をノードとし、出力値をリーフとする回帰木が得られる



回帰木：CART

□ CART (classification and regression tree)

- 木の構造を二分木^{にぶんぎ}に限定した決定木
- 分類基準：ジニ不純度^{ふじゅんど}（Gini impurity）

□ CARTによる識別問題

- 分類前後の集合のジニ不純度 G を求めて、改善度 ΔG が最大^{さいきてき}のものをノードに選ぶことを再帰的に繰り返す

$$G = 1 - \sum_{j=1}^c N(j)^2 \quad \Delta G(D) = G(D) - P_L \cdot G(D_L) - P_R \cdot G(D_R)$$

D : あるノードに属するデータの全体
 $N(j)$: データ中のクラス j の割合

D_L : 左の部分木 (D_R は右の部分木)
 P_L : D_L に属するデータの割合 (P_R は D_R に属する)

回帰木：CART

□ CARTによる回帰問題

- 分類基準として、データの散らばりSSの減り方 ΔSS が最大になるものを選択

$$SS(D) = \sum_{y_i \in D} (y_i - \tilde{y})^2$$

$$\Delta SS(D) = SS(D) - P_L \cdot SS(D_L) - P_R \cdot SS(D_R)$$

\tilde{y} : D に属するデータの平均値

D : あるノードに属するデータの全体

D_L : 左の部分木 (D_R は右の部分木)

P_L : D_L に属するデータの割合 (P_R は D_R に属する)

$SS(D)$ では、データ D の分散を求めているので、
この基準は分割後の分散が最小となるような分割を求めている

回帰木：モデル木

□ モデル木

■ 回帰木のリーフの値を線形回帰式とした木

- ・ 回帰木と線形回帰の双方の利点を活かした方法

■ モデル木の生成手順

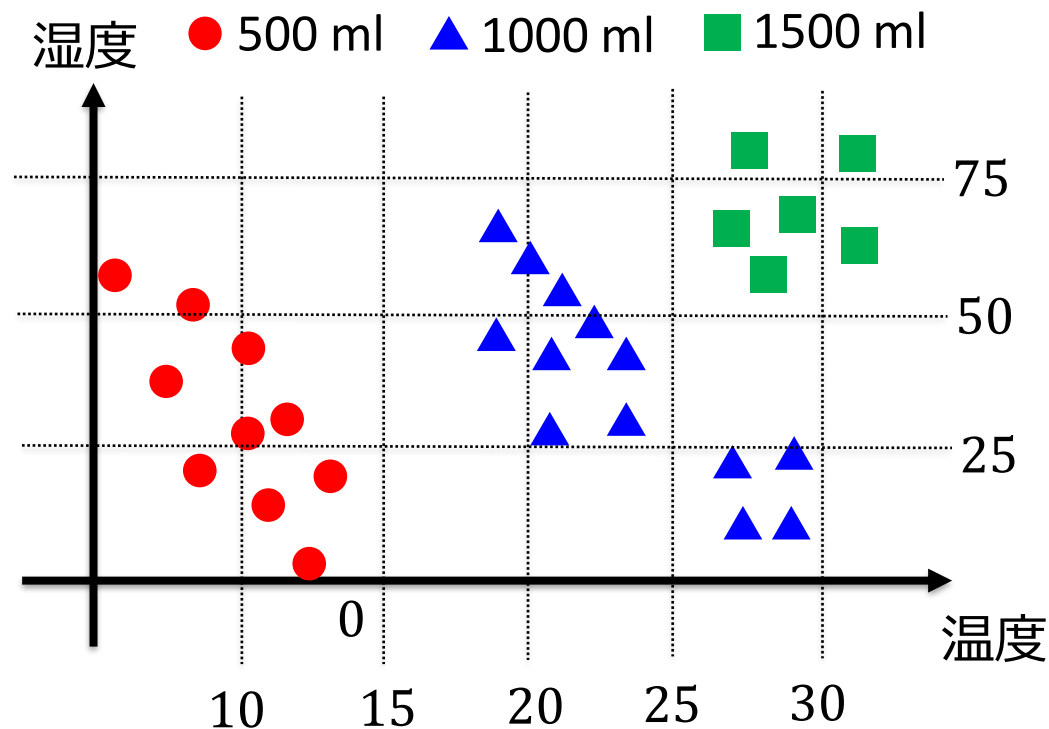
1. 出力値が近い区間を切り出せる特徴を選んでデータ分割
2. 分割後のデータに対して線形回帰式を計算

■ 特定の要因によって振る舞いが異なるデータを分割し、それぞれに対応する規則性を見つけ、かつ、その分割の要因を木構造によって説明できる

- ・ 例えば、季節によって出力に影響を及ぼす要因が異なるデータなどに有効

演習問題6-2（10分間）

- 以下の学習データが与えられたとき、温度と湿度から1日あたりのビールの消費量を予測する回帰木を作成せよ



温度・湿度・1日あたりのビールの消費量の関係