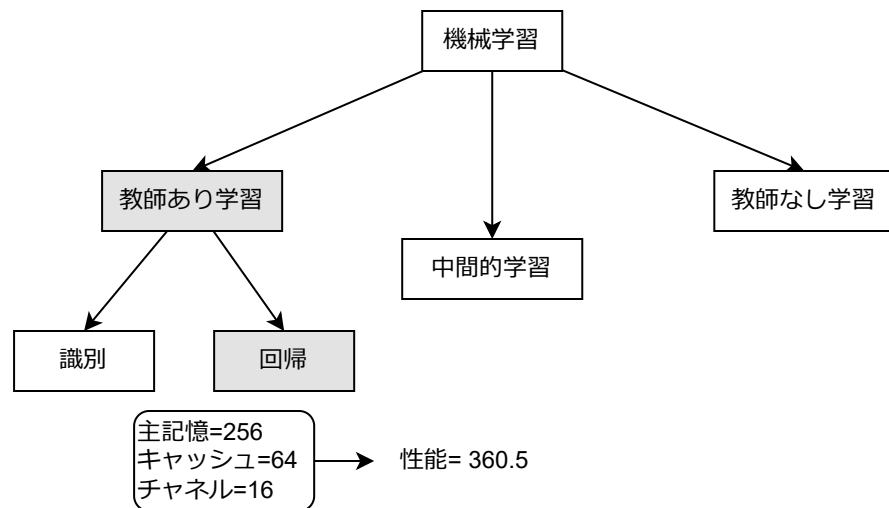


回帰

- 問題設定

- 教師あり学習

- 数値入力 → 数値出力



数値特徴に対する「教師あり・回帰」問題の定義

- 回帰のデータ

- 特徴ベクトル \boldsymbol{x} と正解情報 y のペア

$$\{(\boldsymbol{x}_i, y_i)\}, \quad i = 1, \dots, N$$

- \boldsymbol{x} は数値を要素とする d 次元ベクトル、 y は数値

$$\boldsymbol{x}_i = (x_{i1}, \dots, x_{id})^T$$

$$y \in \mathbb{R}$$

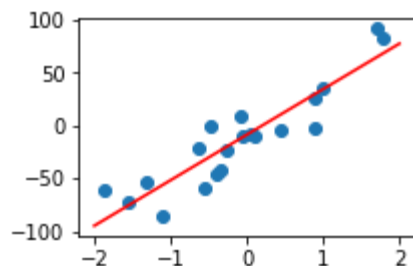
- 回帰問題の正解情報をターゲットとよぶこともある

- 回帰問題は関数 $y = \hat{c}(\boldsymbol{x})$ を求める問題と見なせる

- 未知データに対して予測精度の高い関数を求めたい
- どの特徴が出力値に対してどのような影響を及ぼしているかを知りたい

線形回帰

- 目標: 学習データに対してなるべく誤差の少ない直線を求める



- 問題の定義
 - 入力 \mathbf{x} ($x_0 \equiv 1$ として $d + 1$ 次元に拡張) から出力 y を求める回帰式 $\hat{c}(\mathbf{x})$ を1次式に限定
 - 学習データから回帰式の重み \mathbf{w} を求める

$$\hat{c}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \sum_{j=0}^d w_j x_j$$

線形回帰

- 最小二乗法による重みの推定
 - $E(\mathbf{w})$: 誤差の二乗和を損失関数(最小化の対象)として設定する
 - \mathbf{X} : パターン行列 (特徴ベクトルを転置して行方向に並べたもの。 N 行 $d + 1$ 列)
 - \mathbf{y} : 正解ベクトル (N 次元の列ベクトル)
 - \mathbf{w} : 重みベクトル ($d + 1$ 次元の列ベクトル)

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

- $E(\mathbf{w})$ を \mathbf{w} で偏微分した値を0と置いて、損失関数が極小となる \mathbf{w} を求める

$$\mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = 0 \Leftrightarrow \mathbf{X}^T \mathbf{y} = \mathbf{X}^T \mathbf{X} \mathbf{w} \Leftrightarrow \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

線形回帰

- 線形回帰の精度向上

- 基底関数 $\phi(\mathbf{x}) = (\phi_0(\mathbf{x}), \dots, \phi_b(\mathbf{x}))$ を考える

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

- 基底関数の例

- x が1次元のとき: $\phi(x) = (1, x, x^2, \dots, x^b)$

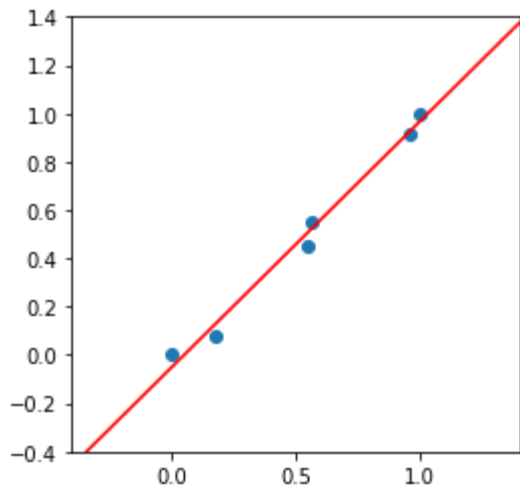
```
from sklearn.preprocessing import PolynomialFeatures
poly = PolynomialFeatures(degree=b)
X2 = poly.fit_transform(X)
```

- \mathbf{x} が2次元(x_1, x_2)のとき(`degree=2`): $\phi(\mathbf{x}) = (1, x_1, x_2, x_1^2, x_1x_2, x_2^2)$

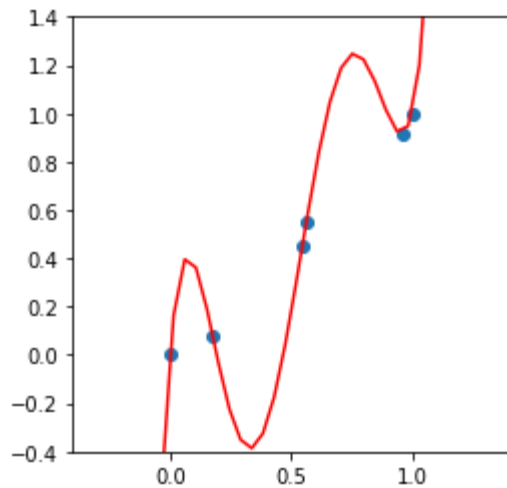
- 重み \mathbf{w} が線形であれば、最小二乗法が適用可能

線形回帰

- 基底関数を用いたときの問題点
 - 汎化性能の低下



(a) 基底関数を用いずに1次式で線形回帰



(b) 基底関数を用いて5次式で線形回帰

回帰モデルの評価

- 回帰モデルの評価法(1)
 - 平均二乗誤差 (MSE: Mean Squared Error)
 - 手法間の評価に有効

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{c}(\mathbf{x}_i))^2$$

- 相関係数
 - 出力と正解とがどの程度似ているかを-1から1までの値で表す

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

$\tilde{c} : \hat{c}(\mathbf{x}_i)$ の平均, $\tilde{y} : y$ の平均

回帰モデルの評価

- 回帰モデルの評価法(2)
 - 決定係数
 - 「正解と予測の差の二乗」と「正解の分散」との比を1から引いたもの
 - 相関係数 R の二乗と等しい

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

正則化

- 正則化の考え方
 - 誤差関数に正則化項を導入 → 複雑な重み \mathbf{w} (過学習) の回避
 - L1ノルム $|\mathbf{w}|$: 0となる重みが多くなる
 - L2ノルム $\|\mathbf{w}\|^2$: 重みを0に近づける
- リッジ回帰
 - 誤差の二乗和にL2ノルム正則化項を加える

$$E(\mathbf{w}) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w}) + \alpha \mathbf{w}^T \mathbf{w} \quad \alpha : \text{正則化項の重み}$$

- \mathbf{w} が解析的に求まる

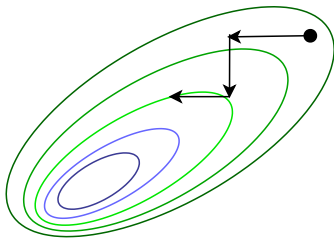
$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

正則化

- ラッソ回帰
 - 誤差の二乗和にL1ノルム正則化項を加える

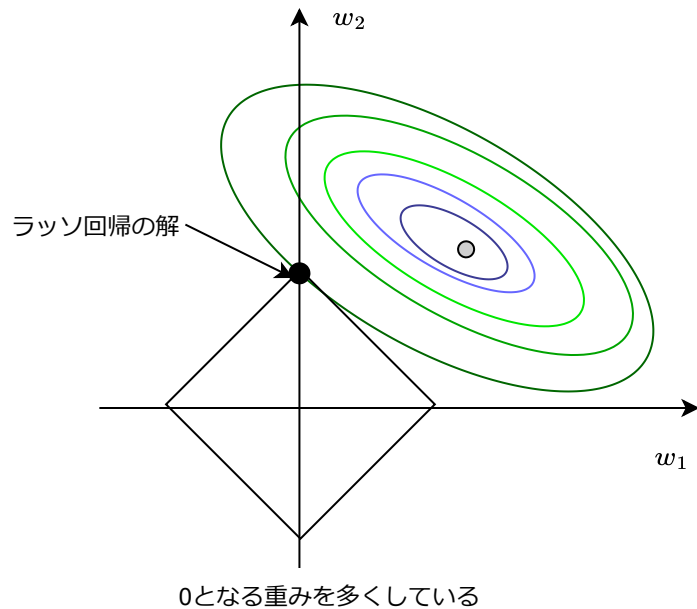
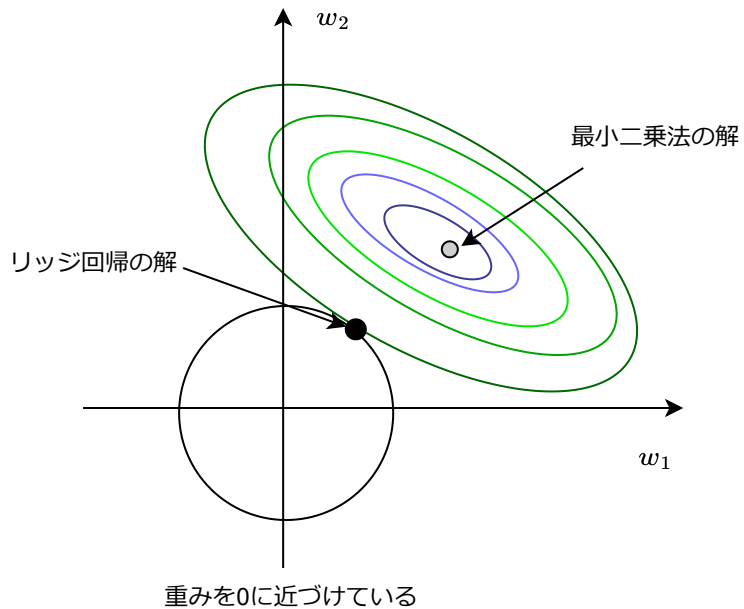
$$E(\mathbf{w}) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w}) + \alpha \sum_{j=1}^d |w_j|$$

- 微分不可能な点があるため、解析的に解を求められない
- 解法の例: coordinate descent algorithm
 - 1つの変数(軸)の値だけを誤差が減る方向に変更することを繰り返す



正則化

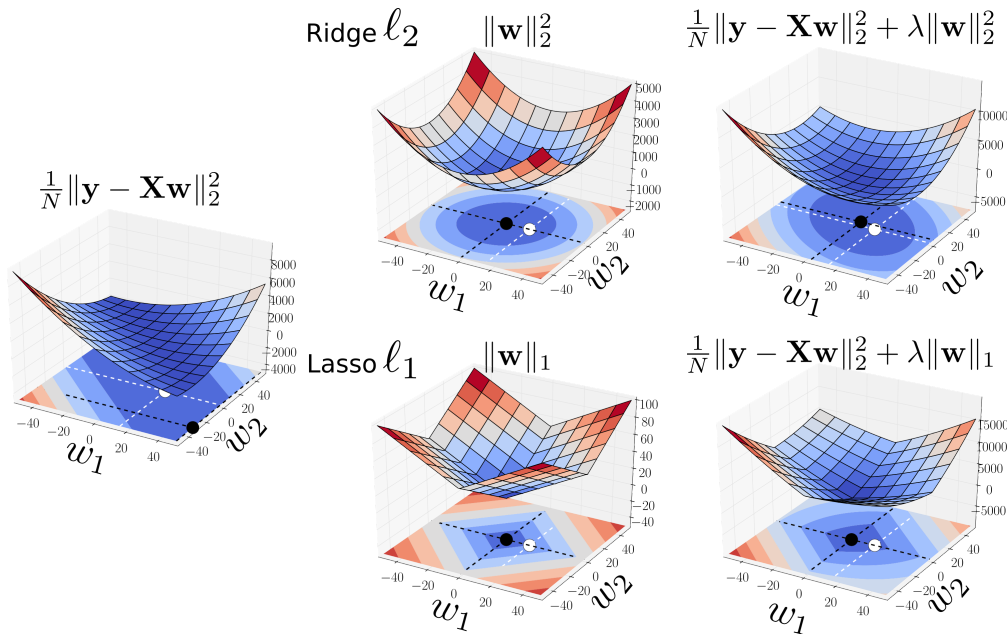
- リッジ回帰とラッソ回帰



正則化

- リッジ回帰とラッソ回帰

- <https://duchesnay.github.io/pystatsml/machine-learning/linear-regression.html>

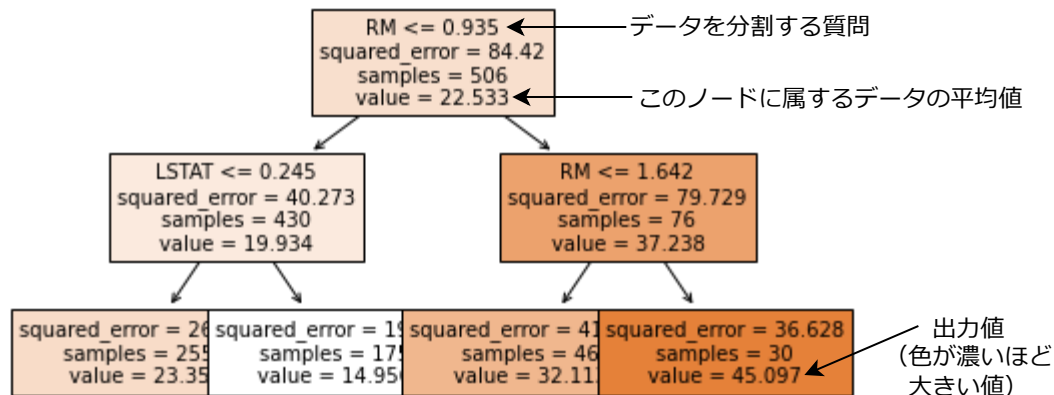


バイアスー分散のトレードオフ

- バイアスと分散
 - バイアス: 正解(=真のモデル)からのズレ
 - 分散: 求まるモデルの安定性
- 単純なモデル
 - 正解をカバーしていないかもしれない → バイアス 大
 - データが多少ぶれても結果は似ている → 分散 小
- 複雑なモデル
 - 正解をカバーしている可能性が高い → バイアス 小
 - データが少し違えば結果が大きく異なる → 分散 大 (正則化項はこれを小さくしようとしている)

回帰木

- 回帰木とは
 - 識別における決定木の考え方を回帰問題に適用
 - ターゲット値の分散が小さくなるように分割



回帰木

- CART (classification and regression tree)

- 木の構造を二分木に限定
- 識別の際のデータの分類基準はジニ不純度
 - 2クラスの場合のジニ不純度

$$\text{Gini}(p) = 2p(1 - p) \quad p : \text{正例の割合}$$

- クラスの出現が等確率のとき最大
- 回帰に用いるときのデータの分類基準: ターゲット値の分散
 - 子ノードの重み付き分散和が最小 = 分割後の分散の減少量が最大 となる特徴を選ぶ

回帰木

- CARTの特徴選択基準

- データ D の分散

$$SS(D) = \sum_{i=1}^N (y_i - \tilde{y})^2$$

- 分割後の分散の減少量が最大となる特徴を選ぶ

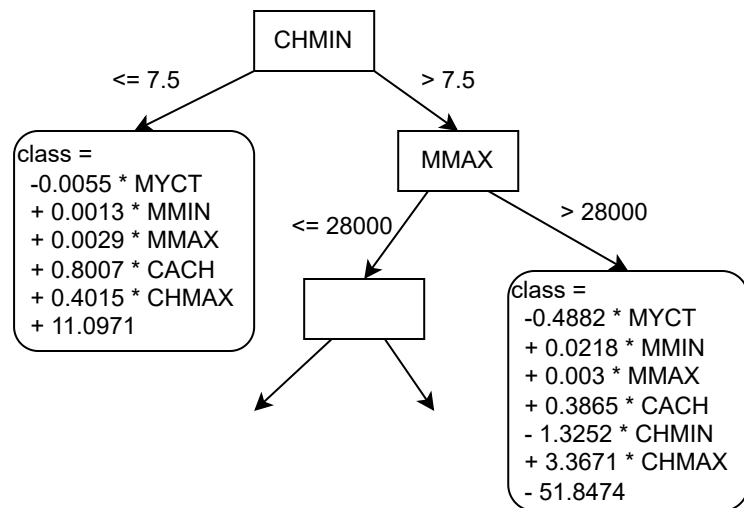
$$\Delta SS(D) = SS(D) - P_L \cdot SS(D_L) - P_R \cdot SS(D_R)$$

- D_L : 左部分木のデータ
- D_R : 右部分木のデータ
- P_L : 左部分木のデータ数の割合
- P_R : 右部分木のデータ数の割合

モデル木

- モデル木とは

- リーフを線形回帰式にした回帰木



まとめ

- 線形回帰
 - 最小二乗法で回帰関数のパラメータを求めることができる
 - 基底関数によって複雑な関数も表現可能
 - 正則化で過学習を回避
- 回帰木
 - 決定木を用いて回帰を行う
 - 振る舞いが異なるデータが混合しているときに有効
 - リーフを線形回帰式にしたものがモデル木