

機械学習基礎

赤穂 昭太郎

akaho@ism.ac.jp

<https://github.com/toddler2009/ml-tutorial>

2025 年 4 月 18 日

目次

第 2 章	教師あり学習	7
2.1	線形モデル	7
2.1.1	重回帰モデル	7
2.1.2	多項式回帰	7
2.1.3	正則化	8
2.1.4	具体例：California Housing データによる回帰問題	9
2.1.5	ベイズ推定と正則化の関係	13
2.2	基底関数とカーネル法による非線形予測	13
2.2.1	基底関数による線形モデルの拡張	13
2.2.2	カーネル法の導入	14
2.2.3	再生核ヒルベルト空間（RKHS）の簡単な導入	15
2.2.4	カーネル回帰と Representer Theorem	15
2.2.5	カーネル法の応用：非数値データへの展開	16
2.2.6	補足：再生核ヒルベルト空間（RKHS）	16
2.2.7	カーネルリッジ回帰の導出	17
2.2.8	補足：線形モデルにおける交差検証誤差の効率的計算	21
2.3	サポートベクターマシン（SVM）	22
2.3.1	マージン最大化による分類	22
2.3.2	ソフトマージンと損失関数の定式化	23
2.3.3	カーネルによる非線形分類	24
2.3.4	サポートベクトルと解のスパース性	24

2.3.5	カーネル SVM の双対定式化と凸最適化の概要	25
2.3.6	SVM とロジスティック回帰の比較：損失関数の観点から . . .	25
2.3.7	サポートベクトルと決定境界の直感的理解	26
2.4	決定木とアンサンブル学習	27
2.4.1	決定木の定義と最適化	27
2.4.2	具体例：Breast Cancer Wisconsin データによる識別問題 . .	28
2.4.3	アンサンブル学習の基本的な考え方	31
2.4.4	ランダムフォレスト	33
2.4.5	勾配ブースティングの具体的アルゴリズム（回帰：二乗誤差 の場合）	35
2.4.6	XGBoost：実用的な勾配ブースティングの高速実装	36
2.4.7	アンサンブル学習における変数重要度とモデルの解釈性 . . .	40
2.4.8	変数重要度の一般的な考え方	41
2.4.9	決定木アンサンブルに特化した重要度の評価	41
2.4.10	補足：Shapley 値の数学的定義	43
2.5	多クラス分類と誤り訂正符号による拡張	45
2.5.1	多クラス分類の基本的アプローチ	45
2.5.2	One-vs-Rest (OvR)	45
2.5.3	One-vs-One (OvO)	47
2.5.4	ECOC（誤り訂正出力符号：Error-Correcting Output Codes）	47
2.5.5	補足：誤り訂正符号としての数式的性質	48
2.5.6	補足：不均衡データにおける多クラス分類の難しさ	49
2.6	最近傍識別器（ k -nearest neighbor 法）	50
2.6.1	基本的な定義	50
2.6.2	ベイズ識別器との性能比較	50
2.6.3	計算アルゴリズム上の工夫	51
2.6.4	k -NN 法の特徴と注意点	51
2.7	そのほかの教師あり学習モデル	53

第 2 章

教師あり学習

2.1 線形モデル

2.1.1 重回帰モデル

単回帰モデルを拡張し，複数の説明変数を用いて目的変数を予測するモデルが重回帰モデルである．モデルは次のように表される：

(59)

$$y = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p + \varepsilon$$

(71)

ここで誤差項 ε は独立同分布の正規分布 $\mathcal{N}(0, \sigma^2)$ に従うと仮定する．この仮定の下では，最小二乗法でパラメータを推定することは，最尤推定と一致する．

最小二乗法では次の目的関数を最小化する：

$$\min_{\beta} \|\mathbf{y} - X\beta\|^2$$

このとき，正規方程式により解は次のように求まる（ $X^T X$ が正則な場合）：

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

2.1.2 多項式回帰

説明変数の多項式を用いることで，非線形な関係も線形回帰の枠組みで扱うことができる．例えば 3 次の多項式回帰は次のようになる：

(290)

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \varepsilon$$

変数変換により、パラメータに対しては線形性が保たれているため、最小二乗法をそのまま適用できる。

2.1.3 正則化

モデルの柔軟性が高すぎると、訓練データに過剰に適合する（過学習）リスクがある。これを防ぐために、パラメータにペナルティを課す「正則化」が用いられる。

リッジ回帰

リッジ回帰では、L2 ノルムによってパラメータの大きさを抑制する。目的関数は
(237) 以下の通り：

$$\min_{\beta} (\|\mathbf{y} - X\beta\|^2 + \lambda \|\beta\|^2)$$

このときの解は以下のように計算される：

$$\beta_{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T \mathbf{y}$$

ここで λ は正則化の強さを調整するハイパーパラメータであり、 I は単位行列である。

スパース回帰 (Lasso)

Lasso 回帰では、L1 ノルムによる正則化を行い、いくつかの係数をゼロにする効果
(241) を持つ。目的関数は以下の通り：

$$\min_{\beta} \left(\|\mathbf{y} - X\beta\|^2 + \lambda \sum_j |\beta_j| \right)$$

これにより、重要な説明変数の選択（変数選択）を同時に行うことができる。

2.1.4 具体例：California Housing データによる回帰問題

データセットの概要

California Housing Dataset は、米国カリフォルニア州の各地域における住宅 (14) に関する統計情報を用いて、**住宅価格の中央値**を予測する回帰問題である*¹。

このデータセットは、1990 年の国勢調査データに基づいており、scikit-learn の `fetch_california_housing` 関数で取得できる。

- サンプル数：20,640
- 特徴量数：8（すべて数値）
- 目的変数：住宅の中央値価格（単位：\$100,000）

特徴量の説明

各サンプルは、ある地域の統計的な特徴を表しており、以下の変数が含まれる：

- **MedInc**：地域の世帯収入の中央値
- **HouseAge**：住宅の築年数の中央値
- **AveRooms**：世帯あたりの平均部屋数
- **AveBedrms**：世帯あたりの平均寝室数
- **Population**：地域の人口
- **AveOccup**：世帯あたりの平均居住者数
- **Latitude**：緯度
- **Longitude**：経度

これらの特徴量はすべて連続値であり、**前処理として標準化（平均 0，分散 1）**を行うことで、多くの回帰モデルが適切に動作する。実際に、リッジ回帰と LASSO 回帰を比較した実験のための Python コードと結果を示す。学習データは 100 個、テストデータは 300 個をランダムに選び、標準化した上で、リッジ回帰と LASSO 回帰の

*¹ もともと有名なデータセットである Boston Housing dataset (ISL にもリストされている) は人種などコンプライアンス問題が議論されるようになり scikit-learn から削除された。

それぞれで正則化パラメータ α を変化させながら (正則化パス追跡), MSE (Mean Square Error) と係数の値の変化をプロットした (図 2.1).

Program 2.1 リッジ回帰と LASSO 回帰

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.datasets import fetch_california_housing
4 from sklearn.linear_model import Ridge, Lasso
5 from sklearn.metrics import mean_squared_error
6 from sklearn.preprocessing import StandardScaler
7
8 # データの読み込み
9 data = fetch_california_housing()
10 X, y = data.data, data.target
11 feature_names = data.feature_names
12
13 # 学習・テストデータの選択
14 k, m = 100, 300
15 np.random.seed(42)
16 indices = np.random.permutation(len(X))
17 train_idx, test_idx = indices[:k], indices[k:k + m]
18 X_train, y_train = X[train_idx], y[train_idx]
19 X_test, y_test = X[test_idx], y[test_idx]
20
21 # 標準化
22 scaler = StandardScaler()
23 X_train_std = scaler.fit_transform(X_train)
24 X_test_std = scaler.transform(X_test)
25
26 # 正則化パラメータ
```

```
27 alphas = np.logspace(-3, 2, 100)
28
29 # Ridge と Lasso の誤差と係数追跡
30 ridge_mse_train, ridge_mse_test = [], []
31 lasso_mse_train, lasso_mse_test = [], []
32 ridge_coefs, lasso_coefs = [], []
33
34 for alpha in alphas:
35     ridge = Ridge(alpha=alpha).fit(X_train_std, y_train)
36     lasso = Lasso(alpha=alpha, max_iter=10000).fit(
37         X_train_std, y_train)
38
39     ridge_mse_train.append(mean_squared_error(y_train,
40         ridge.predict(X_train_std)))
41     ridge_mse_test.append(mean_squared_error(y_test,
42         ridge.predict(X_test_std)))
43
44     lasso_mse_train.append(mean_squared_error(y_train,
45         lasso.predict(X_train_std)))
46     lasso_mse_test.append(mean_squared_error(y_test,
47         lasso.predict(X_test_std)))
48
49     ridge_coefs.append(ridge.coef_)
50     lasso_coefs.append(lasso.coef_)
51
52 # プロット
53 fig, axs = plt.subplots(2, 2, figsize=(12, 10))
54
55 # Ridge: 訓練 / テスト MSE
```

```
51 | axes[0, 0].plot(alphas, ridge_mse_train, label='Train MSE
    | ', linestyle='--')
52 | axes[0, 0].plot(alphas, ridge_mse_test, label='Test MSE')
53 | axes[0, 0].set_xscale('log')
54 | axes[0, 0].set_xlabel('Alpha')
55 | axes[0, 0].set_ylabel('MSE')
56 | axes[0, 0].set_title('Ridge: Train & Test MSE vs Alpha')
57 | axes[0, 0].legend()
58 | axes[0, 0].grid(True)
59 |
60 | # Lasso: 訓練 / テスト MSE
61 | axes[0, 1].plot(alphas, lasso_mse_train, label='Train MSE
    | ', linestyle='--')
62 | axes[0, 1].plot(alphas, lasso_mse_test, label='Test MSE')
63 | axes[0, 1].set_xscale('log')
64 | axes[0, 1].set_xlabel('Alpha')
65 | axes[0, 1].set_ylabel('MSE')
66 | axes[0, 1].set_title('Lasso: Train & Test MSE vs Alpha')
67 | axes[0, 1].legend()
68 | axes[0, 1].grid(True)
69 |
70 | # Ridge: 係数パス
71 | axes[1, 0].plot(alphas, ridge_coefs)
72 | axes[1, 0].set_xscale('log')
73 | axes[1, 0].set_xlabel('Alpha')
74 | axes[1, 0].set_ylabel('Coefficient Value')
75 | axes[1, 0].set_title('Ridge: Coefficient Paths')
76 | axes[1, 0].legend(feature_names, loc='best', fontsize='
    | small')
```

```
77
78 # Lasso: 係数パス
79 axs[1, 1].plot(alphas, lasso_coefs)
80 axs[1, 1].set_xscale('log')
81 axs[1, 1].set_xlabel('Alpha')
82 axs[1, 1].set_ylabel('Coefficient Value')
83 axs[1, 1].set_title('Lasso: Coefficient Paths')
84 axs[1, 1].legend(feature_names, loc='best', fontsize='
    small')
85
86 plt.tight_layout()
87 plt.show()
```

2.1.5 ベイズ推定と正則化の関係

正則化は、ベイズ推定の観点から自然に解釈できる。

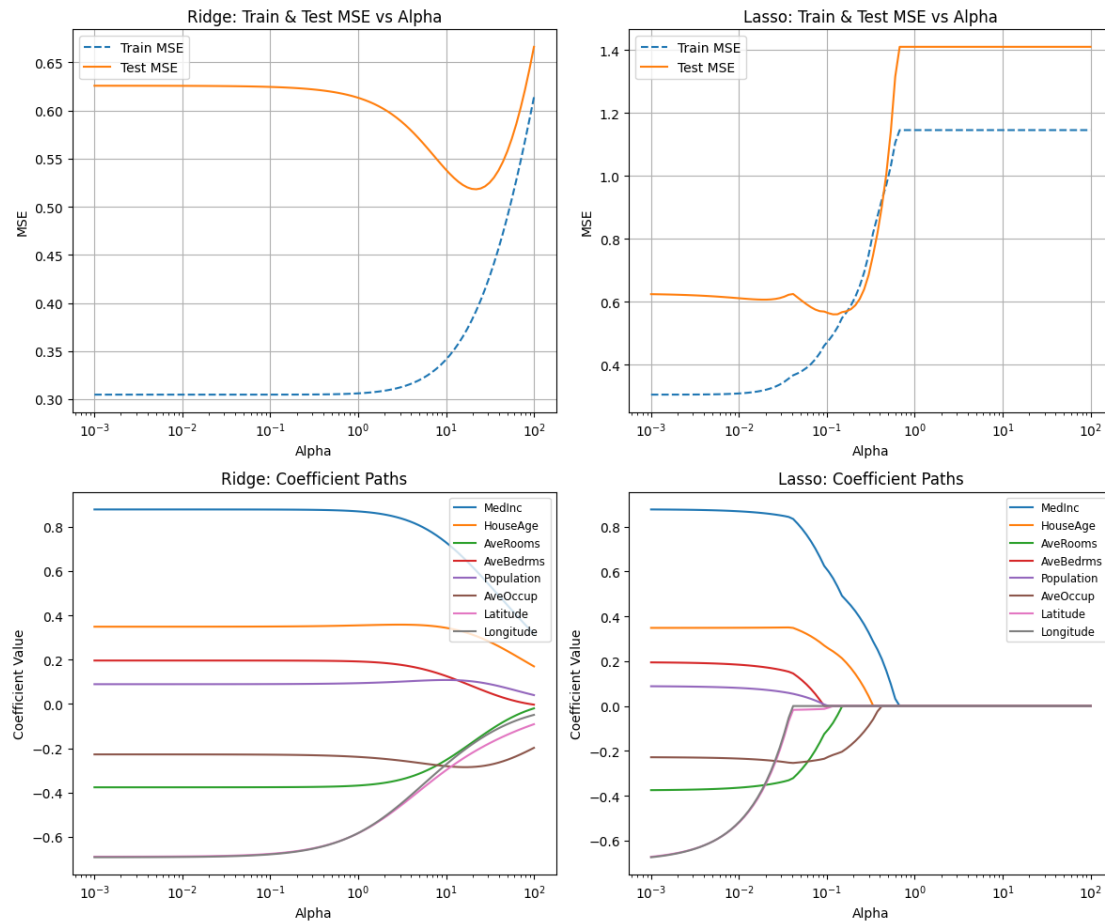
- リッジ回帰は、係数 β に対してゼロ平均の正規分布を事前分布としたときの MAP 推定に対応する。
- Lasso 回帰は、ラプラス分布（尖った分布）を事前分布としたときの MAP 推定に対応する。

このように、正則化はベイズ的には事前知識をモデルに反映させる手法として位置づけられる。

2.2 基底関数とカーネル法による非線形予測

2.2.1 基底関数による線形モデルの拡張

線形モデルは、説明変数に対して線形な関係を仮定するが、入力を変換することで非線形な関係も表現可能である。具体的には、入力 $\mathbf{x} \in \mathbb{R}^d$ を基底関数 $\phi(\mathbf{x}) \in \mathbb{R}^M$

図 2.1 リッジ回帰・Lasso 回帰 ($n = 100$)

(294) によって変換し，その上で線形回帰を行う：

$$f(\mathbf{x}) = \sum_{j=1}^M w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

ここで， $\phi_j(\mathbf{x})$ は入力に対する非線形な変換（例：多項式，ガウス関数，シグモイドなど）である。

この枠組みでは，非線形な入力に対しても予測モデルは依然としてパラメータ \mathbf{w} に対して線形であるため，線形回帰と同様の解析手法が使える。

2.2.2 カーネル法の導入

基底関数を使うと柔軟な表現が可能になるが，基底関数の数が多くなると，計算量やメモリの問題が生じる．そこで，明示的に基底関数を用いずに，内積だけを通じて

学習を行う方法が「カーネル法」である [5].

(382)

入力 \mathbf{x}, \mathbf{x}' に対し, 以下のような関数 k を定義する:

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

このような関数 k を「カーネル関数」と呼ぶ.

2.2.3 再生核ヒルベルト空間 (RKHS) の簡単な導入

カーネル法の理論的基盤として, 「再生核ヒルベルト空間 (Reproducing Kernel Hilbert Space; RKHS)」がある.

RKHS は, 関数空間の一種であり, 各点での評価が内積によって再現できる性質 (再生性) を持つ. すなわち, ある関数 $f \in \mathcal{H}$ に対して, 以下が成り立つ:

$$f(\mathbf{x}) = \langle f, k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}}$$

この性質により, 学習アルゴリズムがカーネル関数を通して関数空間上の最適化問題として定式化できる.

2.2.4 カーネル回帰と Representer Theorem

リッジ回帰をカーネル法に拡張し, 以下の問題を考える:

$$\min_{f \in \mathcal{H}} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|_{\mathcal{H}}^2$$

このような無限次元の最適化問題に対して, 「Representer Theorem (表現定理)」が重要な役割を果たす. この定理は, 最適解が訓練データに対するカーネル関数の線形結合として表されることを保証する:

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$$

これにより, 未知の関数 f を直接求めるのではなく, 係数 α_i を求める有限次元の問題として扱えるようになる.

2.2.5 カーネル法の応用：非数値データへの展開

カーネル関数は、数値ベクトルだけでなく、文字列、グラフ、構造データなどにも定義可能である。たとえば、以下のような応用がある：

- 文字列類似度を測るスペクトラムカーネル（バイオインフォマティクス）
- グラフ構造を対象としたグラフカーネル（化学構造の分類など）
- 木構造や文法に基づいた構文カーネル（自然言語処理）

このように、カーネル法は入力 of 構造に応じて柔軟に拡張できる点が大きな利点である。

2.2.6 補足：再生核ヒルベルト空間（RKHS）

再生核ヒルベルト空間（RKHS）は、内積を備えた関数空間であり、すべての点 \mathbf{x} に対して「点評価関数」 $f \mapsto f(\mathbf{x})$ が連続線形関数であるという性質を持つ。

再生核 $k(\mathbf{x}, \mathbf{x}')$ を持つ RKHS \mathcal{H} では、以下の性質が成り立つ：

- 再生性： $f(\mathbf{x}) = \langle f, k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}}$
- 対称性と正定値性： $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$, および $\sum_{i,j} c_i c_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0$

具体例（ガウスクーネル・RBF コーネル (radial basis function)）：

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$$

このようなカーネルに対応する RKHS は、非常に滑らかな無限次元の関数空間となる。

Representer Theorem の証明スケッチ

Representer Theorem は、次のような形式の最適化問題：

$$\min_{f \in \mathcal{H}} \sum_{i=1}^n L(y_i, f(\mathbf{x}_i)) + \lambda \|f\|_{\mathcal{H}}^2$$

に対して、解 $f^* \in \mathcal{H}$ が以下の形で表せることを保証する：

$$f^*(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$$

証明スケッチ：

- \mathcal{H} の任意の関数 $f \in \mathcal{H}$ を、訓練データ点に対応するカーネル関数の線形結合部分と、その直交補空間（それらに直交する関数）に分解する。
- 正則化項 $\|f\|_{\mathcal{H}}^2$ は直交分解の二乗和になり、訓練データに対する誤差は直交補空間成分は影響せず、目的関数を最小にするには直交補空間成分を 0 にするのが最適となる。
- よって、最適解は訓練データのみ依存するカーネルの線形結合で表現できる。

2.2.7 カーネルリッジ回帰の導出

Representer Theorem によれば、最適な関数 f は次のように表される：

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$$

これをすべての訓練データに対して適用し、ベクトル形式で表すと：

$$\mathbf{f} = K\boldsymbol{\alpha}$$

ここで：

- $K \in \mathbb{R}^{n \times n}$: (i, j) 成分が $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ であるカーネル行列（Gram 行列）
- $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)]^T$
- $\boldsymbol{\alpha} \in \mathbb{R}^n$: 係数ベクトル

正則化付きの二乗誤差損失を最小化する問題：

$$\min_{\boldsymbol{\alpha}} \|\mathbf{y} - K\boldsymbol{\alpha}\|^2 + \lambda \boldsymbol{\alpha}^T K \boldsymbol{\alpha}$$

これを解くと、以下の正則方程式が得られる：

$$(K + \lambda I)\alpha = \mathbf{y}$$

この連立方程式を解くことで、カーネルリッジ回帰モデルが構築できる。

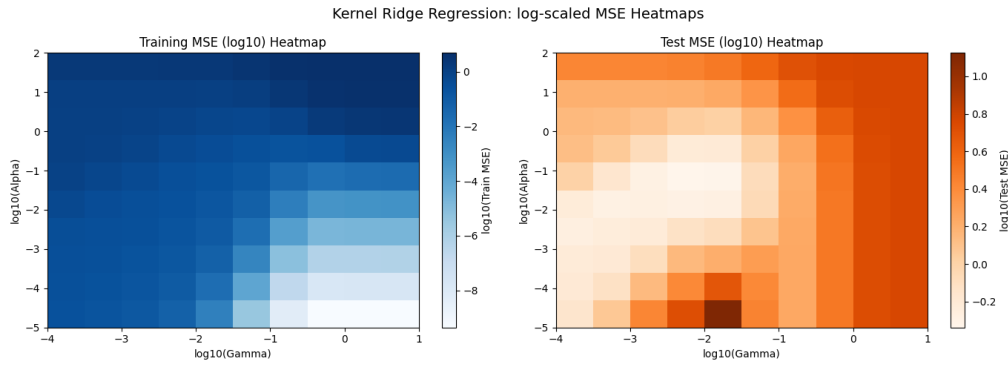
実際に California Housing dataset に対し、カーネルリッジ回帰モデルの正則化パラメータ ($\alpha = \lambda$) と (ガウスカーネルの) カーネル幅の逆数 (γ) を変化させたときの訓練誤差とテスト誤差をヒートマップで表現したのが図 2.2 である (学習データ数は結果をわかりやすくするためにロジスティック回帰と LASSO 回帰の比較のときと同様に 100 個に絞った)。正則化パラメータを 0 に近づけ、カーネル幅を小さくすれば、訓練誤差は 0 に近づく一方、テスト誤差については適切なパラメータが存在することがわかる。

Program 2.2 カーネル回帰

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.datasets import fetch_california_housing
4 from sklearn.kernel_ridge import KernelRidge
5 from sklearn.metrics import mean_squared_error
6 from sklearn.preprocessing import StandardScaler
7
8 # データの読み込み
9 data = fetch_california_housing()
10 X, y = data.data, data.target
11
12 # 学習・テストデータの選択
13 k, m = 100, 300
14 np.random.seed(42)
15 indices = np.random.permutation(len(X))
16 train_idx, test_idx = indices[:k], indices[k:k + m]
17 X_train, y_train = X[train_idx], y[train_idx]
18 X_test, y_test = X[test_idx], y[test_idx]
```

```
19
20 # 標準化
21 scaler = StandardScaler()
22 X_train_std = scaler.fit_transform(X_train)
23 X_test_std = scaler.transform(X_test)
24
25 # ハイパーパラメータのグリッド
26 alphas = np.logspace(-5, 2, 10)
27 gammas = np.logspace(-4, 1, 10)
28
29 # 誤差の格納用
30 train_errors = np.zeros((len(alphas), len(gammas)))
31 test_errors = np.zeros((len(alphas), len(gammas)))
32
33 # モデルの学習と評価
34 for i, alpha in enumerate(alphas):
35     for j, gamma in enumerate(gammas):
36         model = KernelRidge(alpha=alpha, kernel='rbf',
37                               gamma=gamma)
38         model.fit(X_train_std, y_train)
39         y_train_pred = model.predict(X_train_std)
40         y_test_pred = model.predict(X_test_std)
41         train_errors[i, j] = mean_squared_error(y_train,
42                                                  y_train_pred)
43         test_errors[i, j] = mean_squared_error(y_test,
44                                                y_test_pred)
45
46 # log10 変換
47 log_train_errors = np.log10(train_errors)
```

```
45 log_test_errors = np.log10(test_errors)
46
47 # ヒートマップの描画
48 fig, axs = plt.subplots(1, 2, figsize=(14, 5))
49
50 # 訓練誤差 (対数表示)
51 im1 = axs[0].imshow(log_train_errors, origin='lower',
52                     aspect='auto',
53                     extent=[np.log10(gammas[0]), np.
54                             log10(gammas[-1]),
55                             np.log10(alphas[0]), np.
56                                 log10(alphas[-1])],
57                     cmap='Blues')
58 axs[0].set_title("Training MSE (log10) Heatmap")
59
60 # テスト誤差 (対数表示)
61 im2 = axs[1].imshow(log_test_errors, origin='lower',
62                     aspect='auto',
63                     extent=[np.log10(gammas[0]), np.
64                             log10(gammas[-1]),
65                             np.log10(alphas[0]), np.
66                                 log10(alphas[-1])],
67                     cmap='Oranges')
68 axs[1].set_title("Test MSE (log10) Heatmap")
69 axs[1].set_xlabel("log10(Gamma)")
70 axs[1].set_ylabel("log10(Alpha)")
```

図 2.2 カーネルリッジ回帰 ($n = 100$)

```

68 plt.colorbar(im2, ax=axis[1], label="log10(Test MSE)")
69
70 plt.suptitle("Kernel Ridge Regression: log-scaled MSE
    Heatmaps", fontsize=14)
71 plt.tight_layout()
72 plt.show()

```

2.2.8 補足：線形モデルにおける交差検証誤差の効率的計算

重回帰モデルやカーネルリッジ回帰では、最適な出力予測ベクトル $\hat{\mathbf{y}}$ が以下のような形で書ける：

(202)

$$\hat{\mathbf{y}} = H\mathbf{y}$$

ここで、 H は入力行列 X のみに依存する行列（ハット行列、またはスムージング行列と呼ばれる）であり、 \mathbf{y} には依存しない。この H の対角成分 H_{ii} は、 i 番目のデータに対する予測の影響度（自己影響）を表す。

Leave-one-out 交差検証 (LOOCV)

交差検証は、汎化性能を評価するために用いられる手法であり、データを訓練用と検証用に分割して誤差を評価する。特に Leave-one-out 交差検証 (LOOCV) は、 n 個のデータのうち 1 つをテストに用い、残り $n - 1$ 個で学習を行い、これを n 回繰

り返す方法である。

この方法は、汎化誤差に近い評価を得られる一方で、通常は n 回の学習が必要であり計算コストが高い。しかし、線形モデルで $\hat{\mathbf{y}} = H\mathbf{y}$ の形を取る場合には、次の公式により 1 回の学習結果から LOOCV 誤差を計算することができる：

$$\text{CV}_{\text{LOO}} = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{y}_i}{1 - H_{ii}} \right)^2$$

ここで、 $\hat{y}_i = (H\mathbf{y})_i$ は全データを用いて学習したときの i 番目の予測値である。(証明は Wahba[6], 赤穂 [5], 萩原 [4] など参照) これは PRESS (PREdiction Sum of Squares) 統計量とも呼ばれている (一般の k -fold CV についても同様の式が一応存在する)。

■まとめ 線形モデルで予測が $\hat{\mathbf{y}} = H\mathbf{y}$ の形で書ける場合には、LOOCV 誤差は H の対角成分と全データでの予測値を使って効率的に計算できる。この性質は、計算コストの削減だけでなく、解析的な誤差評価の導出にも有用である。

2.3 サポートベクターマシン (SVM)

2.3.1 マージン最大化による分類

(367) サポートベクターマシン (Support Vector Machine; SVM) は、分類問題において「マージン (margin)」と呼ばれる量を最大化することにより、分類器の汎化性能を高めることを目的とした手法である。

ここでは、2 クラス分類 ($y_i \in \{-1, +1\}$) を考える。

分類平面は、次のような線形関数によって定義される：

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

このとき、すべてのデータが以下の不等式を満たす場合、線形分離が可能である：

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

この条件のもとで、マージン (境界から最も近い点までの距離 図 2.3) を最大化す

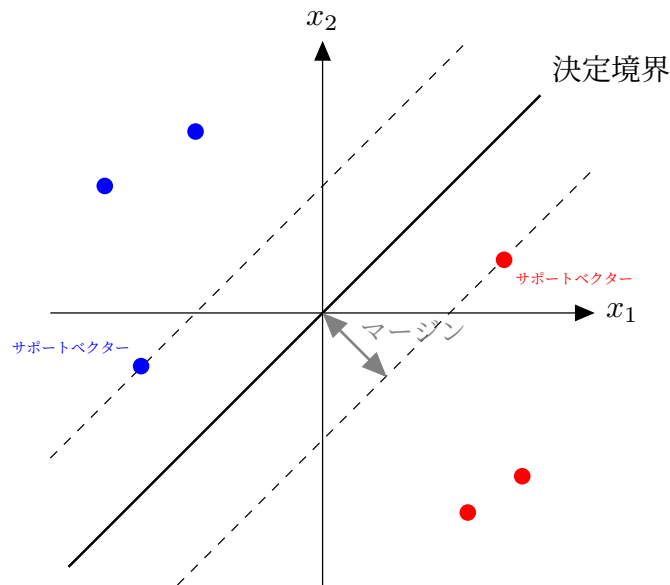


図 2.3 SVM のマージンとサポートベクターの図：破線はマージン境界、黒線は決定境界

るためには、以下の最適化問題を解く：

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad (\forall i)$$

この最適化により、分類平面の傾きを小さく（＝マージンを広く）することができる。

2.3.2 ソフトマージンと損失関数の定式化

実際のデータは完全に線形分離できるとは限らないため、「ソフトマージン」として誤分類を許容する形式に拡張する。これにより損失関数が導入され、以下のような目的関数になる：

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b))$$

ここで：

- 第 1 項：2 次正則化項（モデルの複雑さの抑制）
- 第 2 項：ヒンジ損失関数 (hinge loss)

$$\ell_{\text{hinge}}(y, f(\mathbf{x})) = \max(0, 1 - yf(\mathbf{x}))$$

- $C > 0$: 誤分類に対するペナルティの強さを調整するパラメータ

ヒンジ損失は、分類境界から離れていないデータ点に対してのみ損失を与える「区
分線形関数」であり、サポートベクトルの選定に寄与する。

2.3.3 カーネルによる非線形分類

SVM はカーネル法と組み合わせることで、非線形な分類問題にも対応できる。入
力を明示的に変換せずに、内積の代わりにカーネル関数 $k(\mathbf{x}, \mathbf{x}')$ を用いることで、次
のような非線形分類器が得られる：

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) + b$$

このとき、係数 α_i の多くはゼロになる。これにより、分類器は訓練データのうち
一部の点（サポートベクトル）のみに依存する。

2.3.4 サポートベクトルと解のスパース性

SVM の双対問題を解くことで得られる解の特徴として、**スパース性（疎性）**があ
る。すなわち、学習によって得られる係数 α_i は、多くが 0 であり、非ゼロのものだ
けが分類器に寄与する。

このようなデータ点を「**サポートベクトル**」と呼び、分類境界に近い、あるいは誤
分類されている点に対応する。これにより、学習済みモデルはシンプルで効率的にな
り、予測計算も高速である。

まとめると：

- SVM はマージン最大化により汎化性能の高い分類器を構成する。
- 区分線形のヒンジ損失と 2 次正則化によって損失を定式化する。
- カーネル法によって非線形分類も可能となり、訓練データの一部（サポートベ
クトル）のみがモデルに寄与する。

2.3.5 カーネル SVM の双対定式化と凸最適化の概要

SVM は主問題 (primal problem) を双対問題 (dual problem) に変換することで、効率的に解くことができる。特にカーネルを用いた場合、双対問題の形式が重要になる。

ソフトマージン SVM の双対問題は、以下のように定式化される：

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C, \quad \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

この問題は 凸 2 次計画問題 (quadratic programming, QP) と呼ばれる (最適化分野で深く研究されている枠組み) であり、グローバルな最適解が保証される。

双対解 α を用いて、分類関数は以下のように表現される：

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) + b$$

この形式により、入力空間での非線形な分類境界を効率的に学習することができる。

2.3.6 SVM とロジスティック回帰の比較：損失関数の観点から

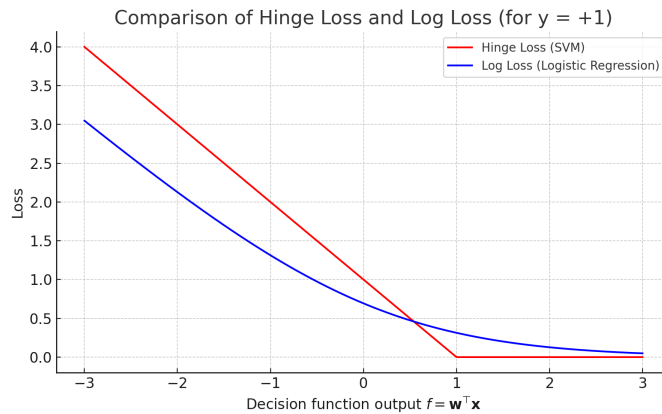


図 2.4 ヒンジ損失と対数損失

SVM とロジスティック回帰は、どちらも分類問題を扱うが、損失関数の設計に大きな違いがある。

- **ロジスティック回帰**：確率的モデルであり、ロジスティック関数を用いてクラスの後確率をモデル化する。損失関数は以下の **対数損失 (log loss)**：

$$\ell_{\log}(y, f) = \log(1 + \exp(-yf))$$

- **SVM**：マージン最大化を目的とし、損失関数は **ヒンジ損失**：

$$\ell_{\text{hinge}}(y, f) = \max(0, 1 - yf)$$

- **違いの直感**：ロジスティック損失は全ての誤差に滑らかにペナルティを与えるが、SVM のヒンジ損失は「マージンの外側」にある点（分類に十分余裕がある点）には損失を与えない。

ロジスティック回帰は出力を確率として解釈可能だが、SVM は分類境界に特化しており、確率の解釈は難しいという違いもある*2。

2.3.7 サポートベクトルと決定境界の直感的理解

SVM では、分類境界 (decision boundary) はマージンの幅を最大化するように決定される。このとき、分類器の決定に直接影響するのは、マージン上またはマージン内にあるデータ点だけである。

これらの点を「**サポートベクトル**」と呼び、それ以外の点（分類に余裕がある点）は分類境界に寄与しない。

- サポートベクトルは、分類平面からの距離が最も近い点。
- サポートベクトルの位置が変われば、分類境界も変わる。

*2 通常は正則化問題はベイズの MAP 解として解釈できるが、SVM の場合は確率分布で解釈しようとした時に式に現れない規格化因子の部分に最適化パラメータが入ってしまうことから純粋に確率としてみなすのは難しいとされている。

- 結果として、モデルは「データの重要な一部」にのみ依存するため、**スパース性**を持ち、計算効率も高くなる。

図示すると、分類境界の両側に平行なマージン線が引かれ、それに接する点がサポートベクトルとして強調される。

2.4 決定木とアンサンブル学習

2.4.1 決定木の定義と最適化

(327)

決定木 (Decision Tree) は、データの特徴ごとに条件分岐させて分類または回帰を行うモデルである。分類では「はい／いいえ」の質問を繰り返してデータを分類し、回帰では数値予測を行う。

決定木は、根 (root) から始まり、各内部ノードで特徴量に基づく条件分岐を行い、葉 (leaf) に到達したところで予測を出力する。

分割基準と目的関数：

ノードでの最適な分割を決めるためには、「純度 (purity)」の尺度を用いて評価を行う。代表的な指標は以下の通り：

- **ジニ不純度 (Gini impurity)：**

$$G(t) = \sum_{k=1}^K p_k(1 - p_k) = 1 - \sum_{k=1}^K p_k^2$$

- **エントロピー (情報利得)：**

$$H(t) = - \sum_{k=1}^K p_k \log p_k$$

ここで p_k はノード t におけるクラス k の出現確率である。

各ノードでは、特徴量としきい値の組み合わせを全探索し、**純度の改善 (不純度の減少)** が最大になる分割を選択する。

剪定 (pruning)：

木を成長させすぎると過学習するため、適切な深さで止める必要がある。これには：

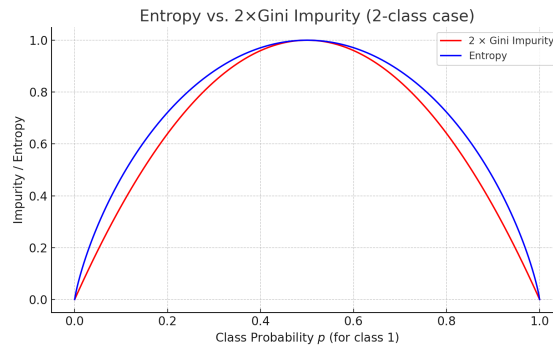


図 2.5 ジニ不順度とエントロピー

- 事前剪定 (max depth, min samples split など)
- 事後剪定 (学習後に不要な枝を削除)

がある。

2.4.2 具体例：Breast Cancer Wisconsin データによる識別問題

データセットの概要

Breast Cancer Wisconsin (Diagnostic) Dataset は、乳がん診断に関する特徴量を用いて、腫瘍が良性か悪性かを識別する分類問題である。UCI Machine Learning Repository で広く使われているベンチマークデータセットの一つである。

- サンプル数：569
- 特徴量数：30（細胞核の形状・テクスチャ・大きさ・滑らかさなど）
- クラスラベル：
 - 0：悪性 (malignant)
 - 1：良性 (benign)

すべての特徴量は数値（実数）であり、スケーリング処理を行うことで多くの分類器に適用可能となる。

決定木の結果

表 2.1 Breast Cancer Wisconsin データセットの特徴量一覧（インデックス付き）

インデックス	カテゴリ	特徴量名（日本語訳）
0	mean	radius_mean（半径の平均）
1	mean	texture_mean（テクスチャの平均）
2	mean	perimeter_mean（周囲長の平均）
3	mean	area_mean（面積の平均）
4	mean	smoothness_mean（平滑度の平均）
5	mean	compactness_mean（緻密さの平均）
6	mean	concavity_mean（凹度の平均）
7	mean	concave_points_mean（凹点数の平均）
8	mean	symmetry_mean（対称性の平均）
9	mean	fractal_dimension_mean（フラクタル次元の平均）
10	se	radius_se（半径の標準誤差）
11	se	texture_se（テクスチャの標準誤差）
12	se	perimeter_se（周囲長の標準誤差）
13	se	area_se（面積の標準誤差）
14	se	smoothness_se（平滑度の標準誤差）
15	se	compactness_se（緻密さの標準誤差）
16	se	concavity_se（凹度の標準誤差）
17	se	concave_points_se（凹点数の標準誤差）
18	se	symmetry_se（対称性の標準誤差）
19	se	fractal_dimension_se（フラクタル次元の標準誤差）
20	worst	radius_worst（最大半径）
21	worst	texture_worst（最大テクスチャ）
22	worst	perimeter_worst（最大周囲長）
23	worst	area_worst（最大面積）
24	worst	smoothness_worst（最大平滑度）
25	worst	compactness_worst（最大緻密さ）
26	worst	concavity_worst（最大凹度）
27	worst	concave_points_worst（最大凹点数）
28	worst	symmetry_worst（最大対称性）
29	worst	fractal_dimension_worst（最大フラクタル次元）

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 from sklearn.datasets import load_breast_cancer
4 from sklearn.tree import DecisionTreeClassifier,
    plot_tree
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import classification_report,
    confusion_matrix, accuracy_score
7
8 # データ読み込み
9 data = load_breast_cancer()
10 X, y = data.data, data.target
11 feature_names = data.feature_names
12 target_names = data.target_names
13
14 # 訓練・テストデータに分割
15 X_train, X_test, y_train, y_test = train_test_split(X, y
    , stratify=y, test_size=0.2, random_state=42)
16
17 # 決定木モデルの構築（深さ無制限）
18 tree_clf = DecisionTreeClassifier(random_state=42)
19 tree_clf.fit(X_train, y_train)
20
21 # テストデータで予測・評価
22 y_pred = tree_clf.predict(X_test)
23 print("Accuracy:", accuracy_score(y_test, y_pred))
24 print("\nClassification Report:")
25 print(classification_report(y_test, y_pred, target_names
    =target_names))
```

```
26
27 # 混同行列の可視化
28 cm = confusion_matrix(y_test, y_pred)
29 plt.figure(figsize=(5, 4))
30 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
              xticklabels=target_names, yticklabels=target_names)
31 plt.xlabel('Predicted')
32 plt.ylabel('Actual')
33 plt.title('Confusion Matrix (Decision Tree)')
34 plt.tight_layout()
35 plt.show()
36
37 # 決定木の可視化（最初の階層）3
38 plt.figure(figsize=(16, 8))
39 plot_tree(tree_clf, feature_names=feature_names,
            class_names=target_names, filled=True, max_depth=3)
40 plt.title("Decision Tree (First 3 Levels)")
41 plt.tight_layout()
42 plt.show()
```

2.4.3 アンサンブル学習の基本的な考え方

アンサンブル学習（Ensemble Learning）は、複数のモデル（基学習器）を組み合わせ、個々のモデルよりも高い精度を実現する手法である。

バギング（Bagging：Bootstrap Aggregating）

バギングは、同じ学習アルゴリズムを用いて、元の訓練データから複数のサブセット（ブートストラップ標本）を作成し、それぞれに学習を行う。

- 各モデルは独立に学習される（並列処理が可能）

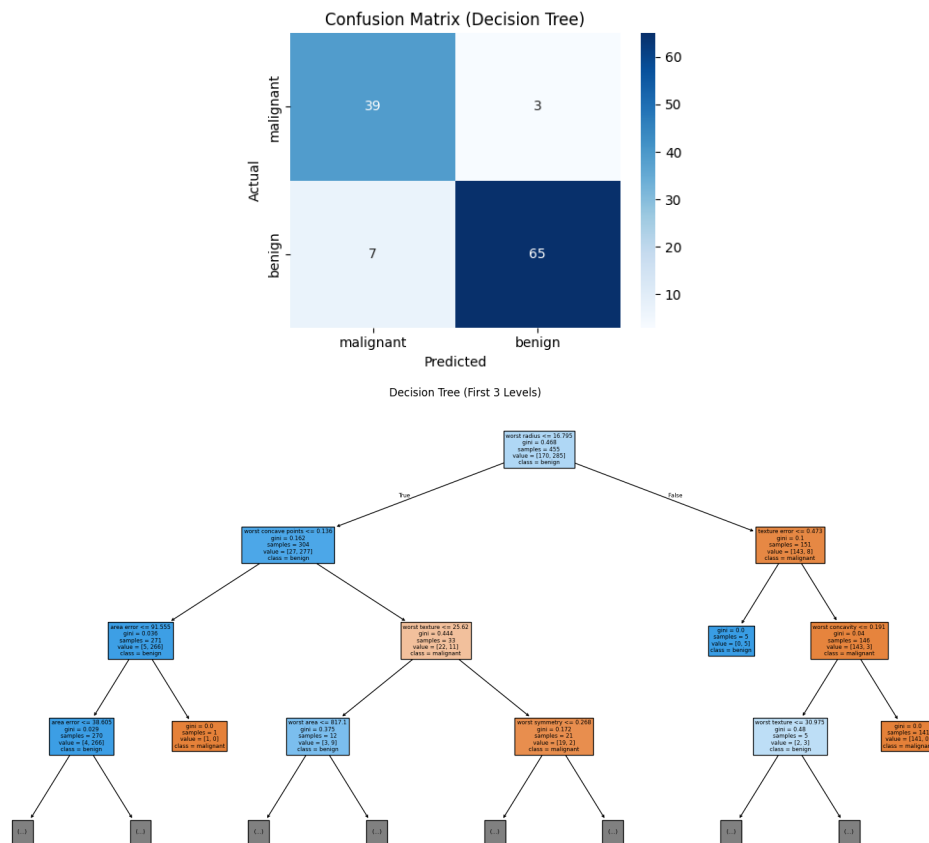


図 2.6 決定木

- 最終的な予測は分類では多数決，回帰では平均をとる
- モデルの分散を低減し，過学習を防ぐ効果がある

代表例：ランダムフォレスト（次節で詳述）

ブースティング (Boosting)

ブースティングは，複数の弱い学習器（性能のよくないモデル）を逐次的に学習させて，それらを組み合わせて性能の良い強い学習器を構成する方法である。

- 各ステップでは，前のモデルの誤りを重視して新しいモデルを学習する
- 予測は加重平均（または加重多数決）で行われる
- バイアスを低減する効果があるが，過学習に注意が必要

代表例：勾配ブースティング（次節で詳述）

バギング vs ブースティング（まとめ）：

- バギング：高分散なモデルに対して効果的（例：決定木）
- ブースティング：高バイアスな問題に対して効果的（例：線形的に分離困難なデータ）
- バギングは独立に学習，ブースティングは逐次的に学習

2.4.4 ランダムフォレスト

1. バギング（Bagging）の具体的手順と特徴

■手順：

1. 学習データからブートストラップ（復元抽出）によって B 個のデータセットを作成
2. 各データセットに学習器（例：決定木）を適用し、 B 個のモデルを構築
3. 新しい入力に対して、 B 個のモデルの出力を平均（回帰）または多数決（分類）して予測値とする

■特徴：

- 高分散な学習器に対して特に有効（例：深い決定木）
- 各学習器は独立であり、並列実行が可能
- バイアスは変化させず、分散の削減によって汎化誤差を改善する

2. ランダムフォレスト（Random Forest）のアルゴリズム

ランダムフォレストは、バギングを基盤としつつ、各決定木の構築時に特徴量のランダム性を導入することで、学習器間の相関を下げ、さらなる性能向上を実現した手法である（Breiman, 2001）。

■追加の工夫：

- 各ノードの分割時に、全特徴量からランダムに選んだ m 個の特徴量のみを分

割候補とする（通常 $m \ll d$ ）

- これにより、木同士の構造の多様性が高まり、**モデル間の相関が低減**する

■アルゴリズム：

1. バギングと同様に、 B 個のブートストラップサンプルを生成
2. 各サンプルに対して、次のルールで決定木を構築：
 - 各ノードでランダムに m 個の特徴量を選択
 - その中で最も良い分割を選んで分岐
 - 完全に成長させる（剪定はしない）
3. 新しい入力に対して、 B 個の木の予測を多数決（分類）または平均（回帰）で統合

■特徴：

- 高精度かつ頑健なモデルで、パラメータチューニングの必要が少ない
- **変数重要度の評価（feature importance）** が可能で、解釈性もある程度確保
- 木の本数 B を大きくすることで安定性が向上
- 高次元データや欠損データにも強く、実用的なアンサンブル学習手法の一つ

まとめ

- バギングは、複数のモデルを用いた平均化により分散を下げるアンサンブル法である。
- ランダムフォレストは、バギングにランダムな特徴選択を加え、モデルの多様性と性能をさらに高めた。
- 実用的かつ頑健な手法であり、分類・回帰ともに広く利用されている。

2.4.5 勾配ブースティングの具体的アルゴリズム（回帰：二乗誤差の場合）

ここでは、損失関数として二乗誤差を用いる回帰問題を考え、勾配ブースティングの具体的な学習ステップを説明する。

■設定 訓練データ： $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$

目的：予測関数 $f(\mathbf{x})$ を学習し、出力 $f(\mathbf{x}_i)$ が y_i に近づくようにする。

損失関数：

$$L(y, f(\mathbf{x})) = \frac{1}{2}(y - f(\mathbf{x}))^2$$

■アルゴリズム（学習ステップ）

1. 初期化：

$$f_0(\mathbf{x}) = \arg \min_c \sum_{i=1}^n \frac{1}{2}(y_i - c)^2 = \frac{1}{n} \sum_{i=1}^n y_i$$

（全体の平均で初期化）

2. **for** $m = 1, 2, \dots, M$ （反復ステップ） **do**：

（a）残差の計算（負の勾配）：

$$r_i^{(m)} = - \left. \frac{\partial L(y_i, f(\mathbf{x}_i))}{\partial f} \right|_{f=f_{m-1}} = y_i - f_{m-1}(\mathbf{x}_i)$$

（b）残差に対する回帰木の学習：回帰木 $h_m(\mathbf{x})$ を学習し、 $r_i^{(m)} \approx h_m(\mathbf{x}_i)$ を満たすようにする。

（c）モデルの更新：

$$f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \eta \cdot h_m(\mathbf{x})$$

（ $\eta \in (0, 1]$ は学習率。小さくすると学習は遅いが過学習しにくくなる。）

3. 最終モデル：

$$f_M(\mathbf{x}) = f_0(\mathbf{x}) + \sum_{m=1}^M \eta \cdot h_m(\mathbf{x})$$

■ポイント

- 各ステップでは、「現在のモデルがうまく予測できていない誤差（残差）」を次のモデルで修正していく。
- 二乗誤差の場合、勾配は単に残差 $y - f(\boldsymbol{x})$ になるため、非常に直感的。
- 弱学習器（例：浅い決定木）を多数重ねることで、複雑で精度の高いモデルが得られる。

2.4.6 XGBoost：実用的な勾配ブースティングの高速実装

XGBoost (eXtreme Gradient Boosting) は、勾配ブースティングを高速かつ高精度に実装した実用的な手法であり、機械学習コンペティションなどで広く用いられている。

XGBoost の主な特徴は以下の通り：

- **2次近似による目的関数の最適化**：勾配（1次導関数）だけでなく、ヘッセ行列（2次導関数）も用いることで、より安定かつ効率的に木の分割を決定する。
- **正則化項の導入**：各木の複雑さ（ノード数や葉の重み）に対するペナルティを目的関数に加え、過学習を防ぐ：

$$\Omega(f) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

（ T ：葉の数， w_j ：葉の出力値）

- **欠損値処理や並列計算の最適化**：自動的に欠損値に対応し、高速な分割探索を並列で行う。
- **扱いやすさと拡張性**：回帰・分類に加え、ユーザー定義の損失関数やカスタム目的関数も利用可能。

要するに、XGBoost は「精度」「速度」「汎用性」を兼ね備えた勾配ブースティングの実装として、実践的な場面で非常に有効な手法である。

複雑な識別器の比較

Breast cancer データセットに対して、以下の代表的な分類モデルを適用し、性能を比較する。

- 正則化ロジスティック回帰：L2 正則化付きの線形分類器
- SVM (RBF カーネル)：非線形分類に対応するサポートベクターマシン
- ランダムフォレスト：多数の決定木によるバギング型アンサンブル
- 勾配ブースティング：残差を逐次的に学習するブースティング型アンサンブル

ハイパーパラメータと交差検証による評価

各モデルに対して、主要なハイパーパラメータ（例：正則化強度、木の深さ、学習率など）を変化させ、**5 分割交差検証 (5-fold cross-validation)** により汎化性能を評価する。

評価指標には、**ROC AUC (受信者操作特性曲線の下面積)** を用いる。これはク (394)
ラス不均衡にも頑健であり、識別性能の全体的な比較に適している。

Program 2.4 識別手法の比較

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 from sklearn.datasets import load_breast_cancer
5 from sklearn.model_selection import cross_val_score,
   StratifiedKFold
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.svm import SVC
9 from sklearn.ensemble import RandomForestClassifier,
   GradientBoostingClassifier
10
```

```

11 # データ読み込み
12 data = load_breast_cancer()
13 X, y = data.data, data.target
14 feature_names = data.feature_names
15
16 # 標準化
17 scaler = StandardScaler()
18 X_std = scaler.fit_transform(X)
19
20 # 交差検証設定
21 cv = StratifiedKFold(n_splits=5, shuffle=True,
    random_state=42)
22
23 # 各モデルとハイパーパラメータ設定
24 models = {
25     'LogisticRegression': [
26         LogisticRegression(C=c, penalty='l2', solver='
            liblinear') for c in [0.01, 0.1, 1, 10, 100]
27     ],
28     'SVM (RBF)': [
29         SVC(C=c, gamma=g, kernel='rbf', probability=True
            ) for c in [0.1, 1, 10] for g in [0.01, 0.1,
                1]
30     ],
31     'RandomForest': [
32         RandomForestClassifier(n_estimators=100,
            max_depth=d, random_state=42) for d in [2,
                4, 6, 8, None]
33     ],

```

```
34     'GradientBoosting': [  
35         GradientBoostingClassifier(n_estimators=100,  
36             learning_rate=lr, max_depth=d, random_state  
37             =42)  
38     ]  
39 }  
40 # 結果保存  
41 results = []  
42  
43 print("=== Cross-Validation Performance (ROC AUC) ===")  
44 for name, model_list in models.items():  
45     for model in model_list:  
46         scores = cross_val_score(model, X_std, y, cv=cv,  
47             scoring='roc_auc')  
48         results.append((name, str(model.get_params()),  
49             np.mean(scores)))  
50  
51     print(f"{name:20s} | mean AUC: {np.mean(scores)  
52         :.4f}")  
53  
54 # グラフ化 (平均スコアの比較)  
55 import pandas as pd  
56  
57 df_results = pd.DataFrame(results, columns=["Model", "  
58     Params", "ROC AUC"])  
59  
60 plt.figure(figsize=(10, 5))  
61 sns.boxplot(data=df_results, x="Model", y="ROC AUC")  
62 plt.title("Model Comparison (ROC AUC, 5-fold CV)")
```

```

57 plt.xticks(rotation=20)
58 plt.tight_layout()
59 plt.show()

```

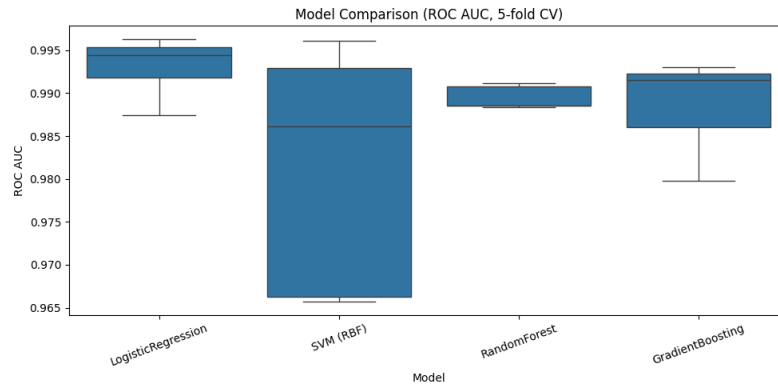


図 2.7 識別手法の比較

2.4.7 アンサンブル学習における変数重要度とモデルの解釈性

■動機：説明性と解釈性の重要性 機械学習モデルが高い精度を持つことは重要だが、なぜその予測がなされたのかを理解することも、実務では非常に重要である。特に以下のような文脈で、モデルの説明性 (explainability) や解釈性 (interpretability) が求められる：

- 医療・金融・法務などの高リスク分野における意思決定の裏付け
- モデル改善のための特徴量の選定や理解
- ステークホルダーへの結果の説明

このような背景から、どの特徴量が予測に大きく寄与したかを評価する指標として

(343) 「変数重要度 (feature importance)」が注目されている*³。

*³ ISL では variance importance という指標が説明されている。これは識別問題の場合、後で出てくる「分割による不純度減少の累積」の方法と同じ。

2.4.8 変数重要度の一般的な考え方

変数重要度とは、各特徴量が予測に与える寄与の大きさを定量的に評価したものである。一般的な評価方法には以下がある：

- **摂動ベース (permutation importance)**：特定の特徴量の値をランダムに入れ替えたときに、モデルの性能がどれだけ低下するかを見る。
- **SHAP (SHapley Additive exPlanations)**：単一の予測に対して、各特徴量がどれだけ予測値に寄与したかを、ゲーム理論に基づく公正な方法で分配する。

SHAP は以下のような特徴を持つ：

- － 予測単位での寄与の分解が可能
- － 全体で平均すればグローバルな重要度になる
- － アンサンブル木モデルとの相性が良く、高速な専用実装 (TreeSHAP) が存在

2.4.9 決定木アンサンブルに特化した重要度の評価

ランダムフォレストや勾配ブースティングなどの決定木ベースのアンサンブル学習モデルでは、モデルの構造を利用した変数重要度の指標が自然に得られる。

■1. 分割による不純度減少の累積 (Gain ベース) 各特徴量について、その特徴でノードを分割した際に得られる不純度の減少量 (例：ジニ不純度の減少、エントロピーの減少) を木全体にわたって合計する：

$$\text{Importance}(x_j) = \sum_{\text{分割に使われたノード}} \text{不純度の減少量}$$

この指標は簡単に得られ、高速に計算可能だが、以下のような注意点がある：

- 値のスケールやカテゴリ数の多さに敏感 (偏りが出ることもある)
- 他の特徴と相関がある場合、重要度が過小または過大評価されることがある

Program 2.5 変数重要性

```
1 # ランダムフォレストと勾配ブースティングの学習
2 rf = RandomForestClassifier(n_estimators=100, max_depth
    =4, random_state=42)
3 gb = GradientBoostingClassifier(n_estimators=100,
    learning_rate=0.1, max_depth=3, random_state=42)
4
5 rf.fit(X_std, y)
6 gb.fit(X_std, y)
7
8 # 特徴量重要度の表示
9 plt.figure(figsize=(10, 4))
10 plt.barh(feature_names, rf.feature_importances_)
11 plt.title("Random Forest Feature Importance")
12 plt.tight_layout()
13 plt.show()
14
15 plt.figure(figsize=(10, 4))
16 plt.barh(feature_names, gb.feature_importances_)
17 plt.title("Gradient Boosting Feature Importance")
18 plt.tight_layout()
19 plt.show()
```

■2. TreeSHAP による厳密な貢献度推定 TreeSHAP は、SHAP のアイデアを決定木アンサンブルに特化して効率的に計算する手法である。各予測に対して、特徴量の貢献度を定量的に割り当てることができる。

特徴：

- 各特徴量の影響を厳密に分離して評価
- モデル全体の平均で集約すれば、グローバルな重要度が得られる

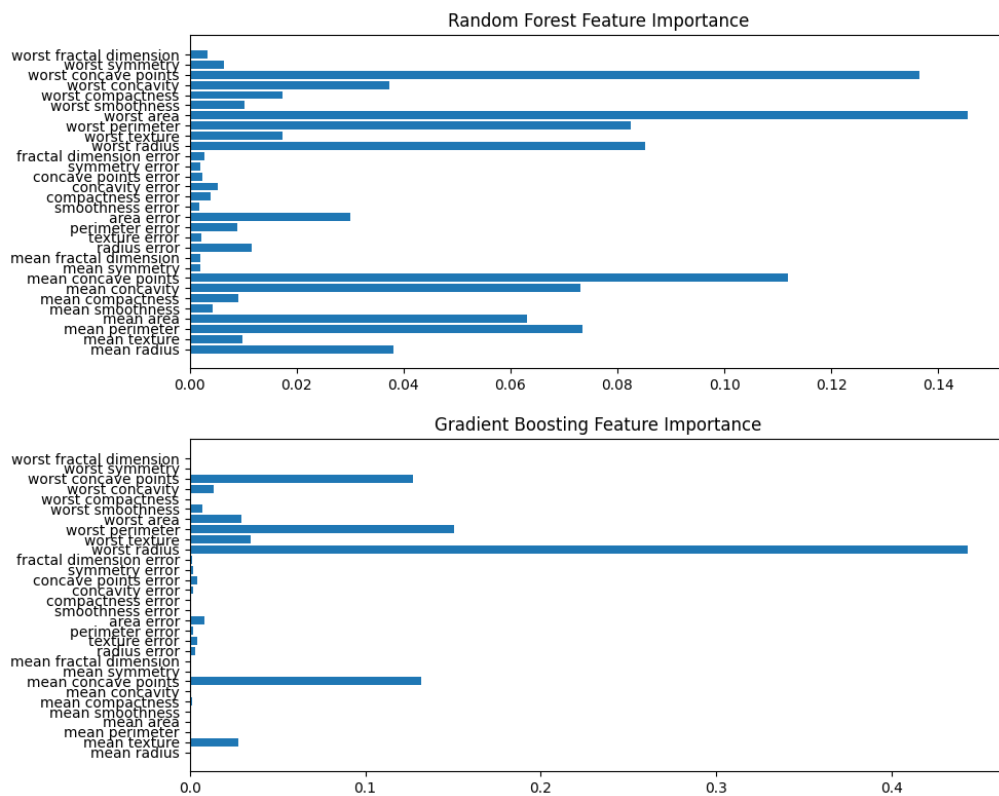


図 2.8 変数重要性

- 実行速度も高く、大規模データにも適用可能

■ まとめ

- 決定木アンサンブルでは、構造情報を活かして変数重要度を定量化できる。
- より正確かつ局所的な解釈には SHAP（特に TreeSHAP）が有効。
- モデルの解釈性は、単なる精度評価を超えて、実用において重要な視点である。

2.4.10 補足：Shapley 値の数学的定義

SHAP における特徴量の貢献度の評価は、ゲーム理論における **Shapley 値** (Shapley value) に基づいている。

Shapley 値は、協力ゲームにおいて、各プレイヤー（ここでは特徴量）が最終的な成果（ここでは予測値）にどれだけ貢献したかを、公平に分配するための理論である。

Shapley 値の定義：

特徴量の集合を $N = \{1, 2, \dots, d\}$ とし、任意の部分集合 $S \subseteq N \setminus \{j\}$ に対する貢献関数（予測値の期待値）を $v(S)$ とする。

特徴量 j に対する Shapley 値 ϕ_j は次のように定義される：

$$\phi_j = \sum_{S \subseteq N \setminus \{j\}} \frac{|S|!(|N| - |S| - 1)!}{|N|!} (v(S \cup \{j\}) - v(S))$$

ここで：

- $v(S)$ ：特徴量集合 S の情報のみを使ったときの予測値の期待値
- $v(S \cup \{j\}) - v(S)$ ：特徴量 j を追加することで得られる **限界貢献**
- 係数は全ての特徴量の順列における順序に基づく重み（公平性の保証）

この定義は、以下の公正性を満たすことが知られている：

- **効率性**：すべての特徴量の貢献度の合計が予測値と一致する
- **対称性**：同じ貢献をする特徴量には同じ値が割り当てられる
- **ダミー性**：貢献しない特徴量の値はゼロになる
- **加法性**：複数モデルに対して Shapley 値を加えると全体の貢献が合計される

SHAP はこの Shapley 値を機械学習の予測に適用し、各予測に対する特徴量の寄与を正確に定量化する方法として利用されている。

Program 2.6 SHAP

```

1 # SHAP 値の可視化
2 import shap
3
4 # Gradient Boosting に対して SHAP 実行 () TreeExplainer
5 explainer = shap.Explainer(gb, X_std)
6 shap_values = explainer(X_std)
7
8 # サマリープロット（グローバルな重要度）

```

```
9 shap.summary_plot(shap_values, X_std, feature_names=  
    feature_names)  
10  
11 # ウォーターフォールプロット（局所的説明）  
12 shap.plots.waterfall(shap_values[0], max_display=10)
```

2.5 多クラス分類と誤り訂正符号による拡張

2.5.1 多クラス分類の基本的アプローチ

多クラス分類とは、ラベルが2つ以上ある分類問題（例：クラス数 $K \geq 3$ ）を扱うものである。多くの機械学習アルゴリズム（特に識別器）は2クラス分類を前提としているため、**多クラス分類に拡張するための方法**が重要となる。

(385)

2.5.2 One-vs-Rest (OvR)

One-vs-Rest（あるいは One-vs-All）は、各クラス k に対して「クラス k vs その他すべて」の2クラス分類器を学習する方法である。

- 合計 K 個の分類器を学習
- 予測時には、すべての分類器の出力を比較し、最もスコアの高いクラスを選択

利点：

- 実装が簡単で汎用的な2値分類器をそのまま利用可能

欠点：

- クラスの不均衡に弱い（ある vs 多数）
- 各分類器のスコアが比較可能であることが前提となる

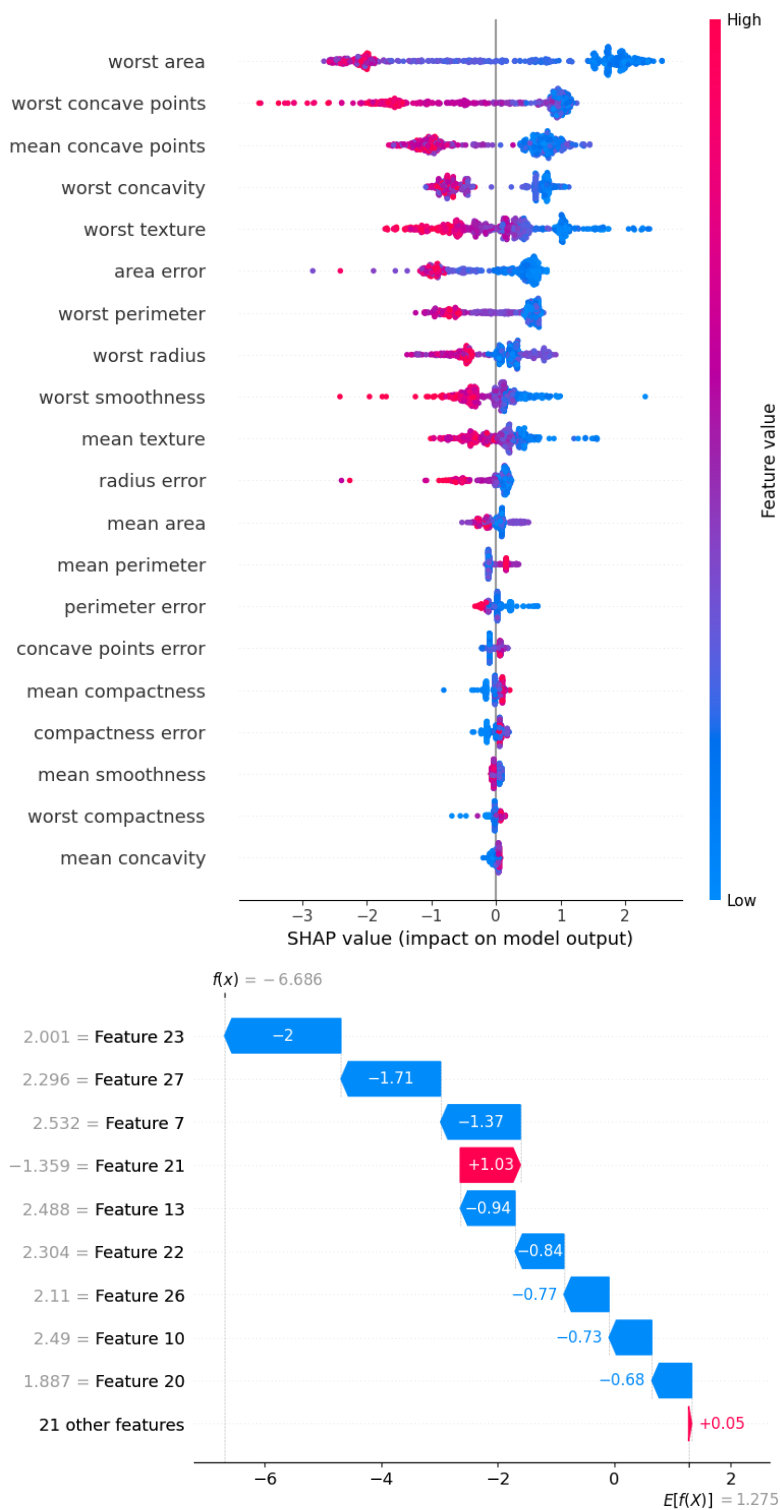


図 2.9 SHAP

2.5.3 One-vs-One (OvO)

One-vs-One は、すべてのクラスのペアに対して分類器を学習する方法である。クラス数が K のとき、分類器の数は：

$$\frac{K(K-1)}{2}$$

- 予測時には、全ての分類器で多数決を行う

利点：

- 各分類器は 2 つのクラスだけに集中できるため、学習がしやすい
- クラス間の分離が鮮明な場合に有効

欠点：

- 分類器数が多くなり、学習・推論コストが高い
- 投票が同数になった場合の扱いが不明確

2.5.4 ECOC（誤り訂正出力符号：Error-Correcting Output Codes）

ECOC は、多クラス分類を複数の 2 値分類問題に分解するための汎用的なフレームワークであり、誤り訂正符号（error-correcting codes）のアイデアに基づいている。
[7]

■基本的な考え方 各クラスに対して、長さ L のバイナリ符号（コードワード）を割り当て、それぞれのビットごとに 2 クラス分類器を学習する。

- 各ビット位置に対して、全体を 2 つのグループ（ラベル）に分けて分類器を学習
- 新しいデータに対する予測は、 L 個の分類器の出力ビット列と各クラスのコードワードの距離（例：Hamming 距離）で比較して、最も近いクラスを選ぶ

■誤り訂正符号としての性質 ECOC の核心は、「コード間距離を大きく取ることで誤分類を訂正できる」ことである。具体的には：

- クラスのコード間の Hamming 距離が d 以上ある場合、 $\lfloor (d-1)/2 \rfloor$ 個までの分類器の誤りを訂正できる
- よって、分類器にある程度の誤りが含まれても、最終予測は正しくなる可能性が高い

■利点と応用

- 複雑なクラス間関係を捉える柔軟な設計が可能
- クラス数が多い場合でも、誤りに対して頑健な分類が可能
- OvR や OvO も ECOC の特別な場合とみなすことができる

■まとめ

- 多クラス分類は2値分類器を工夫して組み合わせることで実現可能
- OvR と OvO は簡便かつ実用的だが、それぞれトレードオフがある
- ECOC は符号理論の考え方を応用し、頑健性の高い多クラス分類手法として汎用的に使える

2.5.5 補足：誤り訂正符号としての数式的性質

ECOC は、符号理論における誤り訂正符号と同様に、「クラスのコードワード同士の距離」を利用して分類誤りの訂正を行う。

■定義：Hamming 距離 2つのバイナリ列 $\mathbf{c}_1, \mathbf{c}_2 \in \{0, 1\}^L$ の間の Hamming 距離 d_H は、異なるビットの数である：

$$d_H(\mathbf{c}_1, \mathbf{c}_2) = \sum_{i=1}^L \mathbb{I}[c_{1i} \neq c_{2i}]$$

■誤り訂正可能な条件 すべてのクラス間のコードワードの最小 Hamming 距離を d_{\min} としたとき, ECOC は最大で

$$t = \left\lfloor \frac{d_{\min} - 1}{2} \right\rfloor$$

個の分類器の誤りを訂正できる。これは、符号理論における誤り訂正の基本原則と一致している。

この性質により、分類器の一部が誤った予測をしても、最終的な予測結果は正しいクラスに復元される可能性が高まる。

2.5.6 補足：不均衡データにおける多クラス分類の難しさ

多クラス分類では、クラス間のデータ数の偏り（不均衡性）が大きな問題となる。これは以下のような影響を及ぼす：

- **学習のバイアス**：多数派クラスに対して偏った分類器が学習されやすい
- **スコアの不公平性**：特に One-vs-Rest では、正例と負例の数の比率が極端になる
- **評価指標の選択が重要**：単純な精度（accuracy）では性能を正しく測れない

このような問題に対処するためには：

- **データの再サンプリング**（オーバーサンプリング／アンダーサンプリング）
- **重み付き損失関数**（クラスごとの誤分類に異なる重みを設定）
- **適切な評価指標**（F1 スコア, AUC, マクロ平均など）

が必要となる。多クラス分類では、単に分類器を構築するだけでなく、**クラス不均衡の影響を考慮した学習と評価**が重要である。

2.6 最近傍識別器 (k -nearest neighbor 法)

2.6.1 基本的な定義

(181)

最近傍法 (nearest neighbor classifier) は、教師あり学習の中でも非常に単純かつ直感的な識別法である。分類対象となる入力ベクトル $\mathbf{x} \in \mathbb{R}^d$ に対して、訓練データ内の「近くにある」サンプルのラベルに基づいてクラスを予測する。

一般に、 k 個の最近傍データ点を用いて多数決を行う手法を **k -最近傍法 (k -NN)** と呼ぶ。訓練データを $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ とすると、分類手順は以下の通りである：

1. 入力点 \mathbf{x} に対し、訓練データの中からユークリッド距離などの距離尺度に基づいて最も近い k 個の点を探索する。
2. それらの点に対応するラベル $\{y_{i_1}, \dots, y_{i_k}\}$ の多数決により、 \mathbf{x} のラベルを決定する。

$k = 1$ のときは **1-最近傍法** と呼ばれ、最も近い 1 点のラベルをそのまま割り当てる。

2.6.2 ベイズ識別器との性能比較

理論的には、ベイズ識別器（事後確率が最大のクラスを選ぶ識別器）が誤識別率の下限を与える。しかし、 k -NN 法は次のような興味深い性質を持つ：

- 訓練データ数 $n \rightarrow \infty$, $k \rightarrow \infty$, $k/n \rightarrow 0$ の極限で、 k -NN 法の誤識別率はベイズ識別器の誤識別率 R^* に対して

$$R_{kNN} \rightarrow R^* \quad (\text{一貫性 consistency})$$

- 特に、1-最近傍法 ($k = 1$) については次のクラシカルな不等式が知られている (Cover and Hart[8])：

$$R_{1NN} \leq 2R^*(1 - R^*)$$

この不等式は、1-NN 法が「ベイズ識別器の 2 倍以下の誤識別率」に抑えられることを保証する。

このように、 k -NN 法は非常に単純な非パラメトリック手法でありながら、理論的にも優れた性質を持つ*4.

2.6.3 計算アルゴリズム上の工夫

k -NN 法では、予測時に毎回すべての訓練データとの距離を計算する必要があり、計算コストが高くなる ($O(n)$). このため、以下のような高速化手法が考案されている:

- **kd-tree (k 次元木)** や **Ball-tree** による空間分割と近傍探索の効率化
- **近似最近傍探索 (Approximate Nearest Neighbor, ANN)**: 高速だが近似的な手法 (例: LSH, HNSW)
- **データの次元削減** (例: PCA) によって距離計算の次元数を削減する

また、次元が高くなると距離尺度の有効性が低下する (次元の呪い). このため、高次元データへの適用には特別な注意が必要である.

2.6.4 k -NN 法の特徴と注意点

- モデルの学習という概念はなく、すべての訓練データがそのまま予測に使われる.
- ハイパーパラメータ k の選択が性能に大きく影響する.
- データ数が多い場合には、計算・メモリのコストが高くなる.

シンプルで実装しやすく、かつ一定の理論的保証もあるため、分類タスクの初期ベースラインや比較対象として広く用いられている.

以下では Breast cancer データセットについて、 k -NN 法を適用する. 20% のテストデータに対する誤差をプロットする.

Program 2.7 kNN

```
1 import numpy as np
2 import matplotlib.pyplot as plt
```

*4 ただし、この不等式は訓練データ数が十分大きい場合の式であり、次元数が多い場合にはこの不等式を実現するのに次元数の指数オーダーのデータ数が必要と言われている.

```
3 from sklearn.datasets import load_breast_cancer
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import StandardScaler
6 from sklearn.neighbors import KNeighborsClassifier
7 from sklearn.metrics import accuracy_score
8
9 # データの読み込み
10 data = load_breast_cancer()
11 X, y = data.data, data.target
12
13 # 訓練・テストに分割（層化抽出）
14 X_train, X_test, y_train, y_test = train_test_split(
15     X, y, stratify=y, test_size=0.2, random_state=42
16 )
17
18 # 特徴量の標準化
19 scaler = StandardScaler()
20 X_train_std = scaler.fit_transform(X_train)
21 X_test_std = scaler.transform(X_test)
22
23 # を変化したときのテスト誤差を計算k
24 k_values = range(1, 31)
25 test_errors = []
26
27 for k in k_values:
28     model = KNeighborsClassifier(n_neighbors=k)
29     model.fit(X_train_std, y_train)
30     y_pred = model.predict(X_test_std)
31     test_error = 1 - accuracy_score(y_test, y_pred)
```

```
32     test_errors.append(test_error)
33
34 # 結果のプロット
35 plt.figure(figsize=(8, 5))
36 plt.plot(k_values, test_errors, marker='o')
37 plt.xlabel("Number of Neighbors (k)")
38 plt.ylabel("Test Error Rate")
39 plt.title("k-NN Classifier: Test Error vs. k (Breast
    Cancer Dataset)")
40 plt.grid(True)
41 plt.tight_layout()
42 plt.show()
```

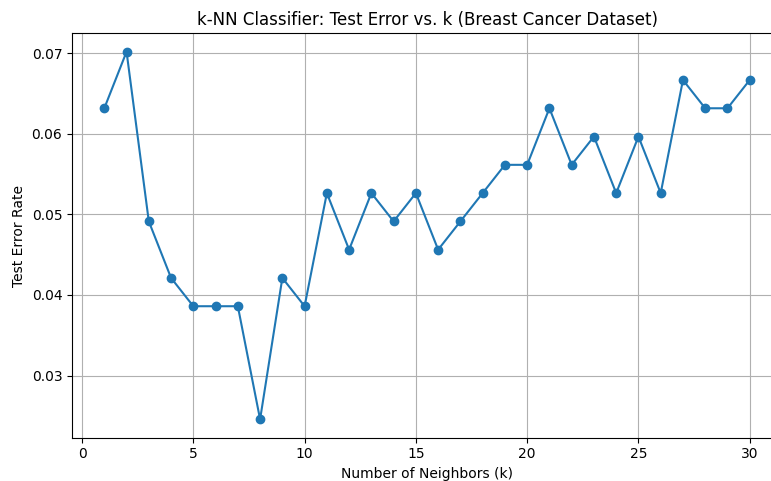


図 2.10 k-Nearest Neighbor 法

2.7 そのほかの教師あり学習モデル

教師あり学習については、本章で紹介した大まかな分類にも無限のバリエーションがあり、そのほか、ベイズモデルに基づくもの（ナীবベイズやガウス過程回帰）や深層学習など本章で扱わなかったモデルについては別の章で取り上げる。

(141) Scikit-learn にも ISL でも取り上げられているものとして Fisher 線形判別や 2 次判別 (Quadratic discriminant Analysis) がある. 多変量解析の文脈で統計学の枠内でも紹介されることが多いのでここでは省略した (正準相関分析や Partial Least Squares など).

(306) ISL で大きく取り上げられているが scikit-learn の基本構成には含まれておらず 認知度も低いものとして一般化加法モデル (GAM=Generalized Additive Model) がある.

参考文献

- [1] James, G., Witten, D., Hastie, T., Tibshirani, R. (2023) An Introduction to Statistical Learning, Second Edition, Springer. <https://www.statlearning.com>
- [2] Bishop, C.M. (2026) Pattern recognition and machine learning, Springer (ビショップ：パターン認識と機械学習（上下），丸善)
- [3] Kamishima, T., Akaho, S., Asoh, H., Sakuma, J. (2012). Fairness-aware classifier with prejudice remover regularizer. In Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2012, Bristol, UK, September 24-28, 2012. Proceedings, Part II 23 (pp. 35-50). Springer Berlin Heidelberg.
- [4] 萩原克幸 (2022), 入門 統計的回帰とモデル選択, 共立出版.
- [5] 赤穂昭太郎 (2008), カーネル多変量解析, 岩波書店.
- [6] Wahba, G. (1990). Spline models for observational data. Society for industrial and applied mathematics.
- [7] 金森敬文, 竹之内高志, 村田昇 (2009). パターン認識 (R で学ぶデータサイエンス 5) 共立出版.
- [8] Cover, T., Hart, P. (1967). Nearest neighbor pattern classification. IEEE transactions on information theory, 13(1), 21-27.
- [9] 青嶋誠, 矢田和善 (2019), 高次元の統計学, 共立出版
- [10] 小西貞則. (2010). 多変量解析入門: 線形から非線形へ. 岩波書店.
- [11] Hyvärinen, A. (2013). Independent component analysis: recent advances. Philosophical Transactions of the Royal Society A: Mathematical, Physical

- and Engineering Sciences, 371(1984), 20110534.
- [12] 赤穂昭太郎 (2018). ガウス過程回帰の基礎. システム/制御/情報, 62(10), 390-395.
- [13] Domingos, P., Pazzani, M. (1997). On the optimality of the simple Bayesian classifier under zero-one loss. Machine learning, 29, 103-130.
- [14] Yedidia, J. S., Freeman, W. T., Weiss, Y. (2001). Bethe free energy, Kikuchi approximations, and belief propagation algorithms. Advances in neural information processing systems, 13(24).
- [15] Srinivas, N., Krause, A., Kakade, S. M., Seeger, M. (2009). Gaussian process optimization in the bandit setting: No regret and experimental design. arXiv preprint arXiv:0912.3995.
- [16] S. Watanabe, M. Opper (2010), Asymptotic equivalence of Bayes cross validation and widely applicable information criterion in singular learning theory. Journal of machine learning research, 11(12).