### 授業準備:Webclassからコードをダウンロードし、 Google colaboratoryで開いておいてください

### 演習授業中の質問対応について



#### 2024/01/24 11:35-12:20



# 教師あり機械学習-決定木-



#### 関数 y=f(x) の f の設定次第で、機械学習にはさまざまなアルゴリズムがある



教師あり学習の種類

2.決定木

3. ランダムフォレスト



演習12

1.サポートベクターマシン

演習13

演習14

### 決定木(decision tree)

#### 決定木(decision tree)は特徴量をもとに条件分岐をすることで分類する手法





### 決定木(decision tree)

決定木の構成要素





# 乳がんのデータを用いて、分類を行う

### 機械学習の流れのまとめ



### 決定木のコードのまとめ



### 決定木のコードのまとめ





### STEP0: Google Colaboratoryの立ち上げ

 STEP0:事前準備

 STEP1:データの用意

 STEP2:学習モデルの選択

 STEP3:データを入れて学習

 STEP4:決定木の図示

 STEP5:予測を行う

 STEP6:モデルの評価

Python基礎 プログラミング基礎

#### 検索google colab Colaboratory へようこそ - Colaboratory - Google

	Colaboratory へようこそ	aboratory へようこそ バル 編集 表示 挿入 ランタイム ツール ヘルプ				
	ファイル 編集 表示 挿入 ランタイム ツール					
	ノートブックを新規作成					
	ノートブックを開く Ctrl+0	一次白13コート.lpynb				
	ノートブックをアップロード	lab へようこそ				
{x} =		こ Colab をよくご存じの場合は、この動画でインタラクティブなF				
07	ドライブにコピーを保存	ドの履歴表示、コマンドパレットについてご覧ください。				
	コピーを GitHub Gist として保存	3 Cool Google				
	GitHub にコピーを保存	Colab Features				
	保存 Ctrl+S					
	ダウンロード -					
	印刷 Ctrl+P					
	Сс	Jabとは				





### STEP0:ライブラリのインポート

### コード13-1 ライブラリとモジュールをインポート

import numpy as np

import pandas as pd

import matplotlib.pyplot as plt

● いつも使うnumpyとpandasライブラリ、pyplotモジュールをインポートする

モジュールをインポート: import <u>(ライブラリ名)</u>.(モジュール名) as 省略形 matplotlib pyplot

STEPO:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価



コード13-2 乳がんデータの読み込み

from sklearn.datasets import load breast cancer

bc = load breast cancer(as frame = True)

今回は分類を行うため、breast\_cancer(乳がん)のデータを読み込む

 ライブラリ
 モジュール
 関数

 関数をインポート from sklearn.datasets import load\_breast\_cancer
 Import load\_breast\_cancer
 Import load\_breast\_cancer

 ● load\_breast\_cancer
 (as\_frame = True/False)
 の引数で、True(真)か

 False(偽)を選択する
 読み込むデータの形式を指定

 True → pdのデータフレーム型
 False→ numpy配列

STEP5:予測を行う

STEP6:モデルの評価



#### STEP1:データの用意

STEP0:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価

#### コード13-3 学習データと検証データに分割

\*最初から学習データセットと検証データセットに分けてモデルの構築と評価を行う

from sklearn.model\_selection import train\_test\_split

x\_train, x\_test, y\_train, y\_test = train\_test\_split(

bc.data, bc.target, test\_size = 0.3, random\_state = 0)

- trainデータ70%とtestデータ30%にデータを分割
- 分割時の乱数シード値を0に指定
- x\_bcとy\_bcを作成せず、bc.data, bc.targetをtrain\_test\_splitに 直接入れてデータ分割を行う
- SVMと同様に、特徴量は30個全て使うことができる





STEP2:学習モデルの選択

STEP0:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価

#### コード13-4 学習モデルに決定木を選択

from sklearn.tree import DecisionTreeClassifier

クラスだけをインポート from <u>sklearn</u>.<u>tree</u>import <u>DesicionTreeClasifier</u> ライブラリ名 モジュール名 クラス名

● 今回のモデルは決定木分析で分類を行うので、

DecisionTreeClassifier()を選択する

DecisionTreeClassifier()クラスからmodel\_tree\_cを新しいインスタン
 スとして作成する



- 引数criterion= で分類のアルゴリズムを指定する。今回は'gini'でジニ不純度が最小になるように分類するCARTアルゴリズム(必ず二つに分岐する構造がシンプルなモデル)を使う
- 引数random\_state=では乱数シード値を固定する。決定木の分割時に、特徴量がランダムに並び替えられたり、分割時にタイブレークの時にどれか一つの分割をランダムに選択する必要があるため。



#### 分類決定木の学習方法

回帰の決定木の場合は RSS を指標にして分割していくが、分類の決定木の場合にはジニ不純度 (gini impurity) をもとに分割方法を決定。 これはクラス分けがどれだけ綺麗かの指標であり、<mark>分類した時にどれだけ「不純なもの」が含まれるかを表す。</mark>

ジニ不純度 G は 
$$G = \sum_{k=1}^K p(k)(1-p(k))$$
 K はクラス数, p(k) はその領域でのクラス k の割合

例えば以下のようにデータを分割すると(赤と青の2クラスの例)



このように、混じり物があるとジニ不純度は大きくなる。 したがって、<mark>ジニ不純度が小さくなるように木の分割を考えていけばよい</mark>。



STEP3:データを入れて学習させる

コード13-5 学習用データで学習させる

model\_tree\_c.fit(x\_train, y\_train)



STEPO:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価

colar





- ●決定木を図示するためにplot\_tree() 関数をインポート
- ●<u>plt.figure(figsize =(xx,xx))</u>で図の大きさを指定で きる。指定しないと文字が潰れて読めないので、大きめの (25,15)で指定する

●plot\_tree(決定木モデル, feature\_names=, class\_names=, filled = ) で引数を指定する

### STEP4:決定木の図示

#### コード13-7 bcのfeature\_namesの内容確認

#### bc.feature\_names



'fractal dimension error', 'worst radius', 'worst texture', 'worst perimeter', 'worst area', 'worst smoothness', 'worst compactness', 'worst concavity', 'worst concave points', 'worst symmetry', 'worst fractal dimension'], dtype='<U23')

bc.feature names.shape

type(bc.feature names)

STEP0:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価

(30, )

<class 'numpy.ndarray'>



#### STEP4:決定木の図示

コード13-8 bcのtarget\_namesの内容確認

bc.target\_names

STEP0:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価

array(['malignant', 'benign'], dtype='<U9')</pre>

bc.target\_names.shape

type(bc.target\_names)



STEP4:決定木の図示

STEPO:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価

colat



- ●plot\_tree(決定木モデル, feature\_names=,
  - class\_names=, filled = )で引数を指定
- ●feature\_namesはbc.feature\_names(特徴量の名前の リスト), class\_namesはbc.target\_names(targetの クラス(良性腫瘍/悪性腫瘍))を指定
- ●引数<u>filled = True</u>を指定すると、カラーになりノードの
   クラスが良性腫瘍(青)/悪性腫瘍(橙)と異なる色で描出



STEPO:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価

リーフノードが1つのクラス(正解値1/0)になるまで分岐される



STEPO:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価

リーフノードが1つのクラス (正解値1/0)になるまで分岐される

### STEP4:決定木の図示

STEP0:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価



### STEP4:決定木の図示







STEP5:予測を行う

コード13-9 検証用データを用いて予測する

print(model\_tree\_c.predict(x\_test))

print(np.array(y\_test))



STEP0:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価



● confusion\_matrix(実際の分類,予測分類,labels=[])で混同行列を出力する。 labels引数=で出力の順番(1:良性腫瘍、0:悪性腫瘍)を指定する

STEP5:予測を行う

array([[97, 11], [4, 59]])

		r		
		予測結果		
		Positive(正)	Negative(負)	
実際の 分類結	Positive (正)	真陽性True Positive <b>97</b>	偽陰性False Negative <mark>11</mark>	
果	Negative (負)	偽陽性False Positive <b>4</b>	真陰性True Negative <b>59</b>	

STEP0:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価









#### ROC 曲線と AUC による分類モデルの良し悪しの比較

**ROC**(Receiver operating characteristic) 曲線はもともとレーダーの測定能力の評価に関する種々の研究や理論が 1970 年以降に医療や機械学習へ応用されたものだ。

さまざまなカットオフで感度・特異度を計算し、横軸に 1 - 特異度、縦軸に感度をプロットしたものだ。 左上が感度 1、特異度 1 という理想的な分類器を示す。

ROC 曲線の下側の面積が AUC (Area under curve) である。ROC 曲線が左上にあればあるほど AUC は 1 に近い値をとる。 この AUC の値をもとに、どちらの分類器が優れるかを判断できる。



(清水先生の講義引用)

STEP6:モデルの評価

### コード13-13 ROC曲線の描出

```
from sklearn.metrics import RocCurveDisplay
RocCurveDisplay.from_predictions(y_test,
    model_tree_c.predict_proba(x_test)[:,1])
plt.axis('square')
plt.show()
```

\*RocCurveDisplay.from\_prediction(実際の分類, Y=1になる 確率)でROC曲線を描出できる \*plt.axis(`square`)で正方形の図に指定 \*plt.show()で図を表示

STEP0:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価





# 剪定 pruning



# 剪定 pruning

決定木分析では、全て純粋なノード(良性腫瘍か悪性腫瘍のどちらか一つ) になるまで分類するため、過学習が起きやすい

<u>= 学習データに対しては100%の正解率になる</u>

**木の深さやノードの作り方を制限する**ことで過学習を抑えることができる


#### STEP2:学習モデルの選択

## DecisionTreeClasifier()の引数

- 引数max\_depth=で木の深さ(何回分類を行うか)の 最大値を指定する。深さの最大値を指定しないと、全 てのデータが分類されるまで分割を繰り返す。大きい と過学習の傾向になる
- 引数min\_samples\_split=では、分岐を作成するため
   に必要な最低サンプルサイズを指定。小さいと過学習の傾向になる

STEP0:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価







- DecisionTreeClassifier()クラスからmodel\_tree\_pを新しいイン スタンスとして作成する
- 引数max\_depth=3で最大3回の深さで分岐するように指定













● confusion\_matrix(実際の分類,予測分類,labels=[])で混同行列を出力する。
 引数labels=で出力の順番(1:良性腫瘍、0:悪性腫瘍)を指定する

STEP5:予測を行う

array([[103, 5], [ 4, 59]])

		予測結果		
		positive(正)	Negative(負)	
実際の 公類	positive (正)	真陽性True Positive <b>103</b>	偽陰性False Negative <b>5</b>	
果	Negativ e (負)	偽陽性False Positive <mark>4</mark>	真陰性True Negative <b>59</b>	

STEP0:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価







ロジスティック回帰のAUC:0.982363315696649 剪定していない決定木のAUC:0.9173280423280423

STEP6:モデルの評価

## コード13-20 ROC曲線の描出

```
from sklearn.metrics import RocCurveDisplay
RocCurveDisplay.from_predictions(y_test,
    model_tree_p.predict_proba(x_test)[:,1])
plt.axis('square')
plt.show()
```

\*RocCurveDisplay.from\_prediction(実際の分類, Y=1になる 確率)でROC曲線を描出できる \*plt.axis('square')で正方形の図に指定 \*plt.show()で図を表示







STEP6:モデルの評価

コード13-21 検証用データでPrecision(適合率)、 Recall (再現率)、F1値を算出 STEPO:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価

from sklearn.metrics import classification\_report

pd.DataFrame(classification\_report(y\_test,

model\_tree\_p.predict(x\_test), output\_dict = True))

		0	1	accuracy	macro avg	weighted avg
	precision	0.921875	0.962617	0.947368	0.942246	0.947607
	recall	0.936508	0.953704	0.947368	0.945106	0.947368
	f1-score	0.929134	0.958140	0.947368	0.943637	0.947453
	support	63.000000	108.000000	0.947368	171.000000	171.000000

#### STEP6:モデルの評価

STEPO:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:決定木の図示 STEP5:予測を行う STEP6:モデルの評価

#### 検証用データでPrecision(適合率)、Recall (再現率)、F1値を算出

	0	1	accuracy	macro avg	weighted avg
precision	0.921875	0.962617	0.947368	0.942246	0.947607
recall	0.936508	0.953704	0.947368	0.945106	0.947368
f1-score	0.929134	0.958140	0.947368	0.943637	0.947453
support	63.000000	108.000000	0.947368	171.000000	171.000000

		0	1	accuracy	macro avg	weighted avg
SVC	precision	0.980769	0.899160	0.923977	0.939964	0.929226
	recall	0.809524	0.990741	0.923977	0.900132	0.923977
	f1-score	0.886957	0.942731	0.923977	0.914844	0.922183
	support	63.000000	108.000000	0.923977	171.000000	171.000000

## 演習13:課題

Webclassで課題を提出してください。締め切りは2024/02/07 23:59まで

<u>DecisionTreeClassifier() で剪定を行い、STEP2~STEP6まで実行してください</u>。 max\_depth=3で作った決定木(正解率0.94736、AUC:0.94319)より分類性能(検証デー タにおける正解率かAUCのどちらか)が優れた決定木を作成するように、剪定をしてくだ さい。

剪定は max\_depth=, と min\_samples\_split= で自分で自由に数値を指定して行って ください。random\_state=0に指定してください。

 webclassにはあなたが作成した優れた決定木を作れる(1)max\_depthと (2)min\_samples\_splitで指定した数値を入力してください。
 (3)優れた指標は正解率かAUCのどちらか、(4)(3)で選んだ指標の値を入力してください
 \*正解率もAUCも優れたモデルになった場合は、AUCを記載してください

## 演習13:課題

Webclassで課題を提出してください。締め切りは2024/02/07 23:59まで

● 決定木について理解を問う選択問題が二問あります。正しい選択肢を一つwebclass で選択してください。

#### (5) 問題1:決定木の分割を行うのに使用される指標はなんですか?

- 1.AUC
- 2.ジニ不純度
- 3.平均二乗誤差
- 4. 適合率
- (6) 問題2:決定木の剪定の目的は何ですか?
  - 1.過学習を防ぐため
  - 2.未学習を防ぐため
  - 3. 学習時間を短縮するため
  - 4. 適切な特徴量を選ぶため

#### 授業準備:Webclassからコードをダウンロードし、 Google colaboratoryで開いておいてください

#### 演習授業中の質問対応について



#### 2024/01/25 10:40-11:25

# 医療とAI・ビッグデータ入門 演習14 教師あり機械学習 -ランダムフォレスト-

統合教育機構 石丸美穂

#### 関数 y=f(x) の f の設定次第で、機械学習にはさまざまなアルゴリズムがある



教師あり学習の種類

2.決定木

3. ランダムフォレスト





1.サポートベクターマシン

演習13



#### アンサンブル学習:「三人寄れば文殊の知恵」作戦

決定木は、解釈しやすいというメリットがある反面、精度はそれほど高くないというデメリットもあった。 単独で使うと精度が低い学習器のことを<mark>弱学習器</mark> (weak learner) というが、だからこそたくさんの弱学習器を集めて「多数決」なり 「平均」なり集団として使うことで精度を大きく上げることができる。これをアンサンブル (ensemble) 学習という。 (音楽の用語で二人以上が同時に演奏すること)

どのようにアンサンブルするか、によって、バギング、ブースティング、スタッキングという3つの考え方がある。 そのうち前者2つを取り上げる。

#### アンサンブルのやり方 1: バギング

バギング (bagging) は bootstrap aggregating の略。

ブートストラップ (bootstrap) というのは, 母集団から 重複を許してランダムにデータを取って標本にするやり方。

学習データから新たに学習データ群を複数作って それぞれの学習データ群でモデルを学習させる。



#### アンサンブルのやり方 2: ブースティング

<u>
「
ースティング</u> (boosting) は, バギングのように並列に複数のモデルを 学習するのではなく, 直列にモデルを学習していく。学習したモデルが うまく予測できなかった学習データに重みをつけてさらにモデルを学習 するというのを繰り返していくイメージ。



#### (清水先生の講義スライドより引用)

#### ランダムフォレスト

決定木のアンサンブル (正確にはバギング)を行ったのが、ランダムフォレスト (random forest) と呼ばれる手法である。

バギングは、ブートストラップ法を使ってサンプル抽出した複数のデータ群に対してそれぞれモデルを構築して、最後に平均や多数決を取る 方法だが、ランダムフォレストにはさらに工夫がある。

それは、それぞれの決定木において**ランダム**に選んだ一部の特徴量のみを使って分割を行うようにしているのだ ( これが名前の由来 )。 つまり,それぞれの決定木を構築する際に、「一部の学習データを使わない」と「一部の特徴量でしか分割を行わない」の 2 つを行うことで、 <u>少しずつ違う決定木</u>を作ろうとしている。

特徵量 一般的には、決定木の数は100、分割に使う 分割する特徴量 学習 X1 X3 X5 X6 特徴量の数は利用可能な特徴量数の平方根をとったもの。  $X_1 X_2 X_3 X_4 X_5 X_6$ 学習データ群1 しかしこれらの値はモデルを作る人が自分で設定するべき  $x_1 x_1 x_4 x_4 x_2 x_6$ ハイパーパラメーター (hyper parameter) である。 ブートストラップ  $X_2 X_3 X_5 X_6$ 学習データ群2  $x_2 x_4 x_5 x_2 x_1 x_1$ 元の学習データ  $x_1 x_2 x_3 x_4 x_5 x_6$  $X_1 X_3 X_4 X_6$ アンサンブル 学習データ群3  $x_1 x_4 x_1 x_5 x_6 x_2$ 決定木がたくさんあるので各決定木の詳細を見るのは難しいが、  $X_2 X_4 X_5 X_6$ 全体での特徴量の重要度を確認することができる。 学習データ群4 (feature importance)  $x_5 x_4 x_1 x_6 x_2 x_2$ (清水先生の講義スライドより引用)

#### ランダムフォレスト (random forest)



- すべてのデータから、n個のサン プリングデータセットを作成する
   n個の決定木を作成する(この時 に特徴量もランダムに選択)
   各決定木モデルで予測
- 4. 多数決で最終予測を行う

## ランダムフォレスト (random forest)

- 回帰・SVM・決定木分析は一つのモデルを作成していた
- 単一のモデルでどれが最適かを選ぶのは難しく、また精度がそれほど高くないこともある
- アンサンブル学習は複数のモデルを作成し、最終予測を行うことで精度 を高めることができる



## 乳がんのデータを用いて、分類を行う

## 機械学習の流れのまとめ



# ランダムフォレストのコードのまとめ

STEP1 データの用意 STEP2 学習モデルの選択(今回は決定木) (モデル名) = RandomForestClassifier() STEP3 データを入れて学習させる モデルは RandomForesteClassifier() (モデル名).fit(特徴量,予測値) STEP4 予測を行う: (モデル名).predict() STEP5 モデルの評価: (モデル名).score (特徴量,予測値) : roc auc score (実測値, y=1になる確率) 特徴量の重要度の図示 STEP6 特徴量重要度の図示を行う

## ランダムフォレストのコードまとめ

<pre>import numpy as np import pandas as pd import matplotlib.pyplot as plt</pre>	STEP0: ライブラリの読み込み					
<pre>from sklearn.datasets import load_brest_cancer bc = load_breast_cancer(as_frame = True)</pre>	STEP1 : データの準備					
<pre>from sklearn.model_selection import train_test_split x_train, x_test, y_train, y_test = train_test_split(bc.data, bc.target, test_size = 0.3, random_state = 0)</pre>						
<pre>from sklearn.ensemble import RandomForestclassifier model_forest = RandomForestClassifier(n_estimators = 100, ma random_state = 0)</pre>	<pre>STEP2. fetures = 5, mx_depth = 3, max_features = 5,</pre>					
<pre>model_forest.fit(x_train, y_train)</pre>	STEP3 : データを入れて学習					
<pre>from sklearn.metrics import confusion_matrix confusion_matrix(y_test, model_forest.predict(x_test), label</pre>	s = [1,0]) STEP4:予測					
<pre>print(model_forest.score(x_test, y_test))</pre>						
<pre>from sklearn.metrics import roc_auc_score roc_auc_score(y_test, modelforest.predict_proba(x_test)[:,1]</pre>	STEP5: モテルの評価					
<pre>forest_importances = pd.DataFrame(model_forest.feature_import bc.feature_names, columns = ['Importance']) forest_importances.plot.bar()</pre>	tances_, index = STEP6:特徵量重要度					



## STEP0: Google Colaboratoryの立ち上げ

STEPO:事前準備
STEP1:データの用意
STEP2:学習モデルの選択
STEP3:データを入れて学習
STEP4:予測を行う
STEP5:モデルの評価
STEP6:特徴量重要度の図示

Python基礎 プログラミング基礎

#### 検索google colab <u>Colaboratory へようこそ - Colaboratory - Google</u>

	Colaboratory へようこそ ファイル 編集 表示 挿入 ランタイム ツ	ール ヘルプ
<ul> <li>4</li> <li>5</li> <li>4</li> <li>6</li> <li>7</li> <li>1</li> </ul>	ノートブックを新規作成 ノートブックを開く C ノートブックをアップロード 名前の変更 ドライブにコピーを保存 コピーを GitHub Gist として保存	<ul> <li>ド + テキスト ▲ ドライブにコピー</li> <li>加ab へようこそ</li> <li>Ic Colab をよくご存じの場合は、この動画でインタラクティブな5</li> <li>ドの履歴表示、コマンドパレットについてご覧ください。</li> <li>3 Cool Google</li> </ul>
5	保存     C       変更履歴     ダウンロード       印刷     C	rri+s rri+P
		Colab とは







ランダムフォレスト	STEP0:事前準備 STEP1:データの用意 STEP2:学習モデルの選択	
STEP1:データの用意	STEP3:データを入れて学習 STEP4:予測を行う	
コード14-3 学習データと検証データに分割	STEP5:モデルの評価 STEP6:特徴量重要度の図示	

\*最初から学習データセットと検証データセットに分けてモデルの構築と評価を行う

from sklearn.model\_selection import train\_test\_split

x\_train, x\_test, y\_train, y\_test = train\_test\_split(

bc.data, bc.target, test\_size = 0.3, random\_state = 0)

- trainデータ70%とtestデータ30%にデータを分割
- 分割時の乱数シード値を"0"に指定
- x\_bcとy\_bcを作成せず、bc.data, bc.targetをtrain\_test\_splitに 直接入れてデータ分割を行う
- svM・決定木と同様に、特徴量は30個全て使うことができる



STEP2:学習モデルの選択



#### コード14-4 学習モデルにランダムフォレストを選択

from sklearn.ensemble import RandomForestClassifier

● 今回のモデルはランダムフォレストで分類を行うので、

RandomForestClassifier()を選択する

● RandomForestClassifier() クラスからmodel\_forestを新しいインスタン スとして作成する



● 引数 n estimators = 100 で作成する決定木の数を指定

引数 max depth = 4 で決定木の最大深さを指定

特徴量の数は利用可能な特徴量数の平方根をとったもの。

しかしこれらの値はモデルを作る人が自分で設定するべき

一般的には、決定木の数は100、分割に使う

ハイパーパラメーター (hyper parameter) である。

(清水先生の講義資料より引用)

● 引数 max features = で一つのモデルに入れる特徴量の最大値を指定

<u>√30 = 5.48 → 5 に指定</u>

# ランダムフォレスト

STEP2:学習モデルの選択

STEPO:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:予測を行う STEP5:モデルの評価 STEP6:特徴量重要度の図示

model\_forest = RandomForestClassifier(n\_estimators = 100,

max depth = 4, max features = 5, random state = 0)

● random state = 0 で乱数シード値を固定

木を構築するときに使用されるサンプルのブートストラップのランダム性と、各ノードで最良 の分割を探すときに考慮する特徴のサンプリングの両方に関係

\*指定していないハイパーパラメータはデフォルト設定のものを使う

- criterion = "gini"
- min\_samples\_split = 2


STEP4:予測を行う

コード14-7 検証用データを用いて予測する

print(model\_forest.predict(x\_test))

print(np.array(y\_test))

STEPO:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:予測を行う STEP5:モデルの評価 STEP6:特徴量重要度の図示

x\_testデータで予測した yが0/1かの値

y\_testの実際の値



● confusion\_matrix(実際の分類,予測分類,labels=[])で混同行列を出力する。
 labels引数=で出力の順番(1:良性腫瘍、0:悪性腫瘍)を指定する

STEP4:予測を行う

### array([[104, 4], [ 3, 60]])

STEPO:事前準備
STEP1:データの用意
STEP2:学習モデルの選択
STEP3:データを入れて学習
STEP4:予測を行う
STEP5:モデルの評価
STEP6:特徴量重要度の図示

			予測結果		
			Positive (1:良性腫瘍)	Negative (0:悪性腫瘍)	
	実際 の 分類 結果	Positive (1:良性腫瘍)	真陽性True Positive <b>104</b>	偽陰性False Negative <b>4</b>	
		Negative (0:悪性腫瘍)	偽陽性False Positive <mark>3</mark>	真陰性True Negative <mark>60</mark>	







#### ROC 曲線と AUC による分類モデルの良し悪しの比較

**ROC**(Receiver operating characteristic) 曲線はもともとレーダーの測定能力の評価に関する種々の研究や理論が 1970 年以降に医療や機械学習へ応用されたものだ。

さまざまなカットオフで感度・特異度を計算し、横軸に 1 - 特異度、縦軸に感度をプロットしたものだ。 左上が感度 1、特異度 1 という理想的な分類器を示す。

ROC 曲線の下側の面積が AUC (Area under curve) である。ROC 曲線が左上にあればあるほど AUC は 1 に近い値をとる。 この AUC の値をもとに、どちらの分類器が優れるかを判断できる。



(清水先生の講義引用)

STEP5:モデルの評価

### コード14-11 ROC曲線の描出

```
from sklearn.metrics import RocCurveDisplay
RocCurveDisplay.from_predictions(y_test,
    model_forest.predict_proba(x_test)[:,1])
plt.axis('square')
plt.show()
```

\*RocCurveDisplay.from\_prediction(実際の分類, Y=1になる 確率)でROC曲線を描出できる \*plt.axis(`square`)で正方形の図に指定 \*plt.show()で図を表示









STEP5:モデルの評価

STEPO:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:予測を行う STEP5:モデルの評価 STEP6:特徴量重要度の図示

#### コード14-12 検証用データでPrecision(適合率)、 Recall (再現率)、F1値を算出

from sklearn.metrics import classification\_report

pd.DataFrame(classification\_report(y\_test,

model\_forest.predict(x\_test), output\_dict = True))

		0	1	accuracy	macro avg	weighted avg
	precision	0.937500	0.971963	0.959064	0.954731	0.959266
	recall	0.952381	0.962963	0.959064	0.957672	0.959064
	f1-score	0.944882	0.967442	0.959064	0.956162	0.959130
	support	63.000000	108.000000	0.959064	171.000000	171.000000

# 特徴量の重要度 Feature Importance

### 特徴量の重要度 Feature Importance

#### <u>ランダムフォレストでは特徴量の重要度を算出できる</u>

特徴量重要度:Mean Decrease Impurity(MDI) その特徴量を使って分割した時の不純度の減少量の木々の平均

- 綺麗に分割できる特徴量ほど重要度が大きくなる、
- 値は特徴量全体から見た相対的な指標であり、すべての特徴量の重
   要度を足すと1になる



STEP6:特徴量重要度の図示

STEPO:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:予測を行う STEP5:モデルの評価 STEP6:特徴量重要度の図示

#### コード14-13 特徴量重要度を取得

print(model\_forest.feature\_importances\_)

[0.0095214 0.01263454 0.01393639 0.06960096 0.00217414 0. 0.12990861 0.1690554 0.00140434 0. 0.05101289 0. 0.00981748 0.1344274 0.00302646 0. 0.00130498 0.00224143 0.00111216 0.00644376 0.11146473 0.00258438 0.08518957 0.00120382 0. 0.00445832 0.01797664 0.14089711 0. 0.01860311]

(モデル名).feature\_importances\_で特徴量の重要度を取得できる

STEP6:特徴量重要度の図示

#### コード14-14 特徴量重要度を見やすくする

```
forest importances = pd.DataFrame(
```

```
model_forest.feature_importances_,
```

```
index = bc.feature names,
```

```
columns=['Importance'])
```

```
STEP0:事前準備
STEP1:データの用意
STEP2:学習モデルの選択
STEP3:データを入れて学習
STEP4:予測を行う
STEP5:モデルの評価
STEP6:特徴量重要度の図示
```

- (モデル名).feature\_importances\_で特徴量の重要度を取得
- DataFrame型にして、index=の引数で行名をbc.feature\_namesで指定
- columns = で列名に名前をつける

STEP6:特徴量重要度の図示

コード14-15 特徴量重要度を見やすくする

forest\_importances

出力すると右のように特徴量とその重要
 度のデータフレームが作成されている

STEPO:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:予測を行う STEP5:モデルの評価 STEP6:特徴量重要度の図示

Importance

mean radius	0.023322
mean texture	0.009979
mean perimeter	0.066447
mean area	0.039999
mean smoothness	0.003435
mean compactness	0.005039
mean concavity	0.093138
mean concave points	0.136197
mean symmetry	0.001898
mean fractal dimension	0.002101
radius error	0.025128
texture error	0.002757



STEP6:特徴量重要度の図示

#### コード14-17 特徴量重要度の図示

forest\_importances.plot.bar()

pandasのDataFrame型では

<u>(データフレーム名).plot.bar()</u>で

棒グラフを描出

STEP0:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:予測を行う STEP5:モデルの評価 STEP6:特徴量重要度の図示



### まとめ

●演習8-11では教師あり機械学習の回帰(線形回帰、ロジスティック回帰) 演習12-14では教師あり機械学習のSVM, 決定木、ランダムフォレストを 行いました

	回帰	SVM	決定木	ランダムフォレスト
利点	・モデルがわか りやすい ・調整するパラ メータが少ない	・未知のデータに も予測性能が高い ・多くの特徴量を 一度に調べられる	・分類の仕方が わかりやすい ・調整するパラ メータが少ない	<ul> <li>・アンサンブル学習であり、たまたまできたモデルに依拠しない</li> <li>・予測性能が高い</li> <li>・特徴量重要度を求めることができる</li> </ul>
欠点	・データ自体の 分布にうまく当 てはめられない こともある	・モデルの中身が よくわからない ・特徴量自体の意 味がなくなる	・過学習になり やすい	・データ数が少ないとき に過学習になる

### まとめ

- ●同じ乳がんデータの分類に対して色々なアルゴリズムを実践してみたけれど、 結局どれが予測性能が高いかは、実際のデータとの相性次第(アンサンブル 学習は高くなりやすい)
- ●今後、医療データで機械学習の分類を行うときには、どのモデルを使って行うのか、特徴を理解するのが大事
- ●また、Pythonで実践することはそれほど難しくなく、コードをポチポチ押す だけで行うことができる。しかし、ハイパーパラメーターのチューニングに よっては全然違う結果になることもあり、内容を理解してチューニングを行 う必要がある

#### 次回からは、さらに複雑なモデルである深層学習(deep learning)を行う

### 演習14:課題

#### Webclassで課題を提出してください。締め切りは2024/02/08 23:59まで

<u>RandomForestClassifier() でハイパーパラメータの調整を行い、STEP2~STEP6</u> <u>まで実行してください</u>。

- 決定木の作成数を1000、木の最大深さ5、最大特徴量を6に設定してください。
- webclass

(1)検証データの正解率、(2)学習データの正解率、(3)AUC を回答してください。 特徴量重要度を求めてください。

(4) 最も重要な特徴量、(5) その値を回答してください。

### 演習14:課題

Webclassで課題を提出してください。締め切りは2024/02/08 23:59まで

● ランダムフォレストについて理解を問う選択問題が二問あります。正しい選択肢を一つwebclassで選択してください。

(6) 問題1:ランダムフォレストの主な利点は何ですか?
 1.計算コストが非常に低い。
 2.個々の決定木の性能が低くても、集合としては高い精度を達成できる。
 3.トレーニングデータの少ない状況でも、常に高い精度を保証する。
 4.ハイパーパラメータのチューニングが不要である。

(7) 問題2: ランダムフォレスト(分類)はどのように最終的な結果を予測を行いますか?

1.すべての決定木の予測を平均化する。 2.最も精度の高い決定木の予測を使用する。 3.すべての決定木の予測の中で最も多数を占めるクラスを最終的な予測とする。 4.それぞれの決定木が重み付けされ、重みが最も高い決定木の予測を使用する。

#### 授業準備:Webclassからコードをダウンロードし、 Google colaboratoryで開いておいてください

#### 演習授業中の質問対応について

Zoom ミーティング			- 0	×
<sup>*</sup> 演習授業中の質 応させていたた	〔問をチュータ ごきます。	一の先生が対	ミーティングチャット	C ×
演習にエラーが出 題があったらリア の <b>挙手</b> を押してくた	たなど問 クション 間 ごさい。	<sup>曹</sup> 質問内容を に送信して	入力して、 ください。	「 <b>全員</b> 」宛て
<ul> <li>Miho Ishimaru</li> <li></li></ul>	<ul> <li>● ▲ 音 ② ● 添 …</li> <li>● ● ○ ● 添 …</li> <li>● ● ● ○ ● ○ ●</li> <li>● ● ●</li> <li>● ● ●</li> <li>● ●</li> <li>● ● ●</li> <li>● ●<th>コーへ 目 へ … <b>桜</b>了 ポード ノート 詳細</th><th>バゼージは誰に表示されますか? 宛先: 全具 ▼ ここにメッセージを入力します ひ ② ① □, ▼ …</th><th>7</th></li></ul>	コーへ 目 へ … <b>桜</b> 了 ポード ノート 詳細	バゼージは誰に表示されますか? 宛先: 全具 ▼ ここにメッセージを入力します ひ ② ① □, ▼ …	7

1



\*本日演習16の授業後に複合領域コースの説明があります

#### ●今までは教師あり機械学習の基礎を実行してきた

#### ●演習15-20では深層学習を実行する

15-16で深層学習の基礎と乳がんデータの分類17-19で画像の分類を深層学習で行う20 機械学習・深層学習の演習



### 機械学習と深層学習の違い



### 深層学習のコードの流れ



ニューラルネットワークとは

- ・ニューロンは、樹状突起、細胞体、軸索からなる
- ・ニューロンは、樹状突起から入力された電気信号が神経細胞内の電位を超えるかどうかの
   している
- ・閾値を超えるとニューロンは興奮状態となり、軸索末端から電気信号が出力される



ニューラルネットワークとは

- ・ニューロンは、樹状突起、細胞体、軸索からなる
- ・ニューロンは、樹状突起から入力された電気信号が神経細胞内の電位を超えるかどうかの
- ・閾値を超えるとニューロンは興奮状態となり、軸索末端から電気信号が出力される



ニューラルネットワークとは

- ・ニューロンは、樹状突起、細胞体、軸索からなる
- ・ニューロンは、樹状突起から入力された電気信号が神経細胞内の電位を超えるかどうかの
- ・閾値を超えるとニューロンは興奮状態となり、軸索末端から電気信号が出力される



ニューラルネットワークとは

- ・ニューロンは、樹状突起、細胞体、軸索からなる
- ・ニューロンは、樹状突起から入力された電気信号が神経細胞内の電位を超えるかどうかの
- ・閾値を超えるとニューロンは興奮状態となり、軸索末端から電気信号が出力される



ニューラルネットワークとは

- ・ニューロンは、樹状突起、細胞体、軸索からなる
- ・ニューロンは、樹状突起から入力された電気信号が神経細胞内の電位を超えるかどうかの
- ・閾値を超えるとニューロンは興奮状態となり、軸索末端から電気信号が出力される



ニューラルネットワークとは

#### 単一の人工ニューロンはこのようなモデルで表すことができる。



前の層

今の層



ニューラルネットワークとは

#### 単一の人工ニューロンはこのようなモデルで表すことができる。

深層学習は、この重み(パラメーター)を最適化して出力(予測結果)を 正解に近づけるように学習する



![](_page_104_Figure_1.jpeg)

x:入力、w:重み、µ:入力合計、 f(x):活性化関数、y:(各ニューロンの)出力

![](_page_105_Figure_1.jpeg)

 $\mu_1 = \mathbf{w}_1 \times \mathbf{x}_1 + \mathbf{w}_3 \times \mathbf{x}_2 + \mathbf{w}_5 \times \mathbf{x}_3$ 

![](_page_106_Figure_1.jpeg)

![](_page_107_Figure_1.jpeg)

 $\mu_3 = w_7 \times y_1 + w_8 \times y_2$


 $y_3 = f_2(\mu_3)$ 

ニューラルネットワークとは



(a) ReLU関数

(b) sigmoid関数



ニューラルネットワークとは



ニューラルネットワークとは



ニューラルネットワークとは

#### 実際に計算してみる





実は、各ニューロンには、重み×入力の他に、バイアス(前のニューロンと 繋がっていない定数)も足される



$$\mu_{1} = w_{1} \times x_{1} + w_{3} \times x_{2} + w_{5} \times x_{3}$$

$$\mu_{2} = w_{2} \times x_{2} + w_{4} \times x_{2} + w_{6} \times x_{3}$$

$$\mu_{3} = w_{7} \times y_{1} + w_{8} \times y_{2}$$

$$\mu_{3} = w_{7} \times y_{1} + w_{8} \times y_{2}$$

$$\mu_{1} = w_{1} \times x_{1} + w_{3} \times x_{2} + w_{5} \times x_{3} + 1 \times b_{1}$$

$$\mu_{2} = w_{2} \times x_{2} + w_{4} \times x_{2} + w_{6} \times x_{3} + 1 \times b_{2}$$

$$\mu_{3} = w_{7} \times y_{1} + w_{8} \times y_{2} + 1 \times b_{3}$$

バイアス(前のニューロンと繋がっていない定数)もµに足す

ニューラルネットワークとは

### 実際に計算してみる



(a) ReLU関数



y = x (x > 0) $y = 0 (x \le 0)$ 







出力層の最後の活性化関数をsigmoid関数にすると、 (ロジスティック回帰分析と同様に) y(正解ラベル)が1になる確率が得られる

ニューラルネットワークとは

#### 特徴量3つ、正解が0or1(病気なら1)で考える

	体 重	年 齢	血 圧	正 解
1	55	28	140	1
2	48	44	130	0
3	77	64	145	0
4	49	42	130	1
5	59	66	128	1
6	74	38	109	0
7	59	46	137	1
8	52	41	140	1
9	38	56	110	0
10	47	53	121	1

					$\overline{}$
	/		//	$\frown$	$\backslash$
	体 重	年 齢	血圧	正 解	$  \rangle  $
1	55	28	140	1	
2	48	44	130	0	
3	77	64	145	0	
4	49	42	130	1	
5	59	66	128	1	
6	74	38	109	0	$  \setminus$
7	59	46	137	1	
8	52	41	140	1	
9	38	56	110	0	
10	47	53	121	1	

特徴量3つ、正解が0or1(病気なら1)で考える



	体 重	年 齢	血 圧	正 解	予 測
1	55	28	140	1	0.6
2	48	44	130	0	
3	77	64	145	0	
4	49	42	130	1	
5	59	66	128	1	
6	74	38	109	0	
7	59	46	137	1	
8	52	41	140	1	
9	38	56	110	0	
10	47	53	121	1	

ニューラルネットワークにはデータが1つずつ入力される 入力層のニューロンの数は1つのデータが持つ説明変数の数

					~
	/				
	体 重	年 齢	血 圧	正 解	$ \rangle$
1	55	28	140	1	
2	48	44	130	0	
3	77	64	145	0	
4	49	42	130	1	
5	59	66	128	1	
6	74	38	109	0	
7	59	46	137	1	
8	52	41	140	1	
9	38	56	110	0	
10	47	53	121	1	

#### 精度を上げるにはどうするか?



	体 重	年齢	自日	正 解	予测
1	55	28	140	1	0.6
2	48	44	130	0	0.3
3	77	64	145	0	0.4
4	49	42	130	1	0.9
5	59	66	128	1	0.6
6	74	38	109	0	0.2
7	59	46	137	1	0.4
8	52	41	140	1	0.9
9	38	56	110	0	0.6
10	47	53	121	1	0.5

ニューラルネットワークにはデータが1つずつ入力される 入力層のニューロンの数は1つのデータが持つ説明変数の数

			/	$\frown$	_
	/	/	/ /	$\frown$	
	体重	年齢	血圧	正 解	
1	55	28	140	1	
2	48	44	130	0	
3	77	64	145	0	
4	49	42	130	1	
5	59	66	128	1	
6	74	38	109	0	
7	59	46	137	1	
8	52	41	140	1	
9	38	56	110	0	
10	47	53	121	1	

#### 精度を上げるにはどうするか? 正解と予測の誤差が小さくなるように 重みとバイアスを更新する



	体 重	年 齢	血圧	正 解	予 測	
1	55	28	140	1	0.6	
2	48	44	130	0	0.3	
3	77	64	145	0	0.4	
4	49	42	130	1	0.9	
5	59	66	128	1	0.6	
6	74	38	109	0	0.2	
7	59	46	137	1	0.4	
8	52	41	140	1	0.9	
9	38	56	110	0	0.6	
10	47	53	121	1	0.5	

ニューラルネットワークにはデータが1つずつ入力される 入力層のニューロンの数は1つのデータが持つ説明変数の数

ニューラルネットワークとは



①入力したデータの予測結果を算出する
 ②正解と予測結果がどれくらい異なっているかという誤差を計算する

③誤差が小さくなるように重みとバイアスを変える

を何度も繰り返すことで、誤差を減らしていき精度を高める



深層学習(乳がんデータ)

#### 深層学習の概念を掴むため、前回まで使ってきた 乳がんデータを使って深層学習を行う





### 深層学習(乳がんデータの分類)コードまとめ





# STEP0: Google Colaboratoryの立ち上げ

STEPO:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:学習結果の図示 STEP5:モデルの評価

Python基礎 プログラミング基礎

#### 検索google colab <u>Colaboratory へようこそ - Colaboratory - Google</u>

	Colaboratory へようこそ ファイル 編集 表示 挿入 ランタイム	ツール ヘルプ
	ノートブックを新規作成	
	ノートブックを開く	ctt+0 澳首15-16コート.lpynp
4	ノートブックをアップロード	lab へようこそ
{x} =		に Colab をよくご存じの場合は、この動画でインタラクティブなラ
07	ドライブにコピーを保存	ドの履歴表示、コマンドパレットについてご覧ください。
	コピーを GitHub Gist として保存 GitHub にコピーを保存	3 Cool Google Colab Features
	保存 変更履歴	Ctrl+S
	ダウンロード	
	印刷	Ctrl+P
		Colab とは



#### 検索google colab <u>Colaboratory へようこそ - Colaboratory - Google</u>

ノートブックを開く





#### STEPO:ライブラリのインポート

STEPO:事前準備 STEP1:データの用意 STEP2:学習モデルの選択 STEP3:データを入れて学習 STEP4:学習結果の図示 STEP5:モデルの評価

## コード15-1 ライブラリとモジュールをインポート

- import numpy as np
- import pandas as pd
- import matplotlib.pyplot as plt

モジュールをインポート:import <u>(ライブラリ名)</u>.<u>(モジュール名)</u>as 省略形 matplotlib pyplot





bc.target

bc.data

array([[1.799e+01, 1.038e+01, 1.228e+02, …, 2.654e-01, 4.601e-01, 1.189e-01], [2.057e+01, 1.777e+01, 1.329e+02, …, 1.860e-01, 2.750e-01, 8.902e-02], [1.969e+01,

colab

as\_frame = Falseで指定したので、numpy配列になっている [



# 深層学習(乳がんデータの分類)

#### STEP1:データの用意

#### コード15-2 乳がんデータ

STEPO:事前準備
STEP1:データの用意
STEP2:学習モデルの選択
STEP3:データを入れて学習
STEP4:学習結果の図示
STEP5:モデルの評価





from sklearn.model\_selection import train\_test\_split

x\_train, x\_test, y\_train, y\_test = train\_test\_split(

bc.data, bc.target, test\_size = 0.3, random\_state = 0)

- trainデータ70%とtestデータ30%にデータを分割
- 分割時の乱数シード値を"0"に指定





- 簡単なニューラルネットワークを行うため、特徴量を1~3番目 (インデックス 番号0~2)の特徴量だけを選択
- (データ名) [行番号,列番号] でnp配列の時は抽出できる
- 今回は全ての行を意味する(: コロン)、列は0~3未満のインデックス番号を 指定する 0:3 で抽出を行う



#### STEP1:データの用意

コード15-5 データの配列構造を確認する

```
print(x_train.shape)
```

```
print(x_test.shape)
```

```
print(x_train3.shape)
```

print(x\_test3.shape)









#### Webclassに課題があります。 締め切りは2024/02/14 23:59まで