

開発者向けオンラインセミナー

機械学習101

堀田 稔

インターシステムズジャパン株式会社

2023年11月29日

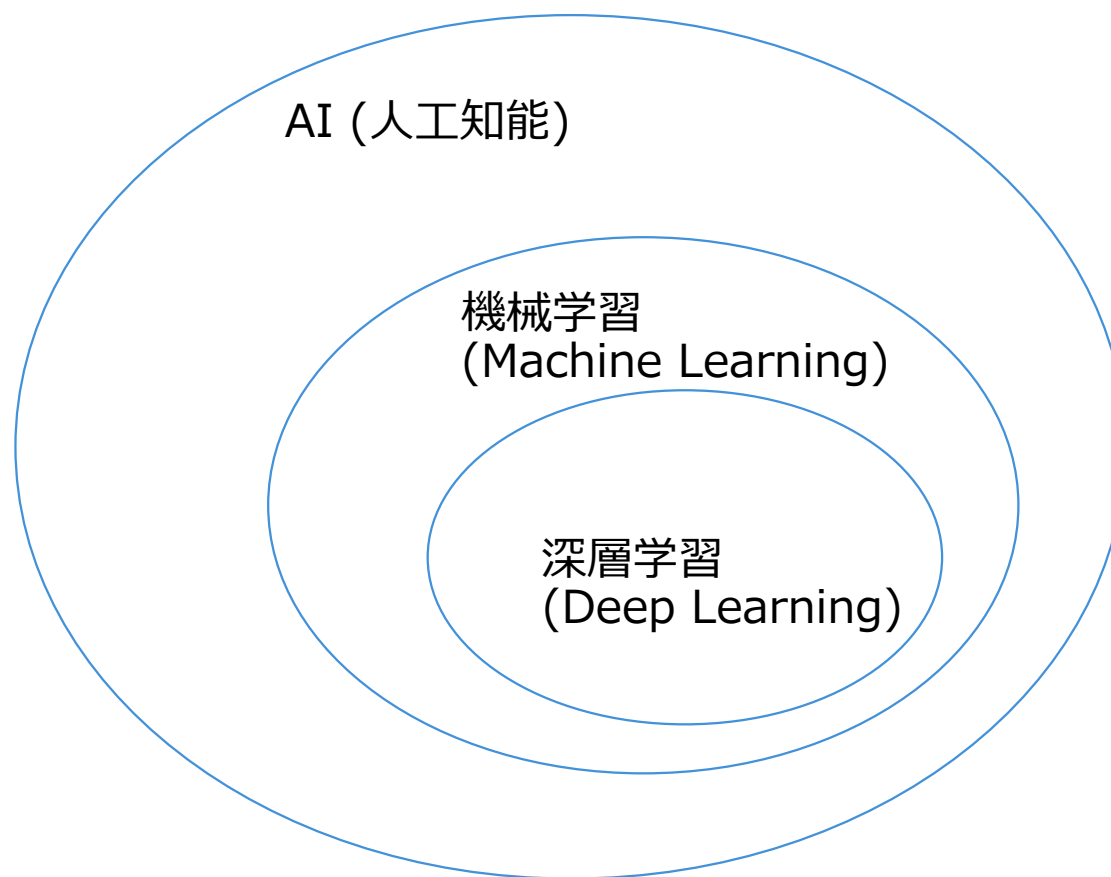


本セミナーの目標



- 機械学習の基本的な仕組みを理解する。
- 機械学習にまつわる用語を理解するための枠組みを紹介する。
- 学習のやり方や性能評価の基本を理解する。
- 機械学習のライブラリを使う利用者にとっても、何が行われているのかを理解するのに必要な知識を得る。
- ※ 説明の関係上数式が出てきますが、苦手な方は無理に理解する必要はありません。

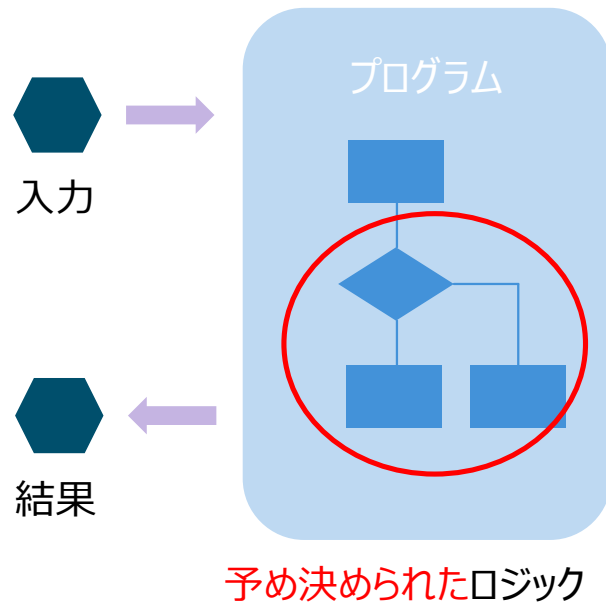
AI - 機械学習 - 深層学習 の関係



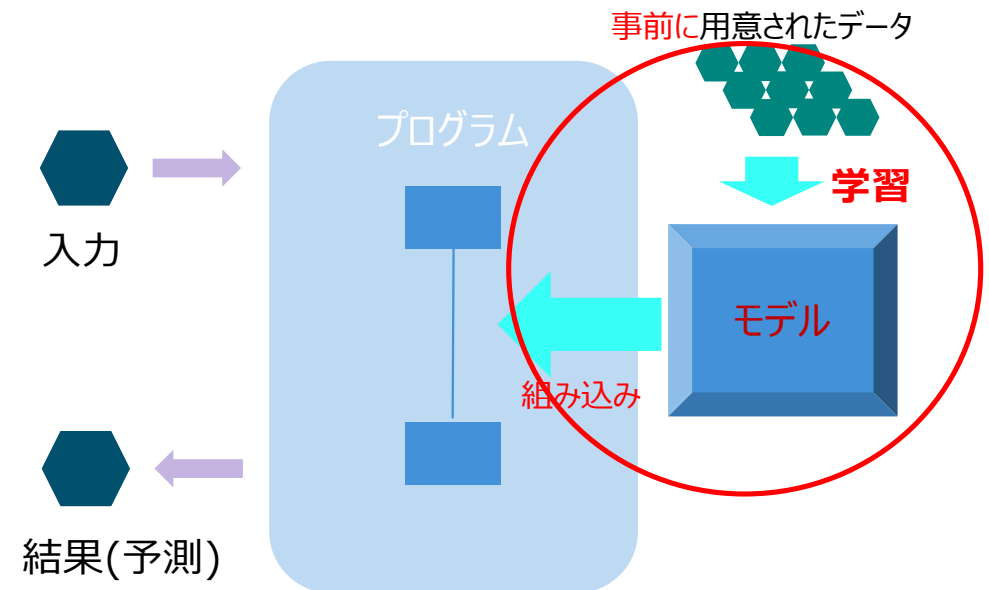
「通常」のプログラムと機械学習を用いたプログラムの違い



通常のプログラム



機械学習を用いたプログラム



機械学習の基本的な考え方

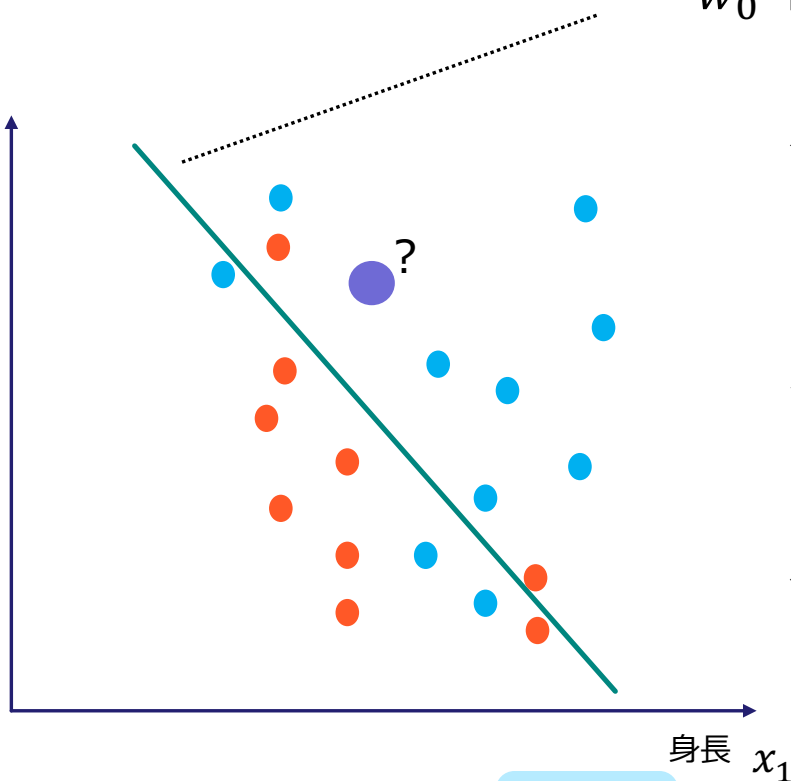


特徴量

x_2
預金額

- 大人
- 子供

分類



特徴量

$$w_0 + w_1 x_1 + w_2 x_2 = 0$$

w_0 w_1 w_2 が直線を決定する

既知のデータから、最も誤りの少ない w_0, w_1, w_2 を求める

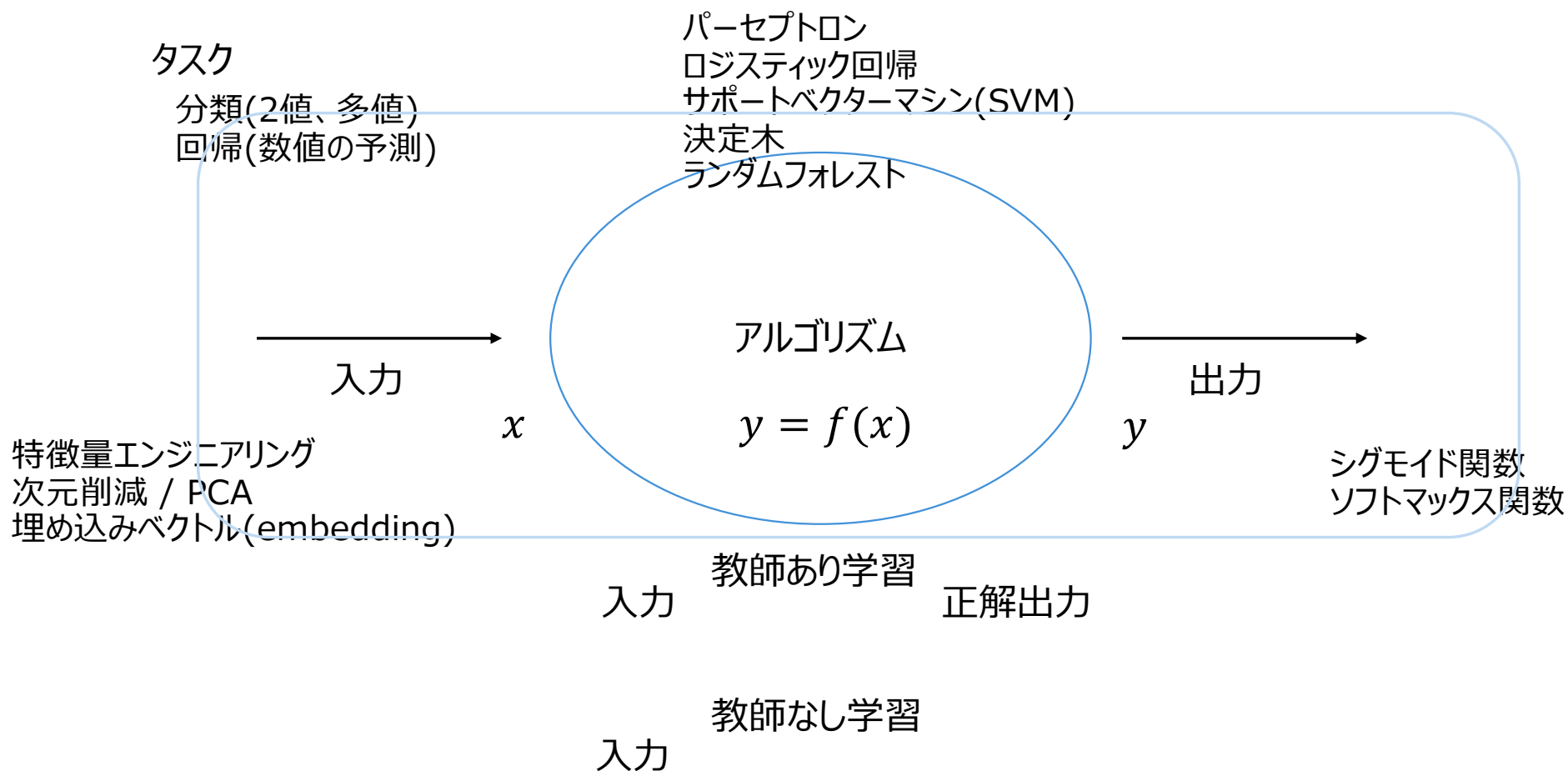
未知のデータ x_1, x_2 対して $w_0 + w_1 x_1 + w_2 x_2$ を計算し、0より大きい小さいかで大人か子供かを判断

モデル

学習

予測

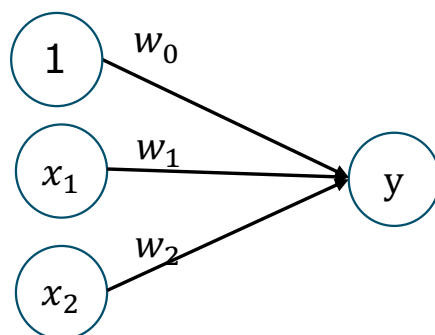
機械学習の流れと用語



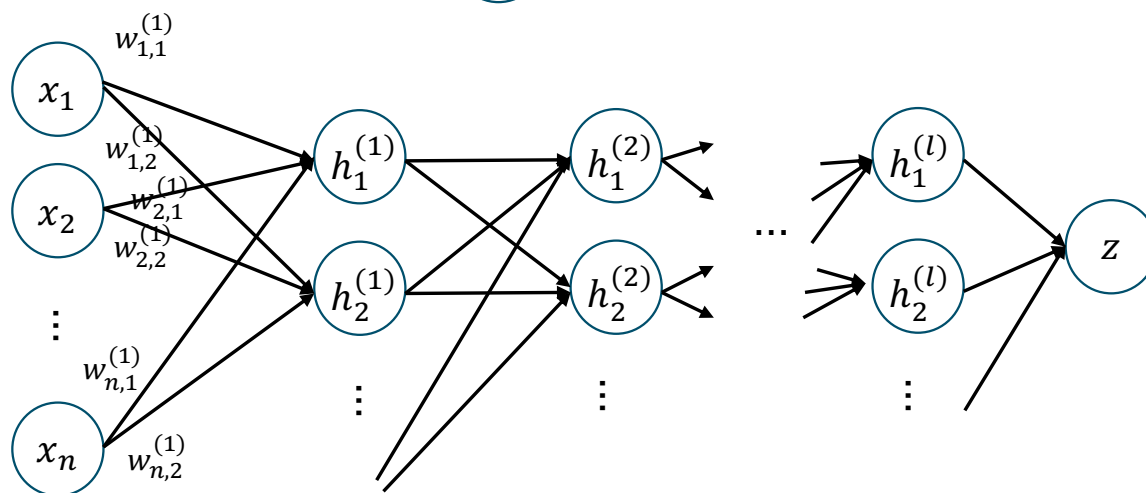
深層学習



単純パーセプトロン



深層学習



- 多様な「関数」が表現できる
- パラメータが多く、学習に大量の計算資源が必要 → GPUの活用
- 層の重ね方や繋ぎ方にバリエーション
 - CNN, RNN, ...

大規模言語モデルのキーワード： 埋め込みベクトル、Transformer

教師データの使い方



教師データ

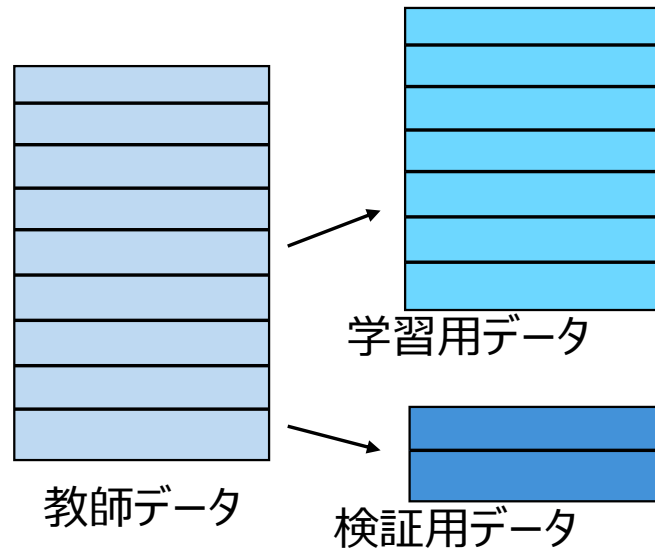
学習
性能評価



性能OK !



学習データに対して最適なモデルになり、未知のデータに対する性能(汎化性能)が得られない(過学習)ことになりやすい。



学習用データ

検証用データ

学習



教師データを学習用データと検証用データに分けて、検証用データで性能を評価。

場合によっては、さらにテスト用データを用意することもある。

性能評価

性能評価 (2値分類タスク)

代表的な指標

$$\text{適合率(Precision)} = \frac{TP}{TP + FP}$$

陽性と予測したデータの中の、真の陽性のデータの割合
高いほど予測の「ノイズ」が少ない

$$\text{再現率(Recall)} = \frac{TP}{TP + FN}$$

真の陽性のデータの中で、陽性と予測できたデータの割合
高いほど予測の「見落とし」が小さい

		予測値	
		Positive(陽性)	Negative(陰性)
真の値	Positive(陽性)	TP: True Positive (真陽性) 陽性を 正しく 陽性と予測	FN: False Negative (偽陰性) 陽性を 誤って 陰性と予測
	Negative(陰性)	FP: False Positive (偽陽性) 陰性を 誤って 陽性と予測	TN: True Negative (真陰性) 陰性を 正しく 陰性と予測

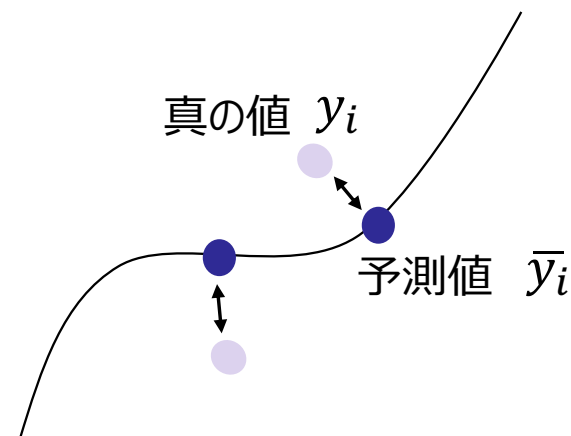
混同行列 (Confusion Matrix)



性能評価 (回帰タスク)

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \bar{y}_i)^2}$$

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \bar{y}_i|$$

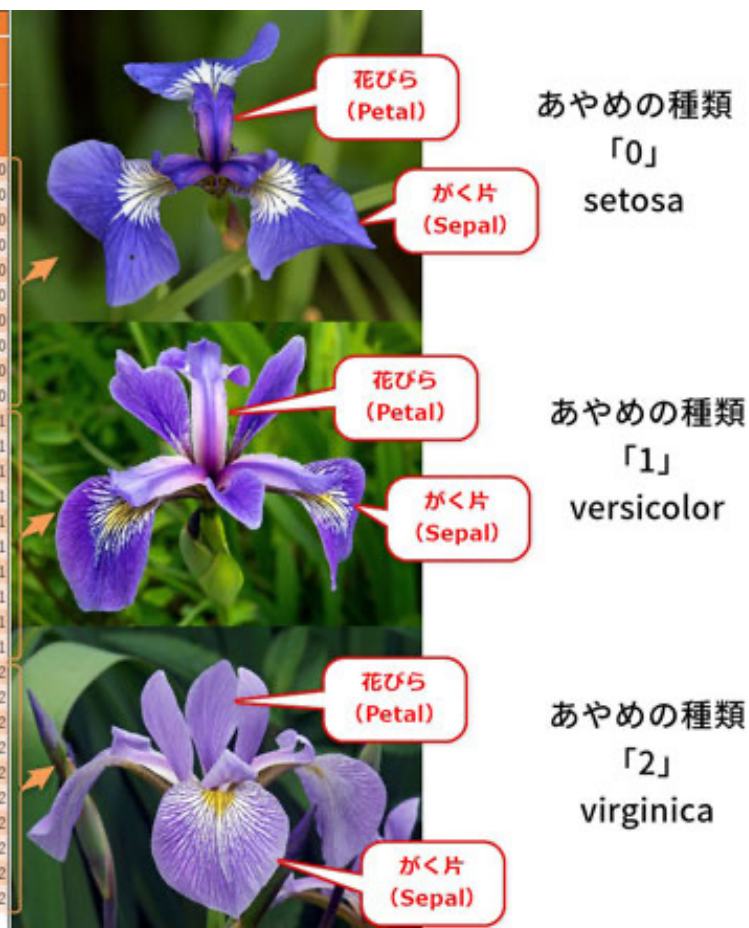


試してみよう：Irisの分類

Pythonで一般に使われる機械学習・データ分析ライブラリ scikit-learnを使って、Iris(アヤメ)の品種分類を試してみましょう

列番号→	1	2	3	4	目的変数→
説明変数→	Sepal Length	Sepal Width	Petal Length	Petal Width	target
行番号→	がく片の長さ (cm)	がく片の幅 (cm)	花びらの長さ (cm)	花びらの幅 (cm)	あやめの種類
0	5.1	3.5	1.4	0.2	0
1	4.9	3.0	1.4	0.2	0
2	4.7	3.2	1.3	0.2	0
3	4.6	3.1	1.5	0.2	0
4	5.0	3.6	1.4	0.2	0
5	5.4	3.9	1.7	0.4	0
6	4.6	3.4	1.4	0.3	0
7	5.0	3.4	1.5	0.2	0
8	4.4	2.9	1.4	0.2	0
9	5.0	3.3	1.4	0.2	0
10	7.0	3.2	4.7	1.4	1
11	6.4	3.2	4.5	1.5	1
12	6.9	3.1	4.9	1.5	1
13	5.5	2.3	4.0	1.3	1
14	6.5	2.8	4.6	1.5	1
15	5.7	2.8	4.5	1.3	1
16	6.3	3.3	4.7	1.6	1
17	4.9	2.4	3.3	1.0	1
18	6.6	2.9	4.6	1.3	1
19	5.7	2.8	4.1	1.3	1
20	6.3	3.3	6.0	2.5	2
21	5.8	2.7	5.1	1.9	2
22	7.1	3.0	5.9	2.1	2
23	6.3	2.9	5.6	1.8	2
24	6.5	3.0	5.8	2.2	2
25	7.6	3.0	6.6	2.1	2
26	4.9	2.5	4.5	1.7	2
27	7.3	2.9	6.3	1.8	2
28	6.7	2.5	5.8	1.8	2
29	5.9	3.0	5.1	1.8	2

※合計データ数は150個。表の右端1個は目的変数（ラベル）。



<https://atmarkit.itmedia.co.jp/ait/articles/2206/13/news032.html> から引用