

# 7.1. Toy datasets

scikit-learn comes with a few small standard datasets that do not require to download any file from some external website.

They can be loaded using the following functions:

<a href="#">load_iris</a> (*[, return_X_y, as_frame])	Load and return the iris dataset (classification).
<a href="#">load_diabetes</a> (*[, return_X_y, as_frame, scaled])	Load and return the diabetes dataset (regression).
<a href="#">load_digits</a> (*[, n_class, return_X_y, as_frame])	Load and return the digits dataset (classification).
<a href="#">load_linnerud</a> (*[, return_X_y, as_frame])	Load and return the physical exercise Linnerud dataset.
<a href="#">load_wine</a> (*[, return_X_y, as_frame])	Load and return the wine dataset (classification).
<a href="#">load_breast_cancer</a> (*[, return_X_y, as_frame])	Load and return the breast cancer wisconsin dataset (classification).

These datasets are useful to quickly illustrate the behavior of the various algorithms implemented in scikit-learn. They are however often too small to be representative of real world machine learning tasks.

## 7.1.1. Iris plants dataset

### Data Set Characteristics:

<b>Number of Instances:</b>
150 (50 in each of three classes)
<b>Number of Attributes:</b>
4 numeric, predictive attributes and the class
<b>Attribute Information:</b>
<ul style="list-style-type: none"><li>sepal length in cm</li><li>sepal width in cm</li><li>petal length in cm</li><li>petal width in cm</li><li><b>class:</b><ul style="list-style-type: none"><li>Iris-Setosa</li><li>Iris-Versicolour</li><li>Iris-Virginica</li></ul></li></ul>
<b>Summary Statistics:</b>
<div><div></div><div>sepal length: 4.3 7.9 5.84 0.83 0.7826</div><div>sepal width: 2.0 4.4 3.05 0.43 -0.4194</div><div>petal length: 1.0 6.9 3.76 1.76 0.9490 (high!)</div><div>petal width: 0.1 2.5 1.20 0.76 0.9565 (high!)</div><div></div></div>
<b>Missing Attribute Values:</b>
None
<b>Class Distribution:</b>
33.3% for each of 3 classes.
<b>Creator:</b>
R.A. Fisher
<b>Donor:</b>
Michael Marshall ( <a href="#">MARSHALL%PLU@io.arc.nasa.gov</a> )
<b>Date:</b>
July, 1988

Toggle Menu

The famous Iris database, first used by Sir R.A. Fisher. The dataset is taken from Fisher’s paper. Note that it’s the same as in R, but not as in the UCI Machine Learning Repository, which has two wrong data points.

This is perhaps the best known database to be found in the pattern recognition literature. Fisher’s paper is a classic in the field and is referenced frequently to this day. (See Duda & Hart, for example.) The data set contains 3 classes of 50 instances each, where each class refers to a type of iris plant. One class is linearly separable from the other 2; the latter are NOT linearly separable from each other.

References

- Fisher, R.A. “The use of multiple measurements in taxonomic problems” Annual Eugenics, 7, Part II, 179-188 (1936); also in “Contributions to Mathematical Statistics” (John Wiley, NY, 1950).
- Duda, R.O., & Hart, P.E. (1973) Pattern Classification and Scene Analysis. (Q327.D83) John Wiley & Sons. ISBN 0-471-22361-1. See page 218.
- Dasarathy, B.V. (1980) “Nosing Around the Neighborhood: A New System Structure and Classification Rule for Recognition in Partially Exposed Environments”. IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. PAMI-2, No. 1, 67-71.
- Gates, G.W. (1972) “The Reduced Nearest Neighbor Rule”. IEEE Transactions on Information Theory, May 1972, 431-433.
- See also: 1988 MLC Proceedings, 54-64. Cheeseman et al’s AUTOCLASS II conceptual clustering system finds 3 classes in the data.
- Many, many more ...

7.1.2. Diabetes dataset

Ten baseline variables, age, sex, body mass index, average blood pressure, and six blood serum measurements were obtained for each of n = 442 diabetes patients, as well as the response of interest, a quantitative measure of disease progression one year after baseline.

Data Set Characteristics:

<b>Number of Instances:</b>
442
<b>Number of Attributes:</b>
First 10 columns are numeric predictive values
<b>Target:</b>
Column 11 is a quantitative measure of disease progression one year after baseline
<b>Attribute Information:</b>
<ul style="list-style-type: none"><li>• age age in years</li><li>• sex</li><li>• bmi body mass index</li><li>• bp average blood pressure</li><li>• s1 tc, total serum cholesterol</li><li>• s2 ldl, low-density lipoproteins</li><li>• s3 hdl, high-density lipoproteins</li><li>• s4 tch, total cholesterol / HDL</li><li>• s5 ltg, possibly log of serum triglycerides level</li><li>• s6 glu, blood sugar level</li></ul>

Note: Each of these 10 feature variables have been mean centered and scaled by the standard deviation times the square root of n\_samples (i.e. the sum of squares of each column totals 1).

Source URL: <https://www4.stat.ncsu.edu/~boos/var.select/diabetes.html>

For more information see: Bradley Efron, Trevor Hastie, Iain Johnstone and Robert Tibshirani (2004) “Least Angle Regression,” Annals of Statistics (with discussion), 407-499. ([https://web.stanford.edu/~hastie/Papers/LARS/LeastAngle\\_2002.pdf](https://web.stanford.edu/~hastie/Papers/LARS/LeastAngle_2002.pdf))

7.1.3. Optical recognition of handwritten digits dataset

Data Set Characteristics:

<b>Number of Instances:</b>
1797
<b>Number of Attributes:</b>
64
<b>Attribute Information:</b>
8x8 image of integer pixels in the range 0..16.
<b>Missing Attribute Values:</b>
None
<b>Creator:</b>
5. Alpaydin (alpaydin '@' boun.edu.tr)
<b>Date:</b>
July; 1998

This is a copy of the test set of the UCI ML hand-written digits datasets  
<https://archive.ics.uci.edu/ml/datasets/Optical+Recognition+of+Handwritten+Digits>

The data set contains images of hand-written digits: 10 classes where each class refers to a digit.

Preprocessing programs made available by NIST were used to extract normalized bitmaps of handwritten digits from a preprinted form. From a total of 43 people, 30 contributed to the training set and different 13 to the test set. 32x32 bitmaps are divided into nonoverlapping blocks of 4x4 and the number of on pixels are counted in each block. This generates an input matrix of 8x8 where each element is an integer in the range 0..16. This reduces dimensionality and gives invariance to small distortions.

For info on NIST preprocessing routines, see M. D. Garris, J. L. Blue, G. T. Candela, D. L. Dimmick, J. Geist, P. J. Grother, S. A. Janet, and C. L. Wilson, NIST Form-Based Handprint Recognition System, NISTIR 5469, 1994.

References

- C. Kaynak (1995) Methods of Combining Multiple Classifiers and Their Applications to Handwritten Digit Recognition, MSc Thesis, Institute of Graduate Studies in Science and Engineering, Bogazici University.
- 5. Alpaydin, C. Kaynak (1998) Cascading Classifiers, Kybernetika.
- Ken Tang and Ponnuthurai N. Suganthan and Xi Yao and A. Kai Qin. Linear dimensionalityreduction using relevance weighted LDA. School of Electrical and Electronic Engineering Nanyang Technological University. 2005.
- Claudio Gentile. A New Approximate Maximal Margin Classification Algorithm. NIPS. 2000.

7.1.4. Linnerrud dataset

Data Set Characteristics:

<b>Number of Instances:</b>
20
<b>Number of Attributes:</b>
3
<b>Missing Attribute Values:</b>
None

The Linnerud dataset is a multi-output regression dataset. It consists of three exercise (data) and three physiological (target) variables collected from twenty middle-aged men in a fitness club:

- **physiological - CSV containing 20 observations on 3 physiological variables:**  
Weight, Waist and Pulse.
- **exercise - CSV containing 20 observations on 3 exercise variables:**  
Chins, Situps and Jumps.

- Tenenhaus, M. (1998). La regression PLS: theorie et pratique. Paris: Editions Technic.

## 7.1.5. Wine recognition dataset

**Data Set Characteristics:**

Number of Instances:																																																																					
178																																																																					
Number of Attributes:																																																																					
13 numeric, predictive attributes and the class																																																																					
Attribute Information:																																																																					
<ul style="list-style-type: none"><li>Alcohol</li><li>Malic acid</li><li>Ash</li><li>Alcalinity of ash</li><li>Magnesium</li><li>Total phenols</li><li>Flavanoids</li><li>Nonflavanoid phenols</li><li>Proanthocyanins</li><li>Color intensity</li><li>Hue</li><li>OD280/OD315 of diluted wines</li><li>Proline</li></ul>																																																																					
<div></div> <ul style="list-style-type: none"><li>class:<ul style="list-style-type: none"><li>class_0</li><li>class_1</li><li>class_2</li></ul></li></ul>																																																																					
Summary Statistics:																																																																					
<div></div> <table><tr><td>Alcohol:</td><td>11.0</td><td>14.8</td><td>13.0</td><td>0.8</td></tr><tr><td>Malic Acid:</td><td>0.74</td><td>5.80</td><td>2.34</td><td>1.12</td></tr><tr><td>Ash:</td><td>1.36</td><td>3.23</td><td>2.36</td><td>0.27</td></tr><tr><td>Alcalinity of Ash:</td><td>10.6</td><td>30.0</td><td>19.5</td><td>3.3</td></tr><tr><td>Magnesium:</td><td>70.0</td><td>162.0</td><td>99.7</td><td>14.3</td></tr><tr><td>Total Phenols:</td><td>0.98</td><td>3.88</td><td>2.29</td><td>0.63</td></tr><tr><td>Flavanoids:</td><td>0.34</td><td>5.08</td><td>2.03</td><td>1.00</td></tr><tr><td>Nonflavanoid Phenols:</td><td>0.13</td><td>0.66</td><td>0.36</td><td>0.12</td></tr><tr><td>Proanthocyanins:</td><td>0.41</td><td>3.58</td><td>1.59</td><td>0.57</td></tr><tr><td>Colour Intensity:</td><td>1.3</td><td>13.0</td><td>5.1</td><td>2.3</td></tr><tr><td>Hue:</td><td>0.48</td><td>1.71</td><td>0.96</td><td>0.23</td></tr><tr><td>OD280/OD315 of diluted wines:</td><td>1.27</td><td>4.00</td><td>2.61</td><td>0.71</td></tr><tr><td>Proline:</td><td>278</td><td>1680</td><td>746</td><td>315</td></tr></table> <div></div>					Alcohol:	11.0	14.8	13.0	0.8	Malic Acid:	0.74	5.80	2.34	1.12	Ash:	1.36	3.23	2.36	0.27	Alcalinity of Ash:	10.6	30.0	19.5	3.3	Magnesium:	70.0	162.0	99.7	14.3	Total Phenols:	0.98	3.88	2.29	0.63	Flavanoids:	0.34	5.08	2.03	1.00	Nonflavanoid Phenols:	0.13	0.66	0.36	0.12	Proanthocyanins:	0.41	3.58	1.59	0.57	Colour Intensity:	1.3	13.0	5.1	2.3	Hue:	0.48	1.71	0.96	0.23	OD280/OD315 of diluted wines:	1.27	4.00	2.61	0.71	Proline:	278	1680	746	315
Alcohol:	11.0	14.8	13.0	0.8																																																																	
Malic Acid:	0.74	5.80	2.34	1.12																																																																	
Ash:	1.36	3.23	2.36	0.27																																																																	
Alcalinity of Ash:	10.6	30.0	19.5	3.3																																																																	
Magnesium:	70.0	162.0	99.7	14.3																																																																	
Total Phenols:	0.98	3.88	2.29	0.63																																																																	
Flavanoids:	0.34	5.08	2.03	1.00																																																																	
Nonflavanoid Phenols:	0.13	0.66	0.36	0.12																																																																	
Proanthocyanins:	0.41	3.58	1.59	0.57																																																																	
Colour Intensity:	1.3	13.0	5.1	2.3																																																																	
Hue:	0.48	1.71	0.96	0.23																																																																	
OD280/OD315 of diluted wines:	1.27	4.00	2.61	0.71																																																																	
Proline:	278	1680	746	315																																																																	
Missing Attribute Values:																																																																					
None																																																																					
Class Distribution:																																																																					
class_0 (59), class_1 (71), class_2 (48)																																																																					
Creator:																																																																					
R.A. Fisher																																																																					
Donor:																																																																					
Michael Marshall ( <a href="mailto:MARSHALL%PLU@io.arc.nasa.gov">MARSHALL%PLU@io.arc.nasa.gov</a> )																																																																					
Date:																																																																					
July, 1988																																																																					

This is a copy of UCI ML Wine recognition datasets. <https://archive.ics.uci.edu/ml/machine-learning-databases/wine/wine.data>

The data is the results of a chemical analysis of wines grown in the same region in Italy by three different cultivators. There are thirteen different measurements taken for different constituents found in the three types of wine.

Original Owners:

Forina, M. et al, PARVUS - An Extendible Package for Data Exploration, Classification and Correlation. Institute of Pharmaceutical and Food Analysis and Technologies, Via Brigata Salerno, 16147 Genoa, Italy.

Citation:

Lichman, M. (2013). UCI Machine Learning Repository [<https://archive.ics.uci.edu/ml>]. Irvine, CA: University of California, School of Information and Computer Science.

## References

(1) S. Aeberhard, D. Coomans and O. de Vel, Comparison of Classifiers in High Dimensional Settings, Tech. Rep. no. 92-02, (1992), Dept. of Computer Science and Dept. of Mathematics and Statistics, James Cook University of North Queensland. (Also submitted to Technometrics).

The data was used with many others for comparing various classifiers. The classes are separable, though only RDA has achieved 100% correct classification. (RDA : 100%, QDA 99.4%, LDA 98.9%, 1NN 96.1% (z-transformed data)) (All results using the leave-one-out technique)

(2) S. Aeberhard, D. Coomans and O. de Vel, "THE CLASSIFICATION PERFORMANCE OF RDA" Tech. Rep. no. 92-01, (1992), Dept. of Computer Science and Dept. of Mathematics and Statistics, James Cook University of North Queensland. (Also submitted to Journal of Chemometrics).

## 7.1.6. Breast cancer wisconsin (diagnostic) dataset

**Data Set Characteristics:**

<b>Number of Instances:</b>
569
<b>Number of Attributes:</b>
30 numeric, predictive attributes and the class
<b>Attribute Information:</b>
<ul style="list-style-type: none"><li>radius (mean of distances from center to points on the perimeter)</li><li>texture (standard deviation of gray-scale values)</li><li>perimeter</li><li>area</li><li>smoothness (local variation in radius lengths)</li><li>compactness (perimeter<sup>2</sup> / area - 1.0)</li><li>concavity (severity of concave portions of the contour)</li><li>concave points (number of concave portions of the contour)</li><li>symmetry</li><li>fractal dimension ("coastline approximation" - 1)</li></ul> <p>The mean, standard error, and "worst" or largest (mean of the three worst/largest values) of these features were computed for each image, resulting in 30 features. For instance, field 0 is Mean Radius, field 10 is Radius SE, field 20 is Worst Radius.</p> <ul style="list-style-type: none"><li><b>class:</b><ul style="list-style-type: none"><li>WDBC-Malignant</li><li>WDBC-Benign</li></ul></li></ul>

**Summary Statistics:**

◀	▶
---	---

radius (mean):	6.981	28.11
texture (mean):	9.71	39.28
perimeter (mean):	43.79	188.5
area (mean):	143.5	2501.0
smoothness (mean):	0.053	0.163
compactness (mean):	0.019	0.345
concavity (mean):	0.0	0.427
concave points (mean):	0.0	0.201
symmetry (mean):	0.106	0.304
fractal dimension (mean):	0.05	0.097
radius (standard error):	0.112	2.873
texture (standard error):	0.36	4.885

perimeter (standard error):	0.757	21.98
area (standard error):	6.802	542.2
smoothness (standard error):	0.002	0.031
compactness (standard error):	0.002	0.135
concavity (standard error):	0.0	0.396
concave points (standard error):	0.0	0.053
symmetry (standard error):	0.008	0.079
fractal dimension (standard error):	0.001	0.03
radius (worst):	7.93	36.04
texture (worst):	12.02	49.54
perimeter (worst):	50.41	251.2
area (worst):	185.2	4254.0
smoothness (worst):	0.071	0.223
compactness (worst):	0.027	1.058
concavity (worst):	0.0	1.252
concave points (worst):	0.0	0.291
symmetry (worst):	0.156	0.664
fractal dimension (worst):	0.055	0.208

Missing Attribute Values:

None

Class Distribution:

212 - Malignant, 357 - Benign

Creator:

Dr. William H. Wolberg, W. Nick Street, Olvi L. Mangasarian

Donor:

Nick Street

Date:

November, 1995

This is a copy of UCI ML Breast Cancer Wisconsin (Diagnostic) datasets. <https://goo.gl/U2Uwz2>

Features are computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe characteristics of the cell nuclei present in the image.

Separating plane described above was obtained using Multisurface Method-Tree (MSM-T) [K. P. Bennett, "Decision Tree Construction Via Linear Programming." Proceedings of the 4th Midwest Artificial Intelligence and Cognitive Science Society, pp. 97-101, 1992], a classification method which uses linear programming to construct a decision tree. Relevant features were selected using an exhaustive search in the space of 1-4 features and 1-3 separating planes.

The actual linear program used to obtain the separating plane in the 3-dimensional space is that described in: [K. P. Bennett and O. L. Mangasarian: "Robust Linear Programming Discrimination of Two Linearly Inseparable Sets", Optimization Methods and Software 1, 1992, 23-34].

This database is also available through the UW CS ftp server:

ftp ftp.cs.wisc.edu cd math-prog/cpo-dataset/machine-learn/WDBC/

References

- W.N. Street, W.H. Wolberg and O.L. Mangasarian. Nuclear feature extraction for breast tumor diagnosis. IS&T/SPIE 1993 International Symposium on Electronic Imaging: Science and Technology, volume 1905, pages 861-870, San Jose, CA, 1993.
- O.L. Mangasarian, W.N. Street and W.H. Wolberg. Breast cancer diagnosis and prognosis via linear programming. Operations Research, 43(4), pages 570-577, July-August 1995.
- W.H. Wolberg, W.N. Street, and O.L. Mangasarian. Machine learning techniques to diagnose breast cancer from fine-needle aspirates. Cancer Letters 77 (1994) 163-171.