

1 The iris data set

This is a classic 'toy' data set used for machine learning testing is the iris data set.

The data set consists of 50 samples from each of three species of Iris (Iris setosa, Iris virginica and Iris versicolor). Four features were measured from each sample: the length and the width of the sepals and petals, in centimetres.

1.1 Loadint the iris data set

The iris data set comes preloaded in scikit learn. Let's load it and have a look at it.

```
import numpy as np
from sklearn import datasets
```

```
iris=datasets.load_iris()
```

```
# The iris dataset is an object that contains a number of elements:
```

```
print (list(iris))
```

OUT:

```
['data', 'target', 'target_names', 'DESCR', 'feature_names']
```

```
# feature_names shows data field titles:
```

```
print (iris.feature_names)
```

OUT:

```
['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)']
```

```
# data is the data for each sample; columns described by feature_names:
```

```
# lets' print just the first 10 roes
```

```
print (iris.data[0:10])
```

OUT:

```
# data is the data for each sample; columns described by feature_names:
```

```
# lets' print just the first 10 roes
```

```
print (iris.data[0:10])
```

```
[[5.1 3.5 1.4 0.2]
 [4.9 3.  1.4 0.2]
 [4.7 3.2 1.3 0.2]
 [4.6 3.1 1.5 0.2]
 [5.  3.6 1.4 0.2]
 [5.4 3.9 1.7 0.4]
 [4.6 3.4 1.4 0.3]
 [5.  3.4 1.5 0.2]
 [4.4 2.9 1.4 0.2]
 [4.9 3.1 1.5 0.1]]
```

```

# target_names lists types of iris identified:

print (iris.target_names)

OUT:

['setosa' 'versicolor' 'virginica']

# target lists the type of iris in each row of data:
# this maps to the target_names

print (iris.target)

OUT:

[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2
 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
 2 2]

# DESCR gives us description of the data set:

print (iris.DESCR)

OUT:

Iris Plants Database
=====

Notes
-----
Data Set Characteristics:
: Number of Instances: 150 (50 in each of three classes)
: Number of Attributes: 4 numeric, predictive attributes and the class
: Attribute Information:
  - sepal length in cm
  - sepal width in cm
  - petal length in cm
  - petal width in cm
  - class:
    - Iris-Setosa
    - Iris-Versicolour
    - Iris-Virginica
: Summary Statistics:

=====  ====  ====  =====  =====  =====
                Min  Max   Mean    SD    Class Correlation
=====  ====  ====  =====  =====  =====
sepal length:   4.3  7.9   5.84   0.83    0.7826
sepal width:    2.0  4.4   3.05   0.43   -0.4194
petal length:    1.0  6.9   3.76   1.76    0.9490 (high!)
petal width:     0.1  2.5   1.20   0.76    0.9565 (high!)
=====  ====  ====  =====  =====  =====

: Missing Attribute Values: None
: Class Distribution: 33.3% for each of 3 classes.
: Creator: R.A. Fisher

```

:Donor: Michael Marshall (MARSHALL%PLU@io.arc.nasa.gov)
:Date: July, 1988

This is a copy of UCI ML iris datasets.
<http://archive.ics.uci.edu/ml/datasets/Iris>

The famous Iris database, first used by Sir R.A Fisher

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 ...

1.2 Data sets in scikit learn

load_boston: boston house-prices dataset (regression).

load_iris: iris dataset (classification).

load_diabetes: diabetes dataset (regression).

load_digits: digits dataset (classification).

load_linnerud: linnerud dataset (multivariate regression).

load_wine: wine dataset (classification).

load_breast_cancer: breast cancer wisconsin dataset (classification).

1.3 Other sources of test data sets

<https://archive.ics.uci.edu/ml/datasets.html>

<https://blog.bigml.com/list-of-public-data-sources-fit-for-machine-learning/>