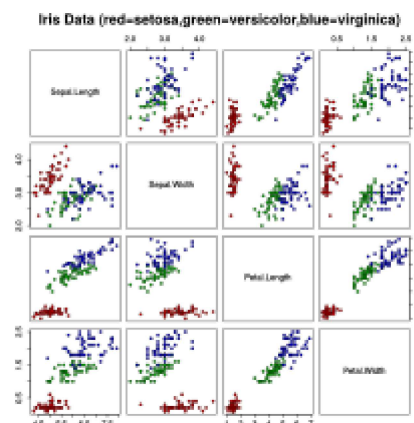


Iris flower data set

The **Iris flower data set** or **Fisher's Iris data set** is a multivariate data set introduced by the British statistician, eugenicist, and biologist Ronald Fisher in his 1936 paper *The use of multiple measurements in taxonomic problems* as an example of linear discriminant analysis.^[1] It is sometimes called **Anderson's Iris data set** because Edgar Anderson collected the data to quantify the morphologic variation of *Iris* flowers of three related species.^[2] Two of the three species were collected in the Gaspé Peninsula "all from the same pasture, and picked on the same day and measured at the same time by the same person with the same apparatus".^[3]

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 centimeters. Based on the combination of these four features, Fisher developed a linear discriminant model to distinguish the species from each other.



Scatterplot of the data set

Contents

Use of the data set

Data set

[R code illustrating usage](#)

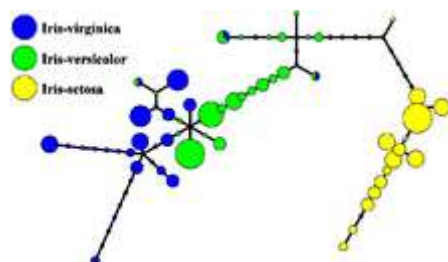
[Python code illustrating usage](#)

See also

References

External links

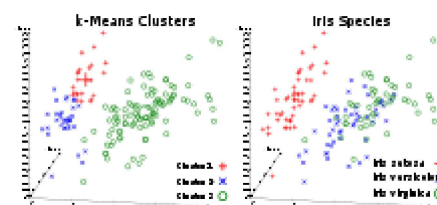
Use of the data set



An example of the so-called "metro map" for the *Iris* data set^[4] Only a small fraction of *Iris-virginica* is mixed with *Iris-versicolor*. All other samples of the different *Iris* species belong to the different nodes.

Based on Fisher's linear discriminant model, this data set became a typical test case for many statistical classification techniques in machine learning such as support vector machines.^[5]

The use of this data set in cluster analysis however is not common, since the data set only contains two clusters with rather obvious separation. One of the clusters contains *Iris setosa*, while the



Unsatisfactory k-means clustering (the data cannot be clustered into the known classes) and actual species visualized using ELKI

other cluster contains both *Iris virginica* and *Iris versicolor* and is not separable without the species information Fisher used. This makes the data set a good example to explain the difference between supervised and unsupervised techniques in data mining: Fisher's linear discriminant model can only be obtained when the object species are known: class labels and clusters are not necessarily the same.^[6]

Nevertheless, all three species of *Iris* are separable in the projection on the nonlinear and branching principal component.^[7] The data set is approximated by the closest tree with some penalty for the excessive number of nodes, bending and stretching. Then the so-called "metro map" is constructed.^[4] The data points are projected into the closest node. For each node the pie diagram of the projected points is prepared. The area of the pie is proportional to the number of the projected points. It is clear from the diagram (left) that the absolute majority of the samples of the different *Iris* species belong to the different nodes. Only a small fraction of *Iris-virginica* is mixed with *Iris-versicolor* (the mixed blue-green nodes in the diagram). Therefore, the three species of *Iris* (*Iris setosa*, *Iris virginica* and *Iris versicolor*) are separable by the unsupervising procedures of nonlinear principal component analysis. To discriminate them, it is sufficient just to select the corresponding nodes on the principal tree.

Data set

The dataset contains a set of 150 records under five attributes - sepal length, sepal width, petal length, petal width and species.

Fisher's *Iris* data

The iris data set is widely used as a beginner's dataset for machine learning purposes. The dataset is included in R base and Python in the machine learning package Scikit-learn, so that users can access it without having to find a source for it.

R code illustrating usage

```
iris
class(iris)
# "data.frame"

iris3
class(iris3)
#"array"
```



Iris setosa

Python code illustrating usage

```
from sklearn.datasets import load_iris

iris = load_iris()
iris
```

This code gives:

```
{'data': array([[5.1, 3.5, 1.4, 0.2],
                [4.9, 3. , 1.4, 0.2],
```



Iris versicolor

```
[4.7, 3.2, 1.3, 0.2],
[4.6, 3.1, 1.5, 0.2],...
'target': array([0, 0, 0, ... 1, 1, 1, ... 2, 2, 2, ...
'target_names': array(['setosa', 'versicolor', 'virginica'], dtype='<U10'),
...}
```

Several versions of the dataset have been published.^[8]

See also

- Classic data sets

References

- R. A. Fisher (1936). "The use of multiple measurements in taxonomic problems". *Annals of Eugenics*. **7** (2): 179–188. doi:10.1111/j.1469-1809.1936.tb02137.x. hdl:2440/15227.
- Edgar Anderson (1936). "The species problem in *Iris*". *Annals of the Missouri Botanical Garden*. **23** (3): 457–509. doi:10.2307/2394164. JSTOR 2394164.
- Edgar Anderson (1935). "The irises of the Gaspé Peninsula". *Bulletin of the American Iris Society*. **59**: 2–5.
- A. N. Gorban, A. Zinovyev. Principal manifolds and graphs in practice: from molecular biology to dynamical systems, International Journal of Neural Systems, Vol. 20, No. 3 (2010) 219–232.
- "UCI Machine Learning Repository: Iris Data Set". *archive.ics.uci.edu*. Retrieved 2017-12-01.
- Ines Färber, Stephan Günnemann, Hans-Peter Kriegel, Peer Kröger, Emmanuel Müller, Erich Schubert, Thomas Seidl, Arthur Zimek (2010). "On Using Class-Labels in Evaluation of Clusterings" (PDF). In Xiaoli Z. Fern; Ian Davidson; Jennifer Dy (eds.). *MultiClust: Discovering, Summarizing, and Using Multiple Clusterings*. **ACM SIGKDD**.
- A.N. Gorban, N.R. Sumner, and A.Y. Zinovyev, Topological grammars for data approximation, Applied Mathematics Letters Volume 20, Issue 4 (2007), 382-386.
- Bezdek, J.C. and Keller, J.M. and Krishnapuram, R. and Kuncheva, L.I. and Pal, N.R. (1999). "Will the real iris data please stand up?". *IEEE Transactions on Fuzzy Systems*. **7** (3): 368–369. doi:10.1109/91.771092.

External links

- "Fisher's Iris Data". *(Contains two errors which are documented)*. UCI Machine Learning Repository: Iris Data Set.

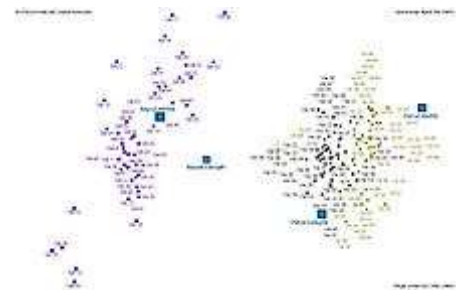
Retrieved from "https://en.wikipedia.org/w/index.php?title=Iris_flower_data_set&oldid=1008814603"

This page was last edited on 25 February 2021, at 05:08 (UTC).

Text is available under the Creative Commons Attribution-ShareAlike License; additional terms may apply. By using this site, you agree to the Terms of Use and Privacy Policy. Wikipedia® is a registered trademark of the Wikimedia Foundation, Inc., a non-profit organization.



Iris virginica



Spectramap biplot of Fisher's iris data set