

Creating your first model

Using SKLearn's datasets

Terry McCann | @SQLShark



sklearn.datasets : Datasets

The `sklearn.datasets` module includes utilities to load datasets, including methods to load and fetch popular reference datasets. It also features some artificial data generators.

User guide: See the [Dataset loading utilities](#) section for further details.

Loaders

<code>datasets.clear_data_home</code> ([data_home])	Delete all the content of the data home cache.
<code>datasets.dump_svmlight_file</code> (X, y, f[, ...])	Dump the dataset in svmlight / libsvm file format.
<code>datasets.fetch_20newsgroups</code> ([data_home, ...])	Load the filenames and data from the 20 newsgroups dataset (classification).
<code>datasets.fetch_20newsgroups_vectorized</code> ([...])	Load the 20 newsgroups dataset and vectorize it into token counts (classification).
<code>datasets.fetch_california_housing</code> ([...])	Load the California housing dataset (regression).
<code>datasets.fetch_covtype</code> ([data_home, ...])	Load the covtype dataset (classification).
<code>datasets.fetch_kddcup99</code> ([subset, data_home, ...])	Load the kddcup99 dataset (classification).
<code>datasets.fetch_lfw_pairs</code> ([subset, ...])	Load the Labeled Faces in the Wild (LFW) pairs dataset (classification).
<code>datasets.fetch_lfw_people</code> ([data_home, ...])	Load the Labeled Faces in the Wild (LFW) people dataset (classification).
<code>datasets.fetch_olivetti_faces</code> ([data_home, ...])	Load the Olivetti faces data-set from AT&T (classification).
<code>datasets.fetch_openml</code> ([name, version, ...])	Fetch dataset from openml by name or dataset id.
<code>datasets.fetch_rcv1</code> ([data_home, subset, ...])	Load the RCV1 multilabel dataset (classification).
<code>datasets.fetch_species_distributions</code> ([...])	Loader for species distribution dataset from Phillips et.
<code>datasets.get_data_home</code> ([data_home])	Return the path of the scikit-learn data dir.
<code>datasets.load_boston</code> ([return_X_y])	Load and return the boston house-prices dataset (regression).
<code>datasets.load_breast_cancer</code> ([return_X_y])	Load and return the breast cancer wisconsin dataset (classification).
<code>datasets.load_diabetes</code> ([return_X_y])	Load and return the diabetes dataset (regression).
<code>datasets.load_digits</code> ([n_class, return_X_y])	Load and return the digits dataset (classification).
<code>datasets.load_files</code> (container_path[, ...])	Load text files with categories as subfolder names.
<code>datasets.load_iris</code> ([return_X_y])	Load and return the iris dataset (classification).
<code>datasets.load_linnerud</code> ([return_X_y])	Load and return the linnerud dataset (multivariate regression).
<code>datasets.load_sample_image</code> (image_name)	Load the numpy array of a single sample image
<code>datasets.load_sample_images</code> ()	Load sample images for image manipulation.
<code>datasets.load_svmlight_file</code> (f[, n_features, ...])	Load datasets in the svmlight / libsvm format into sparse CSR matrix
<code>datasets.load_svmlight_files</code> (files[, ...])	Load dataset from multiple files in SVMLight format
<code>datasets.load_wine</code> ([return_X_y])	Load and return the wine dataset (classification).
<code>datasets.mldata_filename</code> (dataname)	DEPRECATED: mldata_filename was deprecated in version 0.20 and will be removed in version 0.22

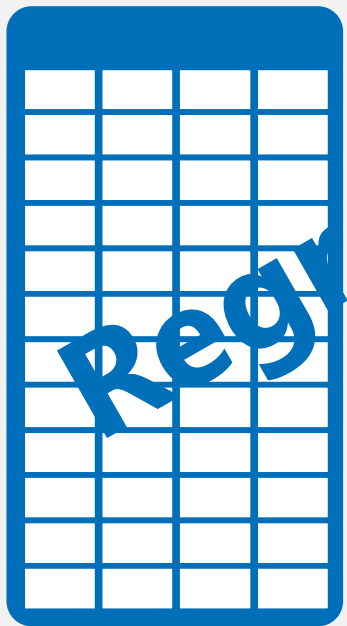
Samples generator

Sample Data

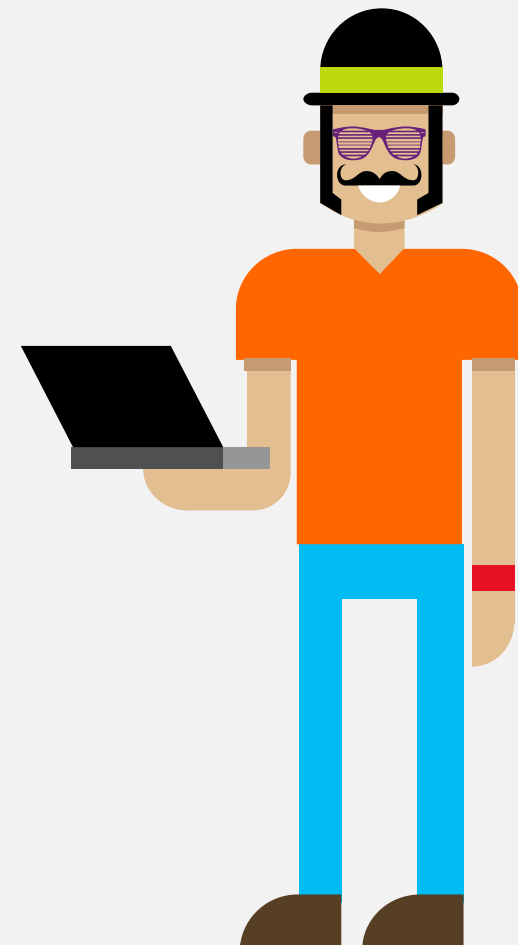
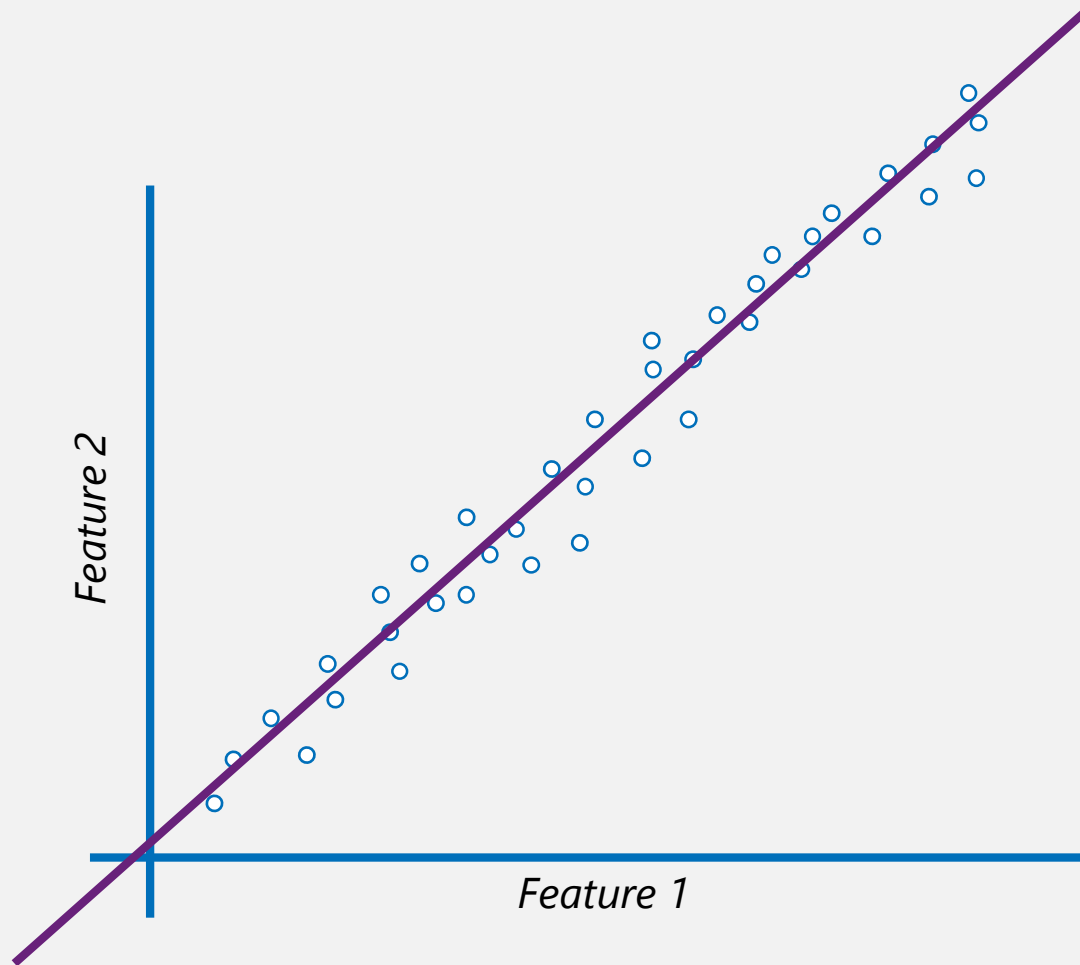


Diabetes data

data



Regression



5.2.3. 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:

Number of Instances:	
442	
Number of Attributes:	
First 10 columns are numeric predictive values	
Target:	Column 11 is a quantitative measure of disease progression one year after baseline
Attribute Information:	
	<ul style="list-style-type: none">• Age• Sex• Body mass index• Average blood pressure• S1• S2• S3• S4• S5• S6

Note: Each of these 10 feature variables have been mean centered and scaled by the standard deviation times n_{samples} (i.e. the sum of squares of each column totals 1).

Source URL: <http://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.
(http://web.stanford.edu/~hastie/Papers/LARS/LeastAngle_2002.pdf)

Demo – Exploring Diabetes Data

Lab 01 – Creating a model

15 minutes