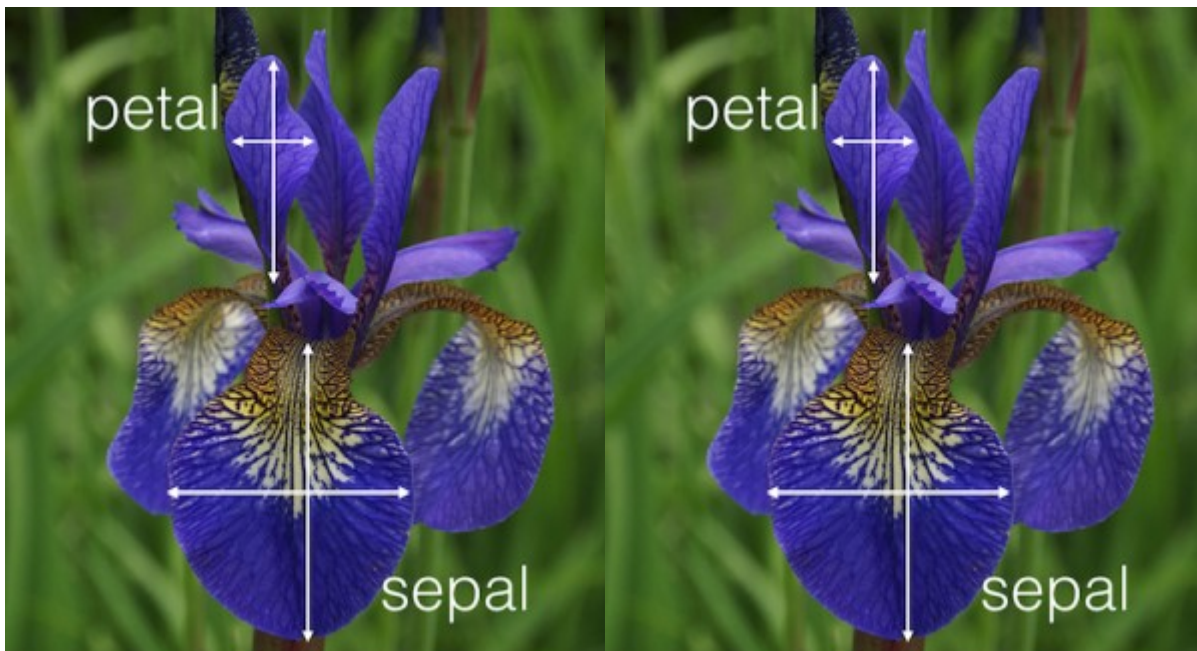


## Agenda

- What is the famous iris dataset, and how does it relate to machine learning?
- How do we load the iris dataset into scikit-learn?
- How do we describe a dataset using machine learning terminology?
- What are scikit-learn's four key requirements for working with data?

## Introducing the iris dataset



- 50 samples of 3 different species of iris (150 samples total)
- Measurements: sepal length, sepal width, petal length, petal width

## Loading the iris dataset into scikit-learn

---

In [2]:

```
# import load_iris function from datasets module  
from sklearn.datasets import load_iris
```

## Machine learning terminology

- Each row is an **observation** (also known as: sample, example, instance, record)
- Each column is a **feature** (also known as: predictor, attribute, independent variable, input, regressor, covariate)

In [5]:

```
# print the names of the four features
print iris.feature_names
['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)']
```

In [6]:

[illegible]

In [7]:

```
# print the encoding scheme for species: 0 = setosa, 1 = versicolor, 2 = virginica
print iris.target_names, '\n'
print iris.target_names[iris.target]
['setosa' 'versicolor' 'virginica']
```

[illegible]

- Each value we are predicting is the **response** (also known as: target, outcome, label, dependent variable)
- **Classification** is supervised learning in which the response is categorical
- **Regression** is supervised learning in which the response is ordered and continuous

## Requirements for working with data in scikit-learn

- ### 1. Features and response are **separate objects**

2. Features and response should be **numeric**
3. Features and response should be **NumPy arrays**
4. Features and response should have **specific shapes**

---

In [8]:

```
# check the types of the features and response
print type(iris.data)
print type(iris.target)
<type 'numpy.ndarray'>
<type 'numpy.ndarray'>
```

---

In [9]:

```
# check the shape of the features (first dimension = number of observations,
second dimensions = number of features)
print iris.data.shape
(150L, 4L)
```

---

In [10]:

```
# check the shape of the response (single dimension matching the number of
observations)
print iris.target.shape
(150L,)
```

---

In [12]:

```
# store feature matrix in "X"
X = iris.data

# store response vector in "y"
y = iris.target
```