# Why Python for Data Science?

- Easy to learn
- Powerful
- Scalable
- Many libraries
  - Machine Learning
  - Deep Learning
  - Visualization
- Python Community

Top 10 Data Science Programming Language by % of Job Ads in which the Language is Mentioned

| Language | % |
| --- | --- |
| Python | 90.4% |
| R | 73.4% |
| SQL | 58.5% |
| Scala | 21.3% |
| SAS | 18.1% |
| Java | 16.0% |
| Matlab | 14.9% |
| Hive | 13.8% |
| C/C++ | 7.4% |
| Pig | 7.4% |

ETS

# Overview of the Python Language: 1/2
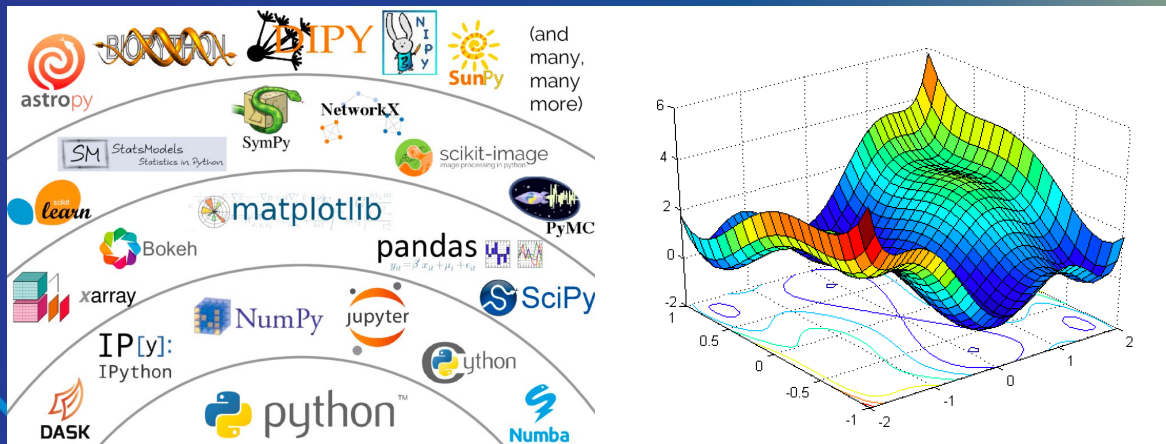
- Data Types
  - None = missing value
- Flow control
  - If-else
  - Loop
  - Function
- Data Structures
  - List []
  - Tuple ()
  - Set {}
  - Dictionary {}
- Printing & formatting
- Importing packages; getting help
- [Python Style Guide](#)
- List comprehension: [x**2 for x in range(10)]

# Overview of the Python Language: 2/2

- Recap: data structures, generators
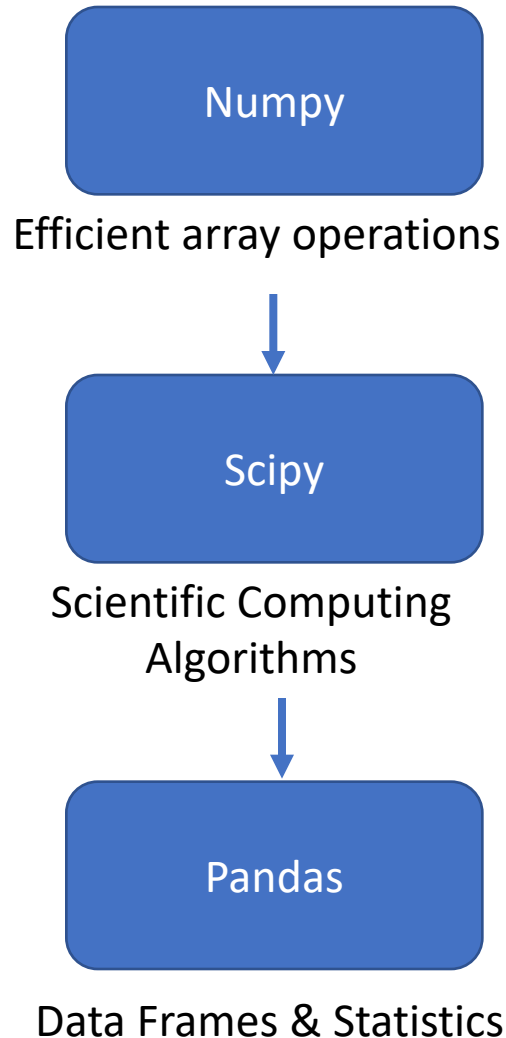- Index ranges, negative index, str as list
- Filtering lists with list comprehension
- Functional programming: map, reduce, filter
- timeit
- Zip

- Further reading:
  - https://docs.python.org/3/tutorial/

# Scientific Computing Packages

# Basic Data Science Packages

Numpy

Efficient array operations

Scipy

Scientific Computing Algorithms

Pandas

Data Frames & Statistics

# Package Highlights

- **Numpy**: [Tutorial Notebook](#)
- **Scipy**: [https://scipy-lectures.org/intro/scipy.html](https://scipy-lectures.org/intro/scipy.html)
  - scipy.linalg - linear algebra (1.6.3)
  - scipy.curvefit – curve fitting (1.6.5)
  - scipy.stats - statistics and random numbers (1.6.6)
    - Statistical testing
- **Pandas** [Tutorial Notebook](#)
  - Series
  - DataFrame
  - Reading a CSV file

- Further reading:
  - [https://numpy.org/doc/stable/user/quickstart.html](https://numpy.org/doc/stable/user/quickstart.html)
  - [10 minutes to pandas](#)

# Interfacing with other Languages

- Calling an R function from Python
  - [Call R for computing -> numpy array](Call R for computing -> numpy array)
  - [Call R -> pandas DataFrame](Call R -> pandas DataFrame)

# Session 2 - Exercises

1.  Complete numpy tutorial assignments in the bottom of the numpy_tutorial.ipynb notebook under the Session 2 materials on the Teams site.

2.  Calculate summary statistics: mean, std, 0%,10%,20%,…, percentiles (look for how to do that in the numpy documentation).

3.  Run a t-test of two Gaussian distribution means using scipy (see https://scipy-lectures.org/intro/scipy.html, Sec. 1.6.6.3.