# *End-to-End Machine Learning: The Python ML Ecosystem*
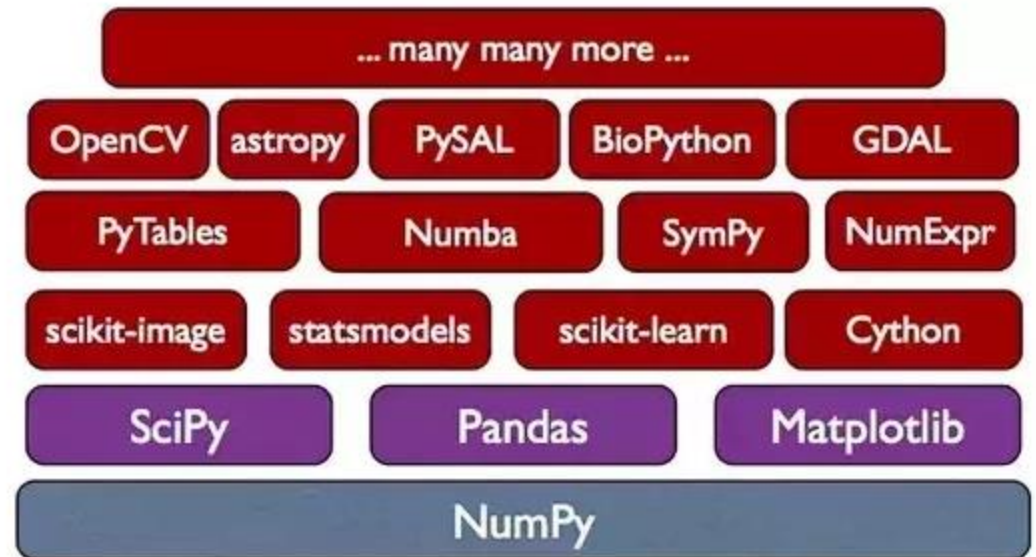
Glenn Bruns

CSUMB

# Learning outcomes

After this lecture you should be able to:

1. List the packages and tools we'll use for machine learning with Python

2. Explain the role of each package/tool

3. Load csv data using Pandas

4. Do basic data exploration with NumPy

5. Edit and run Python programs using Spyder

# The Python machine learning ecosystem

- ☐ NumPy
- ☐ SciPy
- ☐ Scikit-Learn
- ☐ Pandas
- ☐ Matplotlib



source: quora.com/What-is-the-relationship-among-NumPy-SciPy-Pandas-and-Scikit-learn-and-when-should-I-use-each-one-of-them

# What the packages do

NumPy: supports large, multi-dimensional arrays & matrices, plus functions on them

SciPy: builds on NumPy; adds algorithms and convenience functions

Scikit-Learn: builds on SciPy; adds machine learning

Pandas: data manipulation and analysis
- Series – labelled, 1D array
- DataFrame – labelled, 2D data

Matplotlib: plotting

# Reminder: data science process

- ☐ Define goal

- ☐ Data acquisition and cleaning

- ☐ Data exploration

- ☐ Machine learning

- ☐ Present results

# Reading a CSV file

One way, using Pandas:

```python
import numpy as np
import pandas as pd
input_file = "C:/Users/…/titanic-data/train.csv"
dat = pd.read_csv(input_file)
```

dat is a pandas DataFrame

(see pandas.pydata.org/pandas-docs/stable/dsintro.html)

# Examining the data in Spyder

# Basic data exploration

```python
# similar to R's 'summary'
dat.info()

# number of elements
dat.size

# dimensions
dat.shape

# number of rows
len(dat)
dat.shape[0]    # alternative way

# beginning of the data frame
dat.head()
```

# Running code in Spyder

☐  As in R, highlight code and hit ctrl-Enter

☐  To clear your environment, enter `%reset` in the console

# Running dat.info()

```
In [227]: dat.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
PassengerId    891 non-null int64
Survived       891 non-null int64
Pclass         891 non-null int64
Name           891 non-null object
Sex            891 non-null object
Age            714 non-null float64
SibSp          891 non-null int64
Parch          891 non-null int64
Ticket         891 non-null object
Fare           891 non-null float64
Cabin          204 non-null object
Embarked       889 non-null object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.6+ KB
```

# Summary

- ☐ Packages in the Python ML ecosystem

- ☐ Reading CSV data and basic data exploration

- ☐ Edit and run Python programs using Spyder