# Python/Numpy for High-Performance Numerical Processing

Jim Pivarski
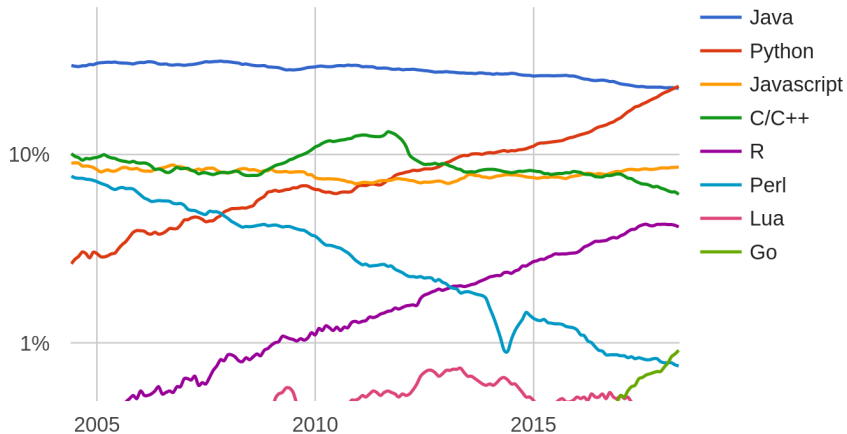
Princeton University

November 15, 2018

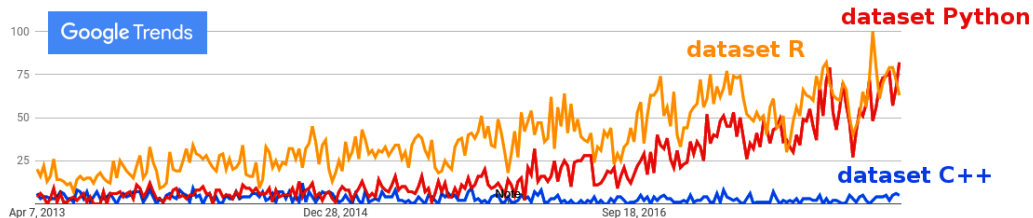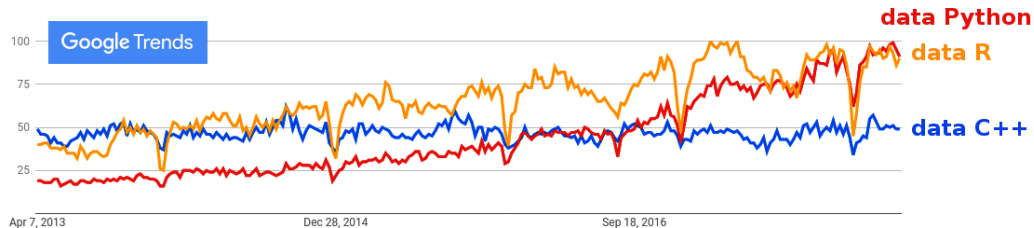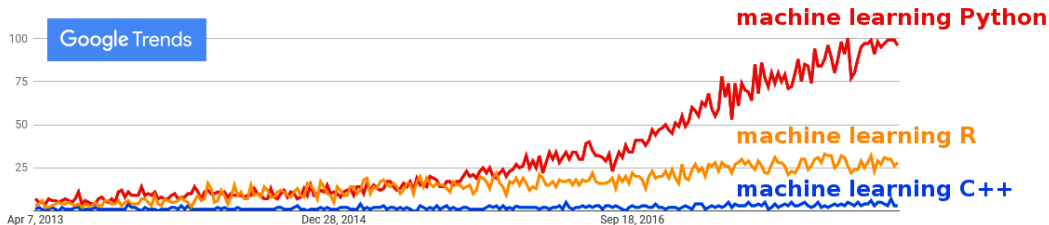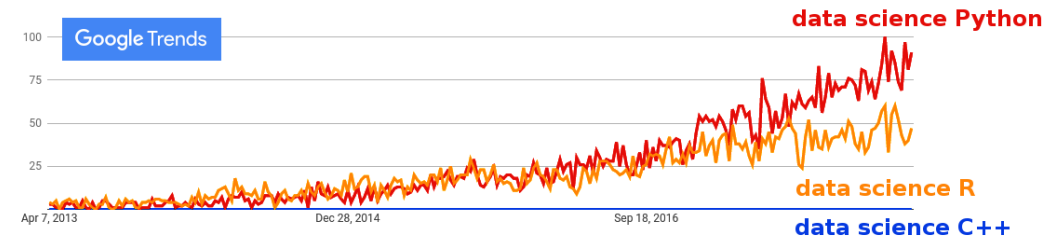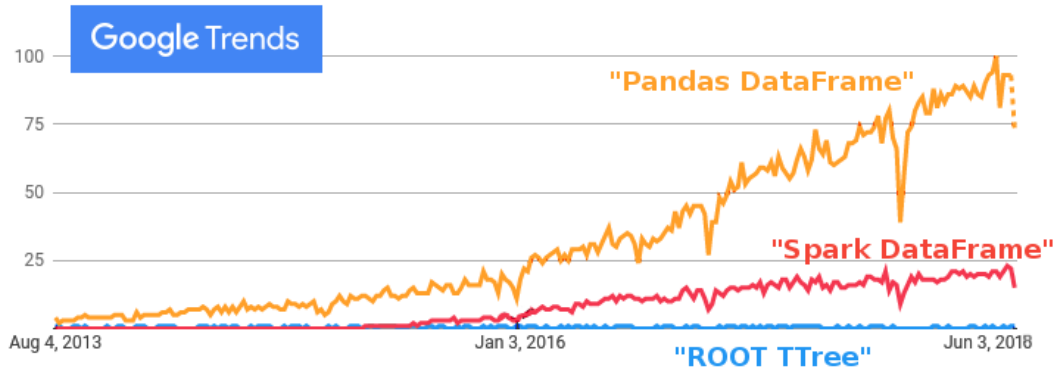**PYPL PopularitY of Programming Language**



http://pypl.github.io/PYPL.html

All of the machine learning libraries I could find either have a Python interface or are primarily/exclusively Python.

Mentions of Software in Astronomy Publications:

Thanks to Juan Nunez-Iglesias, Thomas P. Robitaille, and Chris Beaumont.

Compiled from NASA ADS (code).

## Python's Scientific Stack



The Unexpected
Effectiveness of
Python in Science

Jake VanderPlas @jakevdp
PyCon 2017



python™

## Python's Scientific Stack

If you're used to writing your own code, searching for tools is eye-opening: you learn what's unique about what you do and what isn't.

In science, we often have to scale up analyses to large datasets.

In science, we often have to scale up analyses to large datasets.

10% faster doesn't mean much, but the difference between "five minutes" and "overnight" is life-changing.

In science, we often have to scale up analyses to large datasets.

10% faster doesn't mean much, but the difference between
"five minutes" and "overnight" is life-changing.

That's the scale we're talking about between C and Python.

In science, we often have to scale up analyses to large datasets.

10% faster doesn't mean much, but the difference between
"five minutes" and "overnight" is life-changing.

That's the scale we're talking about between C and Python.

But we also need the interactivity of a dynamic language to *develop* the analysis.
("If we knew what we were doing, it wouldn't be called research.")

Drive to the airport by car, then take a plane.

Small-scale *project organization* in Python, ignoring performance entirely.

Run over *big data* in compiled code, tuning performance until it no longer matters.

**Google** Code

**pyminuit**
*Minuit numerical function minimization in Python*

Search Projects | Search the Web

Project Home | Downloads | Wiki | Issues | Source

## PyMinuit

*Minuit numerical function minimization in Python*

### Minuit

Minuit has been the standard package for minimizing general N-dimensional functions in high-energy physics since its introduction in 1972. It features a robust set of algorithms for optimizing the search, correcting mistakes, and measuring non-linear error bounds. It is the minimization engine used behind-the-scenes in most high-energy physics curve fitting applications.

**New:** more robust installation instructions!

### Python interface

PyMinuit is an extension module for Python that passes low-level Minuit functionality to Python functions. Interaction and data exploration is more user-friendly, in the sense that the user is protected from segmentation faults and index errors, parameters are referenced by their names, even in correlation matrices, and Python exceptions can be passed from the objective function during the minimization process. This extension module also makes it easier to calculate Minos errors and contour curves at an arbitrary number of sigmas from the minimum, and features a new N-dimensional scanning utility.

**License:** GNU General Public License v2

**Labels:** python, minuit, optimization, computation, HEP, math, fitting, physics

**Featured Downloads:**    Show all
⬇ Minuit-1_7_9.tar.gz
⬇ pyminuit-1.0.2.tgz

**Featured Wiki Pages:**    Show all
FunctionReference
GettingStartedGuide
HowToInstall

**Links:**
- Official Minuit homepage at CERN
- Minuit documentation through the ages

**Project owners:**    Join project
jpivarski

# The key to ecosystem development was a common array library

1994 **Python** 1.0 released.

1995 First array package: **Numeric** (a.k.a. Numerical, Numerical Python, NumPy).

2001 Diverse scientific codebases merged into **SciPy**.

2003 **Matplotlib**

2003 Numeric was limited; **numarray** appeared as a competitor with more features (memory-mapped files, alignment, record arrays).

2005 Two packages were incompatible; could not integrate numarray-based code into SciPy. Travis Oliphant merged the codebases as **Numpy**.

2008 **Pandas**

2010 **Scikit-Learn**

2011 **AstroPy**

2012 **Anaconda**

2014 **Jupyter**

2015 **Keras**

1994    **Python** 1.0 released.

1995    First array package: **Numeric** (a.k.a. Numerical, Numerical Python, NumPy).

2001    Diverse scientific codebases merged into **SciPy**.

2003    **Matplotlib**

2003    Numeric was limited; **numarray** appeared as a competitor with more features (memory-mapped files, alignment, record arrays).

2005    Two packages were incompatible; could not integrate numarray-based code into SciPy. Travis Oliphant merged the codebases as **Numpy**.

2008    **Pandas**

2010    **Scikit-Learn**

2011    **AstroPy**

2012    **Anaconda**

2014    **Jupyter**

2015    **Keras**

> The scientific Python ecosystem could have failed before it started if the Numeric/numarray split hadn't been resolved!

```
>>> import numpy
>>> a = numpy.arange(12)
>>> a
array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11])
>>> a.shape = (3, 4)
>>> a
array([[ 0,  1,  2,  3],
       [ 4,  5,  6,  7],
       [ 8,  9, 10, 11]])
>>> a.sum(axis=0)
array([12, 15, 18, 21])
>>> a.min(axis=1)
array([0, 4, 8])
>>> a**2
array([[  0,   1,   4,   9],
       [ 16,  25,  36,  49],
       [ 64,  81, 100, 121]])
>>> numpy.sqrt(a)
array([[0.        , 1.        , 1.41421356, 1.73205081],
       [2.        , 2.23606798, 2.44948974, 2.64575131],
       [2.82842712, 3.        , 3.16227766, 3.31662479]])
```

# The Numpythonic mindset

Although you can write Python `for` loops over Numpy arrays, you don't reap the benefit unless you express your calculation in Numpy universal functions (ufuncs).

```
pz = numpy.empty(len(pt))
for i in range(len(pt)):
    pz[i] = pt[i]*numpy.sinh(eta[i])
```

VS

```
pz = pt * numpy.sinh(eta)
```

$\mathcal{O}(N)$ Python bytecode instructions, type-checks, interpreter locks.

$\mathcal{O}(1)$ Python bytecode instructions, type-checks, interpreter locks.

$\mathcal{O}(N)$ statically typed, probably vectorized native bytecode operations on contiguous memory.

Although you can write Python `for` loops over Numpy arrays, you don't reap the benefit unless you express your calculation in Numpy universal functions (ufuncs).

```
pz = numpy.empty(len(pt))
for i in range(len(pt)):
    pz[i] = pt[i]*numpy.sinh(eta[i])
```

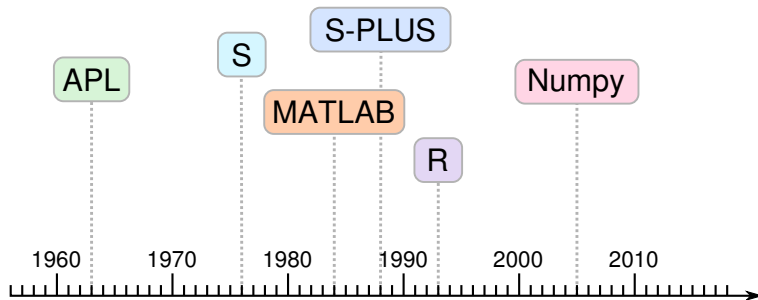vs

```
pz = pt * numpy.sinh(eta)
```

$\mathcal{O}(N)$ Python bytecode instructions, type-checks, interpreter locks.

$\mathcal{O}(1)$ Python bytecode instructions, type-checks, interpreter locks.

$\mathcal{O}(N)$ statically typed, probably vectorized native bytecode operations on contiguous memory.
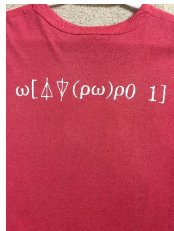
In other words, a Single (Python) Instruction on Multiple Data.
Conceptually similar to SIMD, the program flow of GPUs.

**APL**, "A Programming Language" introduced the idea of single commands having sweeping effects across large arrays.



All members of the APL family are intended for interactive data analysis.
Numpy, however, is a library in a general-purpose language, not a language in itself.

APL pioneered conciseness;
discovered the mistake of being too concise.



$\omega[\not\Delta\not\Psi(\rho\omega)\rho0\ 1]$

Conway's Game of Life was one line of code:

$\texttt{life} \leftarrow \{\uparrow 1\ \ \omega \vee .\wedge 3\ \ 4 = +/,^{-}1\ \ 0\ \ 1 \circ .\Theta^{-}1\ \ 0\ \ 1 \circ .\Phi \subset \omega\}$

"Map" was implicit, "reduce" was a slash, functions were symbols. For example:

| APL | Numpy |
|---|---|
| $\texttt{m} \leftarrow +/(3 + \iota 4)$ | `m = (numpy.arange(4) + 3).sum()` |

As an array abstraction, Numpy presents a high-level way
for users to think about vectorization.

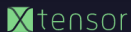Vectorization is key to using GPUs and modern CPUs efficiently.

| | |
|---|---|
| 10:00 — Intro talk | 1-intro.pdf |
| Just Numpy | 2-just-numpy.ipynb |
| 11:00 | |
| 12:00 — Lunch | |
| 1:00 — Numpy ecosystem talk | 3-ecosystem.pdf |
| Pandas | 4-pandas.ipynb |
| 2:00 — Dask & multiprocessing | 5-dask.ipynb |
| Coffee break | |
| 3:00 — Numba, Cython, pybind11 | 6-compilers.ipynb |
| CuPy, Numba-GPU, PyCUDA | 7-gpu.ipynb |
| ctypes & low-level hackery | 8-low-level.ipynb |
| 4:00 | |

Skills-based Numpy tutorial with a couple of exercises in the morning: how to think in SIMD.

Overview of libraries in the afternoon: where to look for solutions to your problems.