



Python's Data Science Stack

Jake VanderPlas @jakevdp
JSM, July 31, 2016

Python is *not* a statistical computing language!

**Python is *not* a statistical
computing language!**

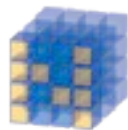
**. . . and this may be its
greatest strength as a
language for statistical
computing.**

A Quick Tour of Python's Data Science Stack

Python's Data Science Stack



IP[y]:
IPython



NumPy



Cython

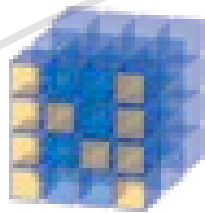


python™



NumPy = Numerical Python

Efficient array storage, manipulation, and computation



NumPy ython

IP[y]:
IPython

 python™



NumPy = Numerical Python

Efficient array storage, manipulation, and computation

```
In [1]: import numpy as np
```

```
# Create a 5x5 uniform random matrix
```

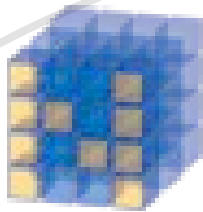
```
M = np.random.rand(5, 5)
```

```
# Compute the SVD
```

```
U, S, VT = np.linalg.svd(M)
```

```
print(S)
```

```
[ 2.46102945  0.94542853  0.53550015  0.20705388  0.13071452]
```



NumPy  ython

IP[y]:
IPython

 python™





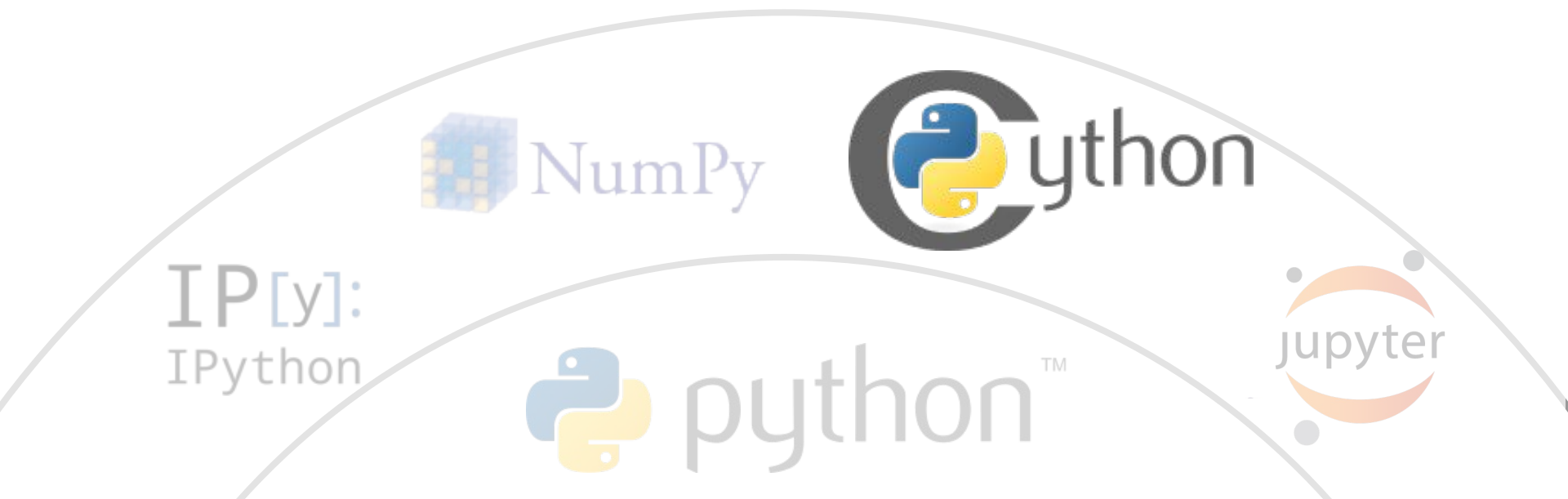
IP[y]:
IPython



Cython = C + Python

Super-set of the Python language that allows easy interfacing with C & Fortran libraries (e.g. BLAS, LAPACK, etc.) and also fast Python code.

Drives many of the packages in the data science stack.



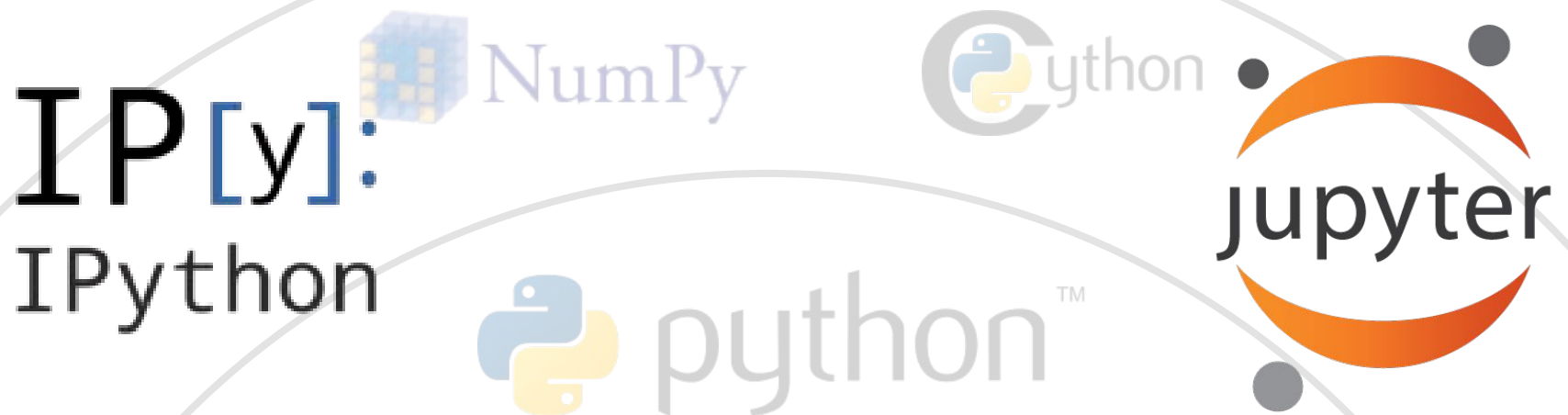


IP[y]:
IPython

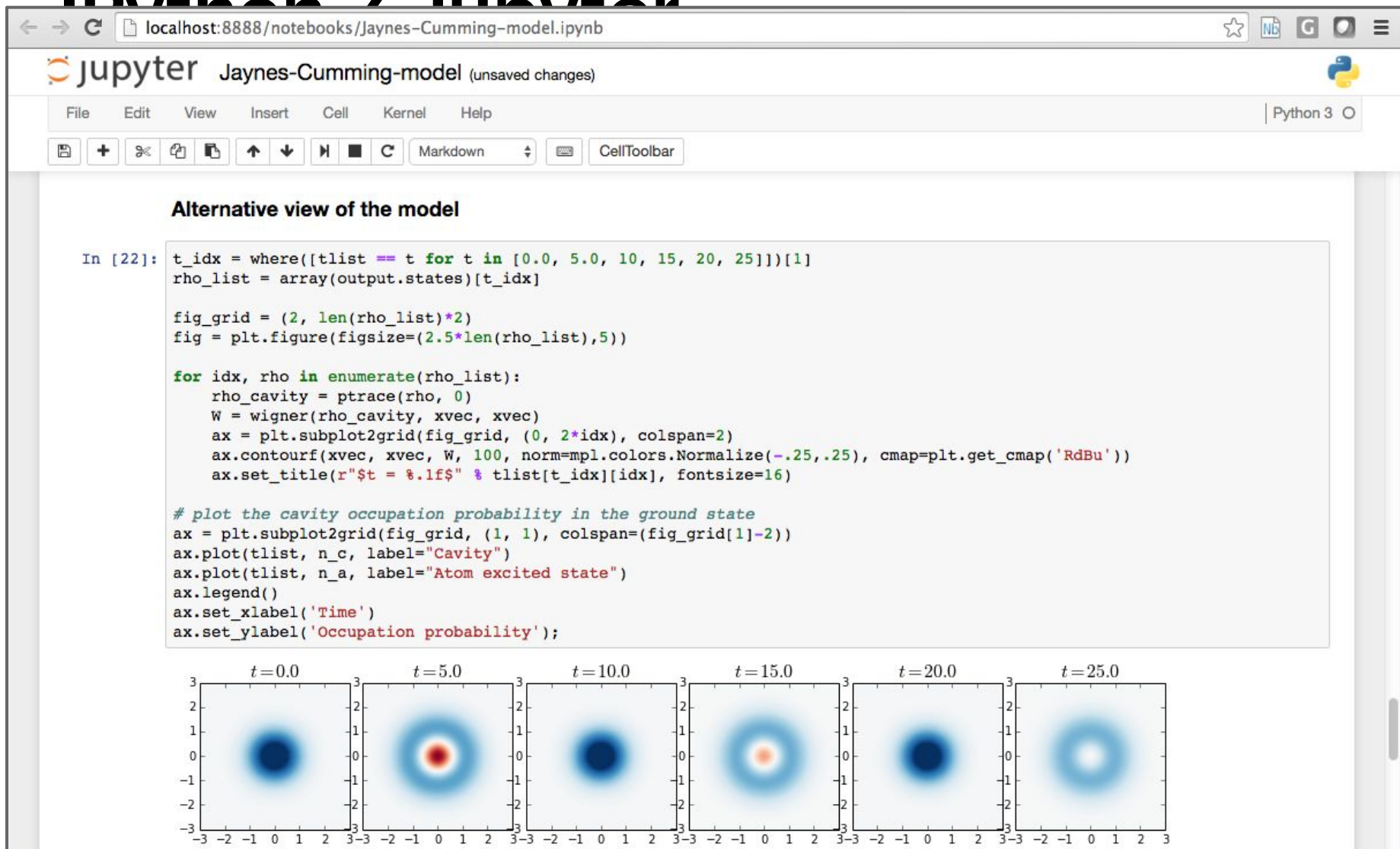


IPython / Jupyter

Terminal, development environment, Notebooks, and more for efficient use of Python in day-to-day work



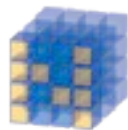
IPython / Jupyter



IPython

python™

IP[y]:
IPython

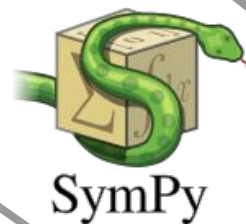


NumPy



python™





SciPy

Provides an interface to common scientific computing Tasks, including wrappers of many NetLib packages.



SciPy

Provide
Tasks

List from <http://docs.scipy.org/doc/scipy/reference/>

- Special functions (**scipy.special**)
- Integration (**scipy.integrate**)
- Optimization (**scipy.optimize**)
- Interpolation (**scipy.interpolate**)
- Fourier Transforms (**scipy.fftpack**)
- Signal Processing (**scipy.signal**)
- Linear Algebra (**scipy.linalg**)
- Sparse Eigenvalue Problems with ARPACK
- Compressed Sparse Graph Routines (**scipy.sparse.csgraph**)
- Spatial data structures and algorithms (**scipy.spatial**)
- Statistics (**scipy.stats**)
- Multidimensional image processing (**scipy.ndimage**)
- File IO (**scipy.io**)

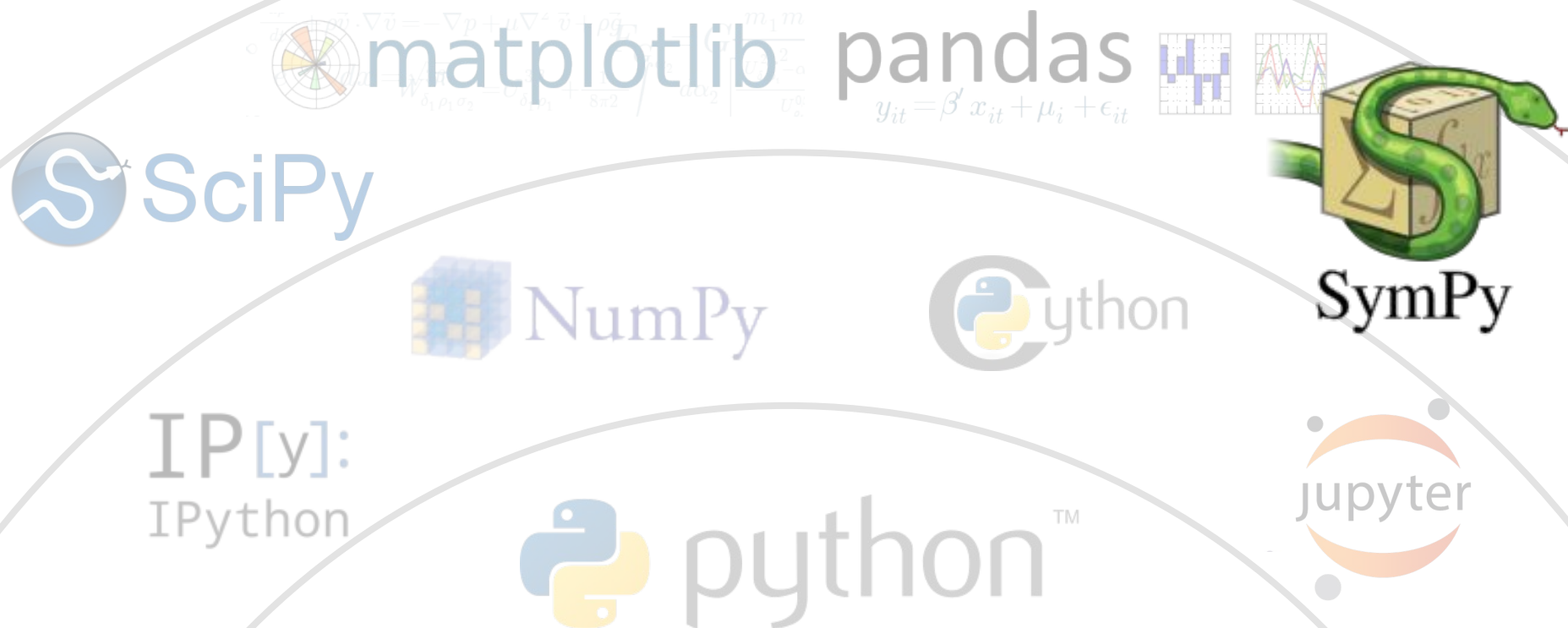
IP[y]:
IPython





Sympy

Library for symbolic computation: algebraic operations, differentiation & integration, optimization, etc.



Sympy

Library for symbolic computation: algebraic operations, differentiation & integration, optimization, etc.

Polynomials and rational functions

SymPy does not expand brackets automatically. The function `expand` is used for this.

```
In [6]: a=(x+y-z)**6  
a
```

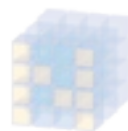
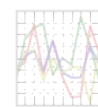
```
Out[6]: (x + y - z)6
```

```
In [7]: a=expand(a)  
a
```

```
Out[7]: x6 + 6x5y - 6x5z + 15x4y2 - 30x4yz + 15x4z2 + 20x3y3 - 60x3y2z + 60x3yz2 - 20x3z3  
        + 15x2y4 - 60x2y3z + 90x2y2z2 - 60x2yz3 + 15x2z4 + 6xy5 - 30xy4z + 60xy3z2  
        - 60xy2z3 + 30xyz4 - 6xz5 + y6 - 6y5z + 15y4z2 - 20y3z3 + 15y2z4 - 6yz5 + z6
```



pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$


NumPy



ython



SymPy

IP[y]:
IPython

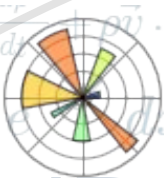


python™

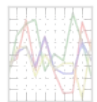
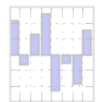


matplotlib

Matlab-inspired plotting and visualization



matplotlib



NumPy



Cython



SymPy

IP[y]:
IPython



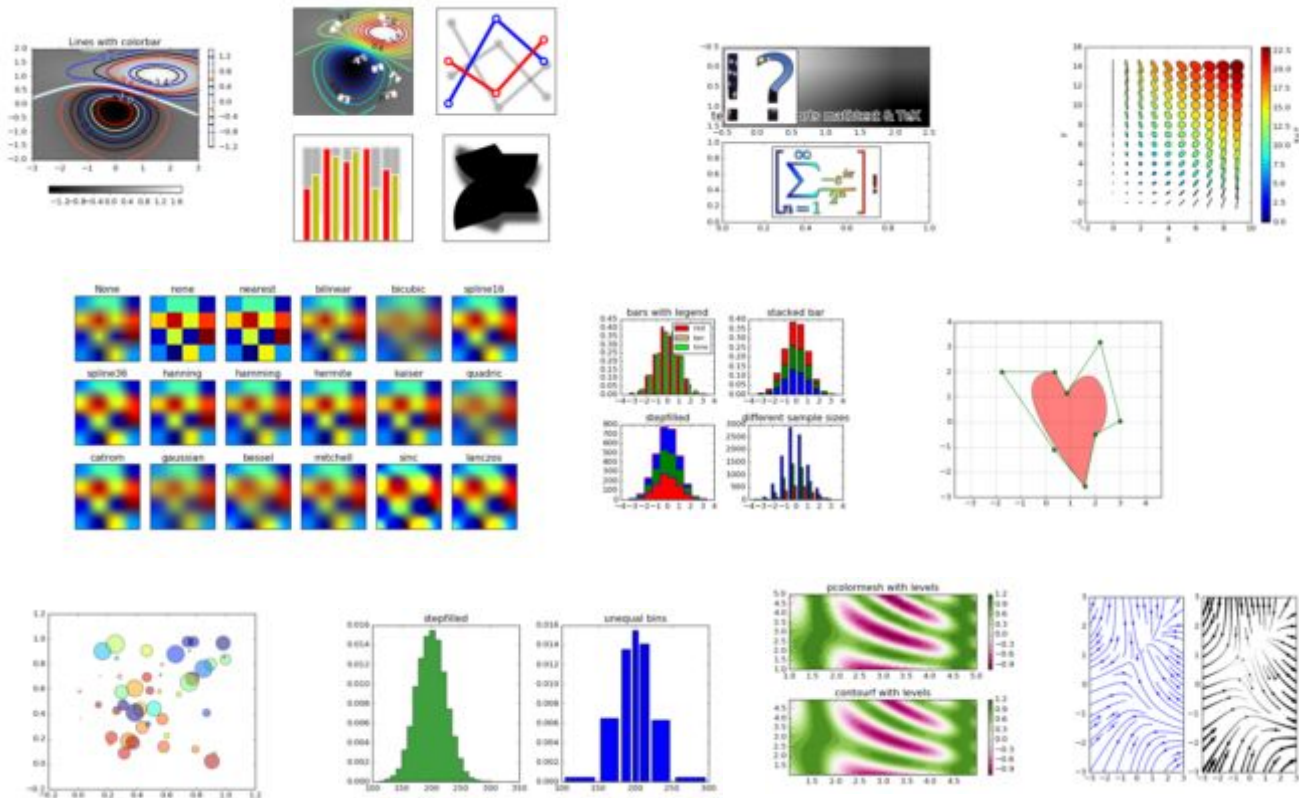
python™



matplotlib

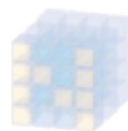
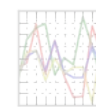
Matlab-inspired plotting and visualization

From <http://matplotlib.org/gallery.html>





pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$


NumPy



ython



SymPy

IP[y]:
IPython

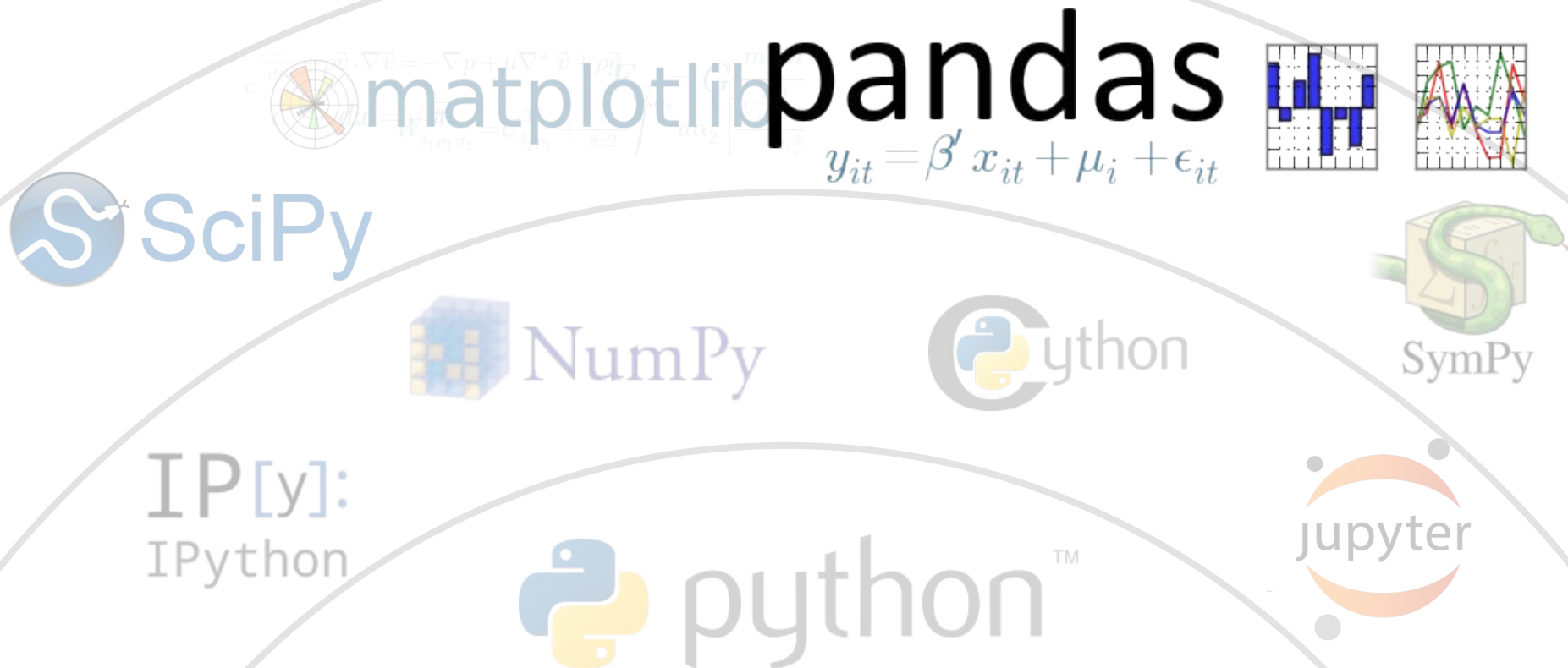


python™



Pandas

R-inspired DataFrames & associated functionality
(data munging & cleaning, group-by & transformations,
and much more)



```
In [1]: import pandas as pd
data = pd.read_csv('iris.csv')
data.head()
```

Out[1]:

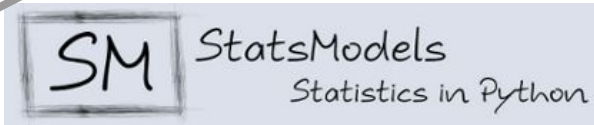
	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa

```
In [2]: data.groupby('Species').mean()
```

Out[2]:

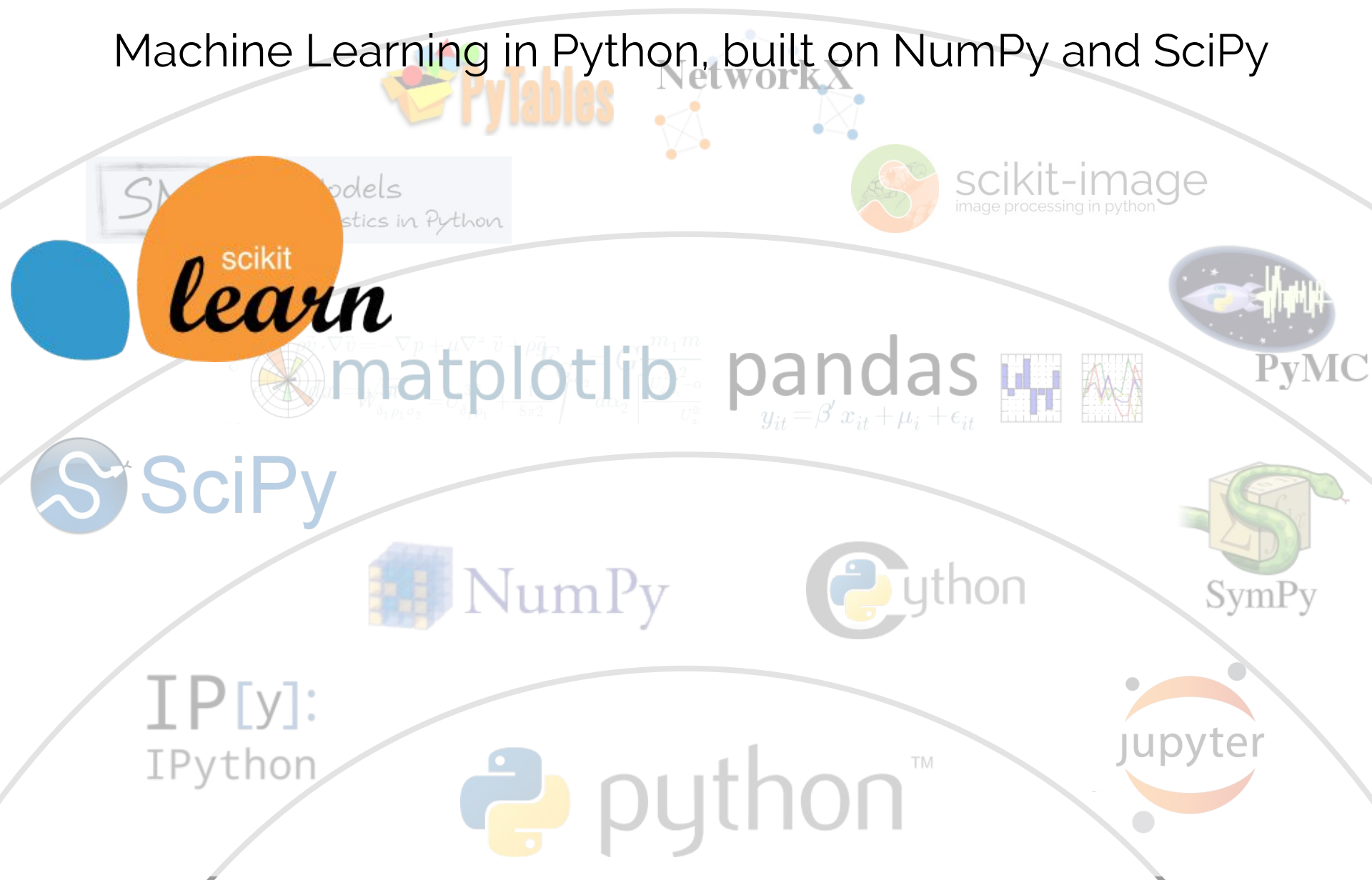
	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width
Species				
setosa	5.006	3.428	1.462	0.246
versicolor	5.936	2.770	4.260	1.326
virginica	6.588	2.974	5.552	2.026





Scikit-Learn

Machine Learning in Python, built on NumPy and SciPy



Scikit-Learn

Machine Learning in Python, built on NumPy and SciPy

```
In [3]: from sklearn.ensemble import RandomForestClassifier

features = data.drop('Species', axis=1)
labels = data['Species']

model = RandomForestClassifier()
model.fit(features, labels)

predicted = model.predict(features.iloc[:5])
print(predicted)

['setosa' 'setosa' 'setosa' 'setosa' 'setosa']
```

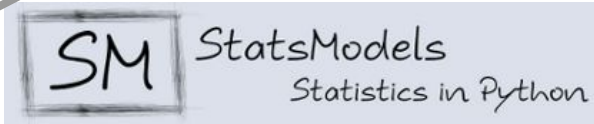
IP[y]:
IPython

 python™

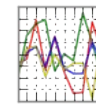
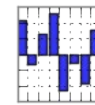
 jupyter



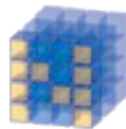
scikit-image
image processing in python



pandas
 $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$



PyMC



NumPy



SymPy

IP[y]:
IPython

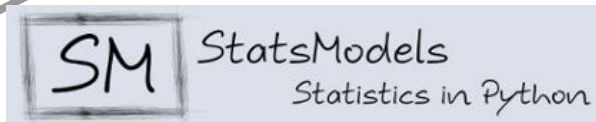




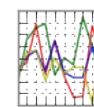
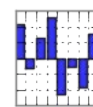
(and many, many more)



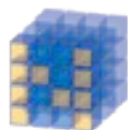
scikit-image
image processing in python



pandas
 $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$



PyMC



NumPy



SymPy

IP[y]:
IPython



Recent-ish Developments

- **Dask**: Parallelization of Data & Computation
- **Numba**: LLVM compilation of Python code
- **Jupyter Lab**: interactive & extensible polyglot development environment
- **Altair**: Declarative Visualization based on Vega-Lite

Dask: Parallel Computation for Distributed Arrays & DataFrames

With minimal changes to your NumPy & Pandas expressions, parallelize your computations over distributed data!

Dask: Parallel Computation for Distributed Arrays & DataFrames

A straightforward NumPy computation:

```
In [3]: import numpy as np

# create an array of normally-distributed random numbers
a = np.random.randn(1000)

# multiply this array by a factor
b = a * 4

# find the minimum value
b_min = b.min()
print(b_min)

-11.4051061336
```

Dask: Parallel Computation for Distributed Arrays & DataFrames

Dask uses the same expressions ...

```
In [4]: import dask.array as da

# create a dask array from the above array
a2 = da.from_array(a, chunks=200)

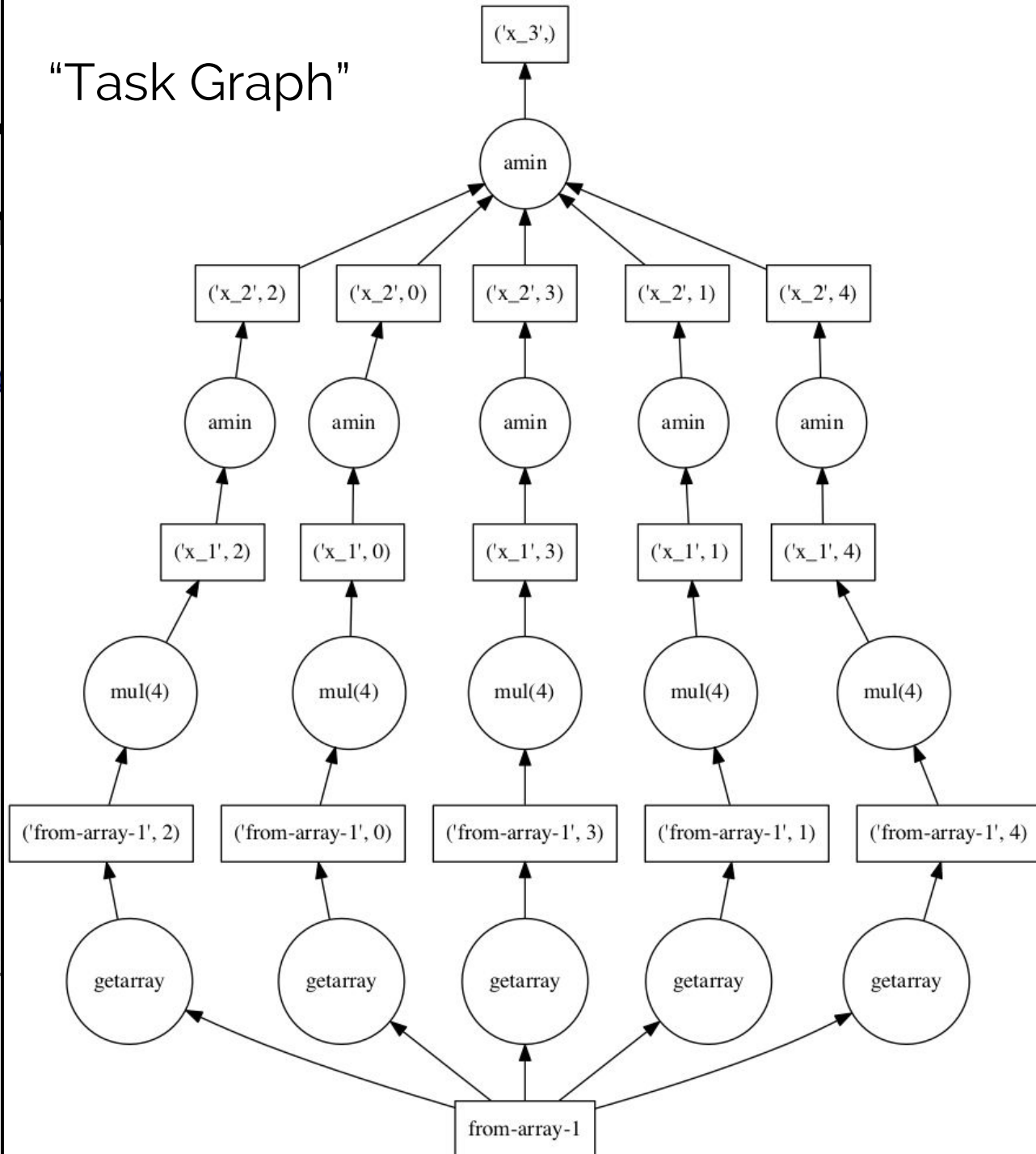
# multiply this array by a factor
b2 = a2 * 4

# find the minimum value
b2_min = b2.min()
print(b2_min)

dask.array<x_3, shape=(), chunks=(), dtype=float64>
```

Task Graph

"Task Graph"



Dask: Parallel Computation for Distributed Arrays & DataFrames

```
In [6]: b2_min.compute()
```

```
Out[6]: -11.405106133564583
```

Numba: JIT-compilation of Python code

With a simple decorator, Python is compiled to LLVM and executes at near C/Fortran speed!

```
def fib(n):  
    a, b = 0, 1  
    for i in range(n):  
        a, b = b, a + b  
    return a
```

```
%timeit fib(50)
```

```
100000 loops, best of 3: 3.83 µs per loop
```

Still some features missing, but very promising (see my blog posts for some examples).

<http://numba.pydata.org/>

Numba: JIT-compilation of Python code

With a simple decorator, Python is compiled to LLVM and executes at near C/Fortran speed!

```
@numba.jit
def fib(n):
    a, b = 0, 1
    for i in range(n):
        a, b = b, a + b
    return a

%timeit(fib(50))
```

1 loops, best of 3: 468 ns per loop

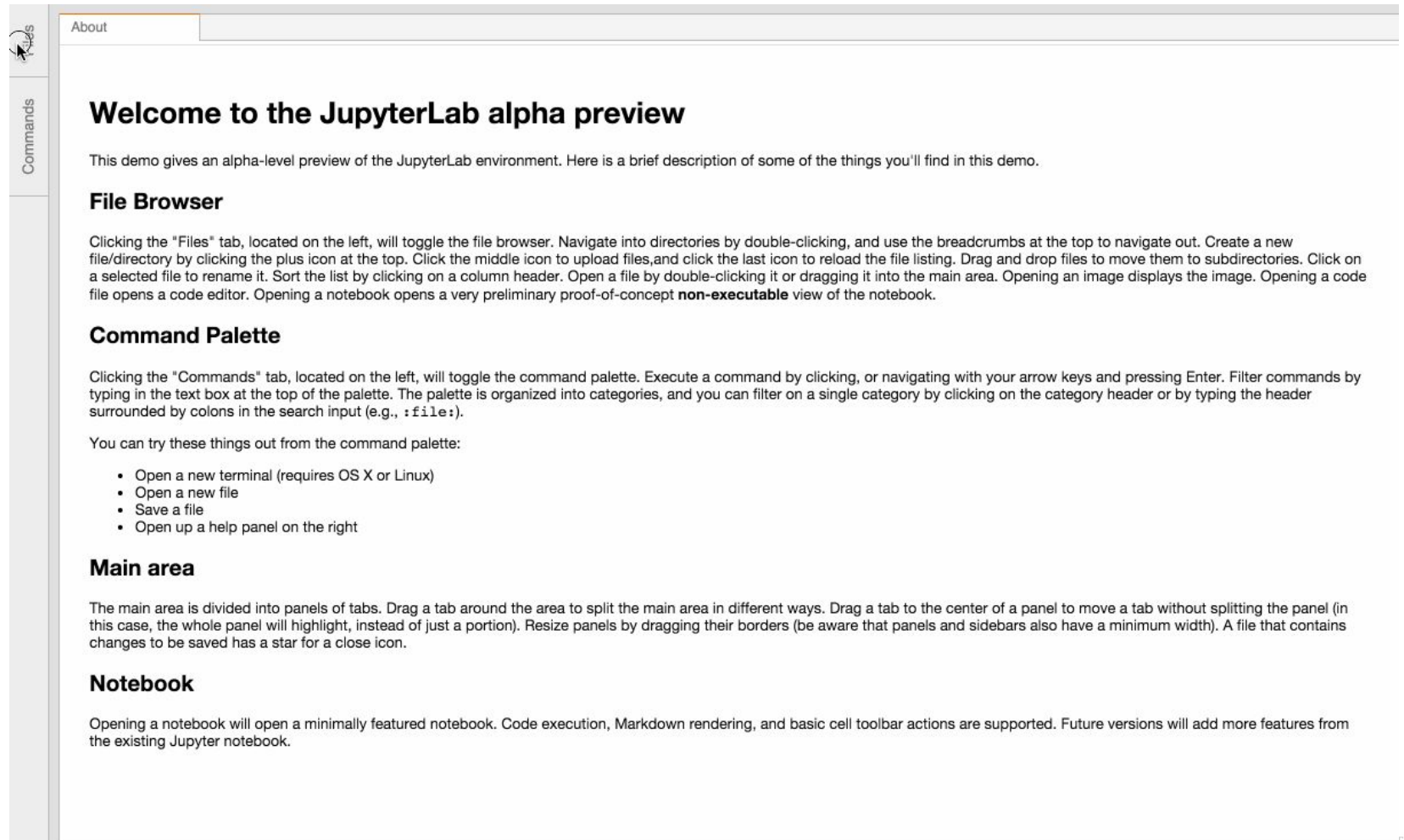
20x speedup!

Still some features missing, but very promising (see my blog posts for some examples).

<http://numba.pydata.org/>

Jupyter Lab

Jupyter beyond notebooks: extensible cross-platform interactive computing environment (release soon!)



Thank You!



Email: jakevdp@uw.edu



Twitter: [@jakevdp](https://twitter.com/jakevdp)



Github: [jakevdp](https://github.com/jakevdp)



Web: <http://vanderplas.com>



Blog: <http://jakevdp.github.io>