

DUSTIN MICHELS
DEVX, CARLETON
16 OCTOBER 2017

PYTHON FOR DATA SCIENCE

I. INTRODUCTION

WHAT IS DATA SCIENCE?



Figure P-1. Drew Conway's Data Science Venn Diagram

“DATA SCIENCE COMPRISES THREE DISTINCT AND OVERLAPPING AREAS: THE SKILLS OF A STATISTICIAN WHO KNOWS HOW TO MODEL AND SUMMARIZE DATASETS (WHICH ARE GROWING EVER LARGER); THE SKILLS OF A COMPUTER SCIENTIST WHO CAN DESIGN AND USE ALGORITHMS TO EFFICIENTLY STORE, PROCESS, AND VISUALIZE THIS DATA; AND THE DOMAIN EXPERTISE—WHAT WE MIGHT THINK OF AS “CLASSICAL” TRAINING IN A SUBJECT—NECESSARY BOTH TO FORMULATE THE RIGHT QUESTIONS AND TO PUT THEIR ANSWERS IN CONTEXT.”

(JAKE VANDERPLAS, PYTHON FOR DATA SCIENCE HANDBOOK, XI)

I. INTRODUCING DATA SCIENCE

DRAWING ON WORK OF JAKE VANDERPLAS

- ▶ “Why Astronomers Love Python”, PyCon 2017, <https://youtu.be/IWl6d7mkru4>.
- ▶ “Python for Data Science Handbook” (2017), O’Reily.



KEY IDEA: PUTTING THE “SCIENCE” IN COMPUTER SCIENCE

- ▶ Programming is integral part of scientific work, but not always done scientifically...

KEY IDEA: PUTTING THE “SCIENCE” IN COMPUTER SCIENCE

- ▶ Programming is integral part of scientific work, but not always done scientifically...
- ▶ Python and its packages are becoming popular choice
 - ▶ Quality
 - ▶ Transparency
 - ▶ Reproducibly
 - ▶ Collaboration

TOOLS

▶ **Interactivity / collaboration**

- ▶ IPython
- ▶ Jupyter

▶ **Data wrangling / analysis**

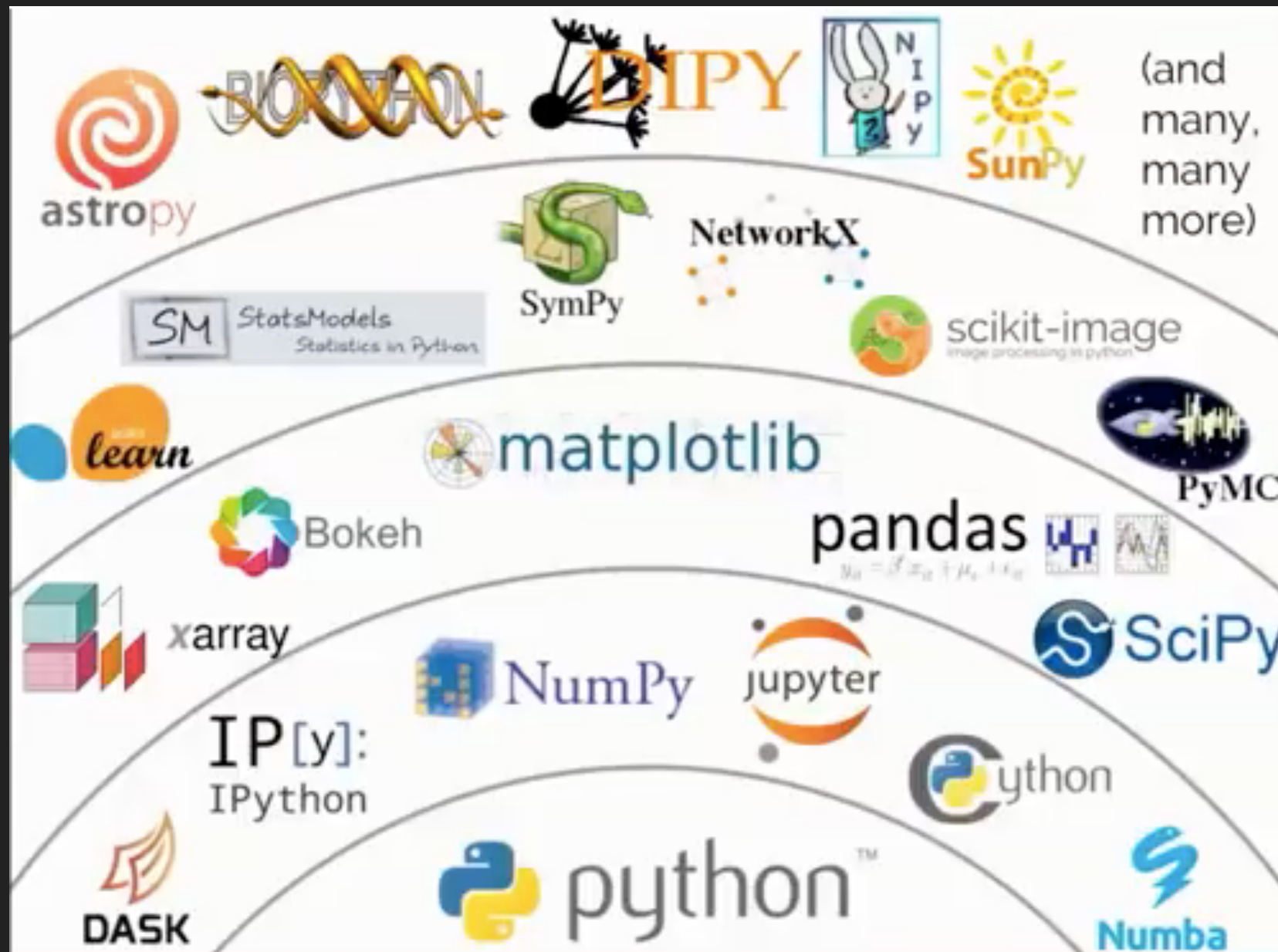
- ▶ Numpy
- ▶ Pandas

▶ **Visualization**

- ▶ Matplotlib
- ▶ Seaborn

I. INTRODUCING DATA SCIENCE

“WHY ASTRONOMERS LOVE PYTHON” (PYCON, 2017)



TODO #1

1. Go to: <https://tinyurl.com/devx-data>
2. Git clone or download folder

TODO #2

Option A:

Conda/ Miniconda

Option B:

```
$ pip3 install --upgrade pip`
```

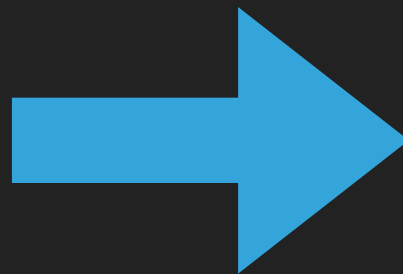
```
$ pip3 install numpy  
matplotlib ipython jupyter
```

II. INTERACTIVITY / COLLABORATION IPYTHON / JUPYTER

IPYTHON: IMPROVED “REPL”

```
>>> 2 + 2
4
>>> print("hello")
hello
>>> def hello():
...     print("hello there!")
...
>>> hello()
hello there!
>>> █
```

python REPL



```
[In [1]: 2 + 2
Out[1]: 4

[In [2]: print("hello")
hello

[In [3]: def hello():
...:     print("hello there!")
...:

[In [4]: hello()
hello there!]
```

ipython REPL

IPYTHON: KEY FEATURES

- ▶ Input/output numbering!
- ▶ Syntax highlighting!
- ▶ Multi-line entry / recall!
- ▶ Tab completion!
- ▶ Easily pull up documentation with "?" !

IPYTHON: METHOD DISCOVERY

```
[In [8]: h = "hello there!"]
```

```
[In [9]: h.]
```

capitalize()	encode()	format()	isalpha()	islower()
casefold()	endswith()	format_map()	isdecimal()	isnumeric()
center()	expandtabs()	index()	isdigit()	isprintable()
count()	find()	isalnum()	isidentifier()	isspace()

```
[In [9]: h.capitalize?
```

```
Docstring:
```

```
S.capitalize() -> str
```

```
Return a capitalized version of S, i.e. make the first character  
have upper case and the rest lower case.
```

```
Type:      builtin_function_or_method
```

```
[In [10]: h.capitalize()  
Out[10]: 'Hello there!']
```

IPYTHON: KEYBOARD SHORTCUTS

Keystroke	Action
Ctrl-a	Move cursor to the beginning of the line
Ctrl-e	Move cursor to the end of the line
Ctrl-b (or the left arrow key)	Move cursor back one character
Ctrl-f (or the right arrow key)	Move cursor forward one character

Keystroke	Action
Backspace key	Delete previous character in line
Ctrl-d	Delete next character in line
Ctrl-k	Cut text from cursor to end of line
Ctrl-u	Cut text from beginning of line to cursor
Ctrl-y	Yank (i.e., paste) text that was previously cut
Ctrl-t	Transpose (i.e., switch) previous two characters



THANKS JAKE!

IPYTHON: REFERENCING PAST IN/OUT

```
[In [1]: 2 + 2
Out[1]: 4

[In [2]: Out
Out[2]: {1: 4}

[In [3]: Out[1]
Out[3]: 4

[In [4]: _
Out[4]: 4

[In [5]: _ + 3
Out[5]: 7

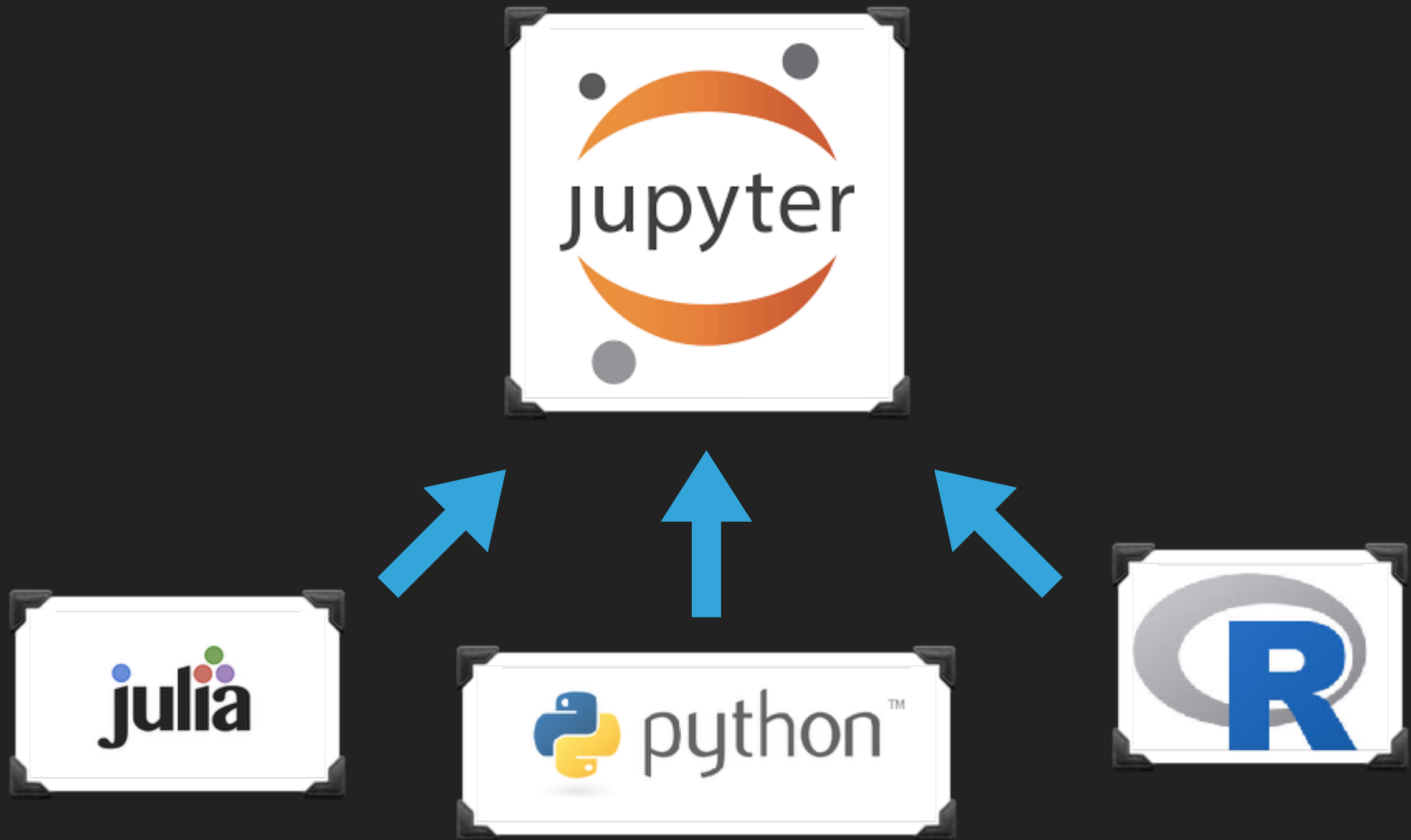
[In [6]: __ + _
Out[6]: 11

In [7]: █
```


IPYTHON: MAGIC STUFF!

- ▶ Magic Stuff!
 - ▶ `%paste`
 - ▶ `%timeit`
 - ▶ `%magic`
- ▶ See: <https://jakevdp.github.io/PythonDataScienceHandbook/01.03-magic-commands.html>

JUPYTER



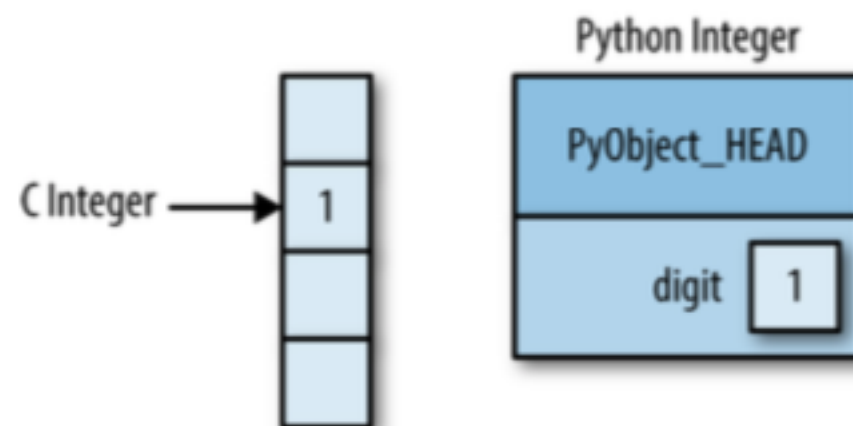
JUPYTER

1. Basic navigation: `enter`, `shift-enter`, `up/k`, `down/j`
2. Saving the notebook: `s`
3. Change Cell types: `y`, `m`, `1-6`, `t`
4. Cell creation: `a`, `b`
5. Cell editing: `x`, `c`, `v`, `d`, `z`
6. Kernel operations: `i`, `0` (press twice)

III. DATA WRANGLING / ANALYSIS

NUMPY / PANDAS

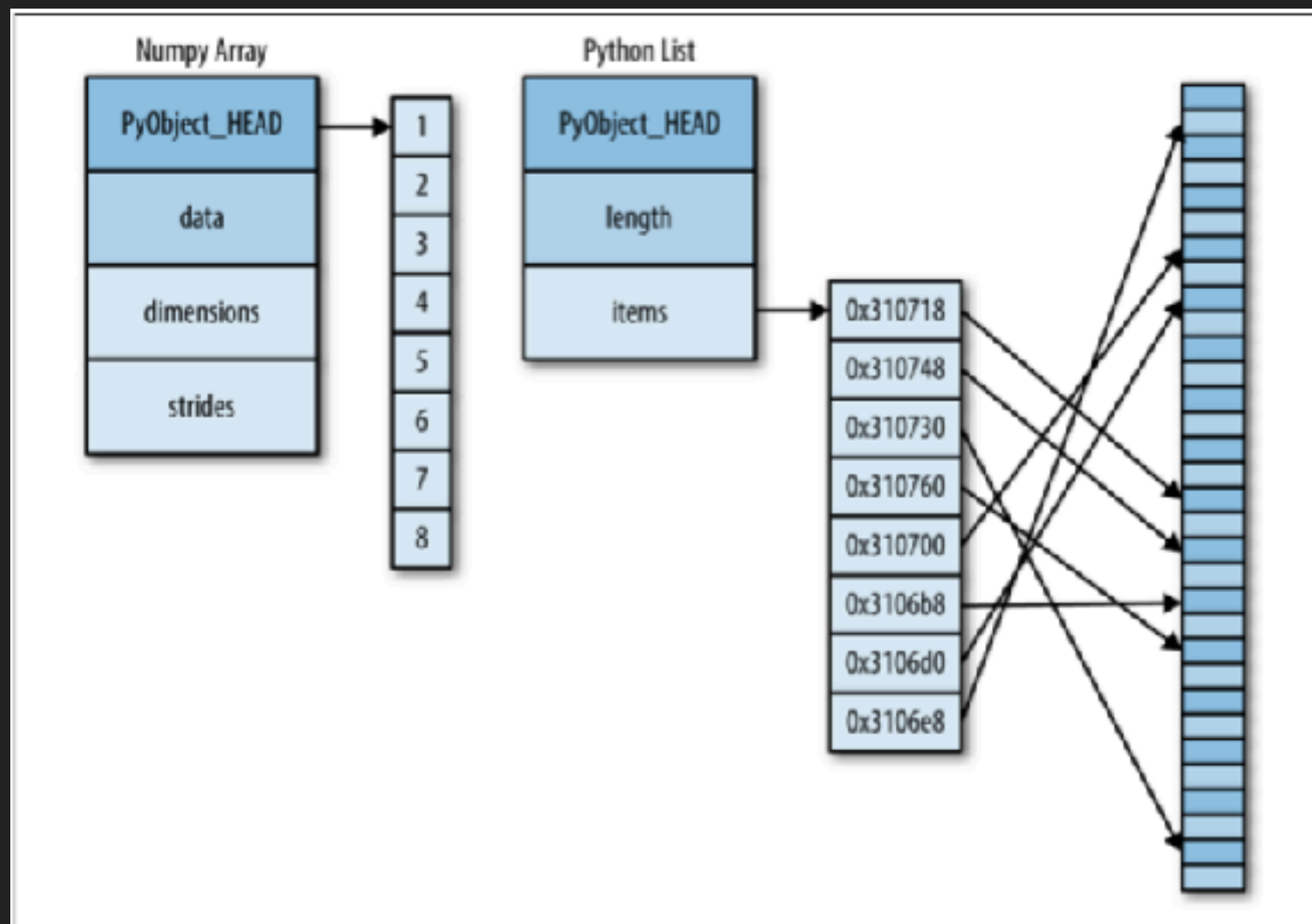
NUMPY



THANKS JAKE!

III. DATA WRANGLING / ANALYSIS

NUMPY



THANKS JAKE!

NUMPY: DATA TYPES

Table 2-1. Standard NumPy data types

Data type	Description
<code>bool_</code>	Boolean (True or False) stored as a byte
<code>int_</code>	Default integer type (same as C <code>long</code> ; normally either <code>int64</code> or <code>int32</code>)
<code>intc</code>	Identical to C <code>int</code> (normally <code>int32</code> or <code>int64</code>)
<code>intp</code>	Integer used for indexing (same as C <code>ssize_t</code> ; normally either <code>int32</code> or <code>int64</code>)
<code>int8</code>	Byte (−128 to 127)
<code>int16</code>	Integer (−32768 to 32767)
<code>int32</code>	Integer (−2147483648 to 2147483647)
<code>int64</code>	Integer (−9223372036854775808 to 9223372036854775807)
<code>uint8</code>	Unsigned integer (0 to 255)
<code>uint16</code>	Unsigned integer (0 to 65535)
<code>uint32</code>	Unsigned integer (0 to 4294967295)
<code>uint64</code>	Unsigned integer (0 to 18446744073709551615)
<code>float_</code>	Shorthand for <code>float64</code>
<code>float16</code>	Half-precision float: sign bit, 5 bits exponent, 10 bits mantissa
<code>float32</code>	Single-precision float: sign bit, 8 bits exponent, 23 bits mantissa
<code>float64</code>	Double-precision float: sign bit, 11 bits exponent, 52 bits mantissa
<code>complex_</code>	Shorthand for <code>complex128</code>
<code>complex64</code>	Complex number, represented by two 32-bit floats
<code>complex128</code>	Complex number, represented by two 64-bit floats

THANKS JAKE!

NUMPY: UFUNCS

Table 2-2. Arithmetic operators implemented in NumPy

Operator	Equivalent ufunc	Description
+	<code>np.add</code>	Addition (e.g., $1 + 1 = 2$)
-	<code>np.subtract</code>	Subtraction (e.g., $3 - 2 = 1$)
-	<code>np.negative</code>	Unary negation (e.g., -2)
*	<code>np.multiply</code>	Multiplication (e.g., $2 * 3 = 6$)
/	<code>np.divide</code>	Division (e.g., $3 / 2 = 1.5$)
//	<code>np.floor_divide</code>	Floor division (e.g., $3 // 2 = 1$)
**	<code>np.power</code>	Exponentiation (e.g., $2 ** 3 = 8$)
%	<code>np.mod</code>	Modulus/remainder (e.g., $9 \% 4 = 1$)



THANKS JAKE!

NUMPY: AGGREGATION FUNCTIONS

Table 2-3. Aggregation functions available in NumPy

Function Name	NaN-safe Version	Description
<code>np.sum</code>	<code>np.nansum</code>	Compute sum of elements
<code>np.prod</code>	<code>np.nanprod</code>	Compute product of elements
<code>np.mean</code>	<code>np.nanmean</code>	Compute mean of elements
<code>np.std</code>	<code>np.nanstd</code>	Compute standard deviation
<code>np.var</code>	<code>np.nanvar</code>	Compute variance
<code>np.min</code>	<code>np.nanmin</code>	Find minimum value
<code>np.max</code>	<code>np.nanmax</code>	Find maximum value
<code>np.argmin</code>	<code>np.nanargmin</code>	Find index of minimum value
<code>np.argmax</code>	<code>np.nanargmax</code>	Find index of maximum value
<code>np.median</code>	<code>np.nanmedian</code>	Compute median of elements
<code>np.percentile</code>	<code>np.nanpercentile</code>	Compute rank-based statistics of elements
<code>np.any</code>	N/A	Evaluate whether any elements are true
<code>np.all</code>	N/A	Evaluate whether all elements are true

THANKS JAKE!

IV. VISUALIZATION

MATPLOTLIB / SEABORN