# Introduction to DS Tools

grossberg@cs.ccny.cuny.edu

# What can we use to do Data Science?

grossberg@cs.ccny.cuny.edu

# Excel?

# Iris Data Set
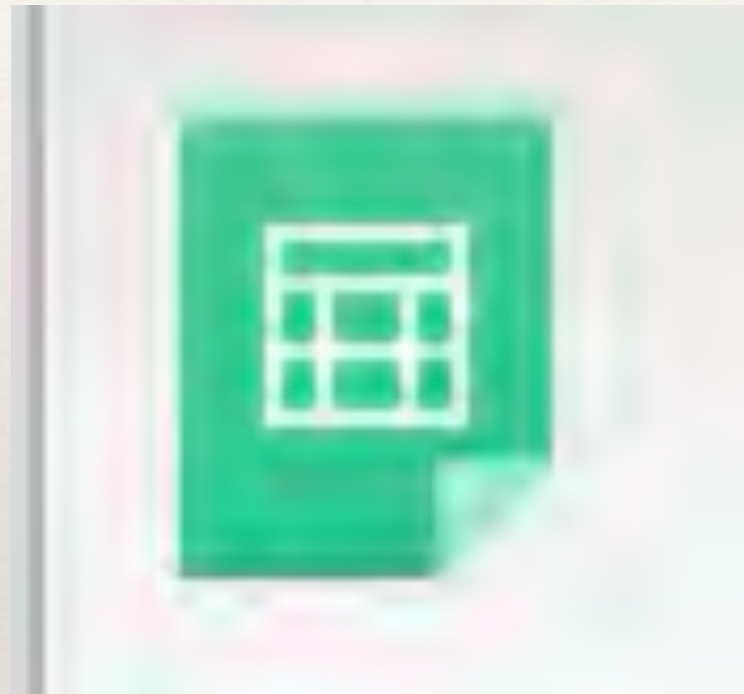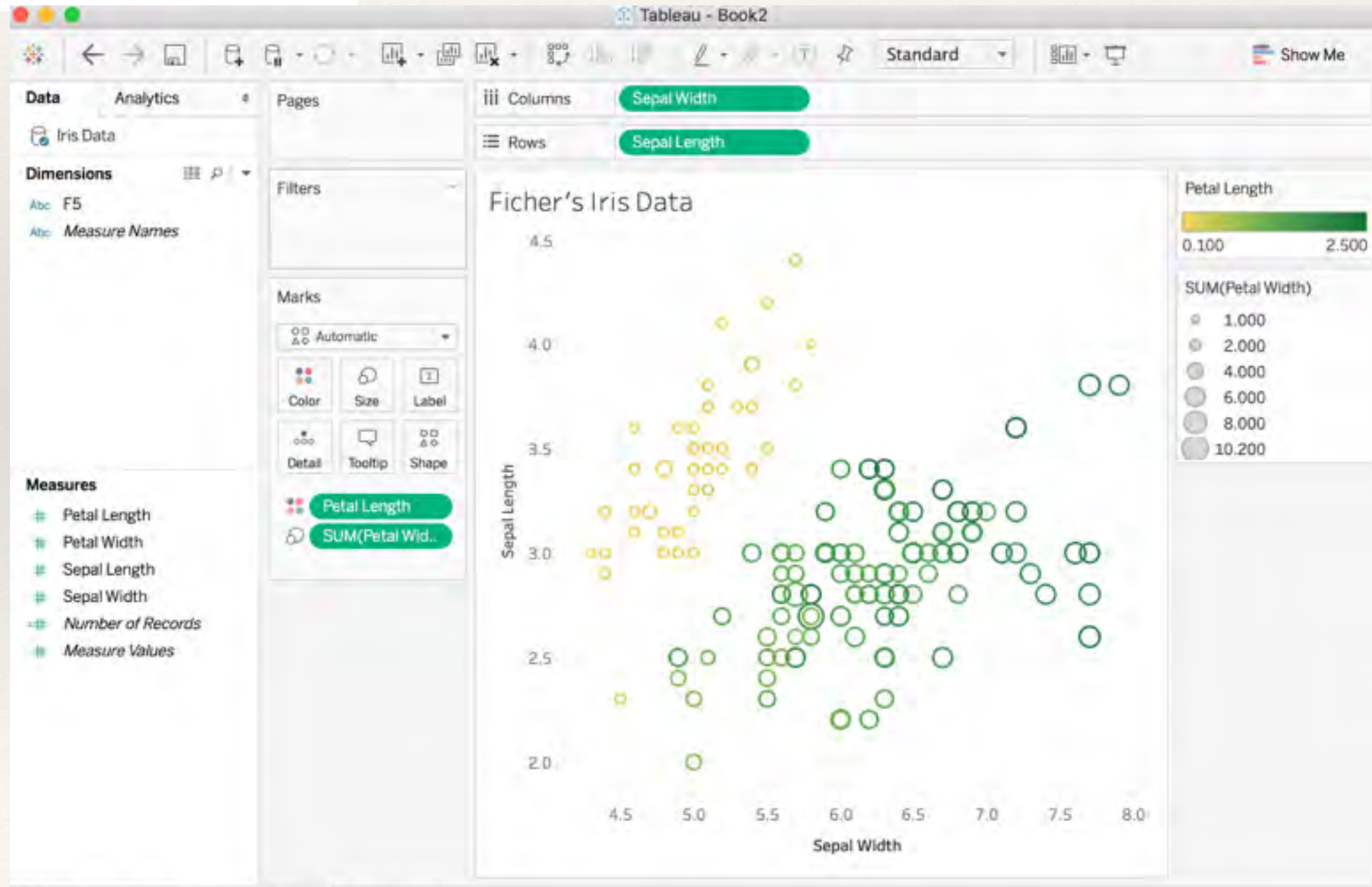
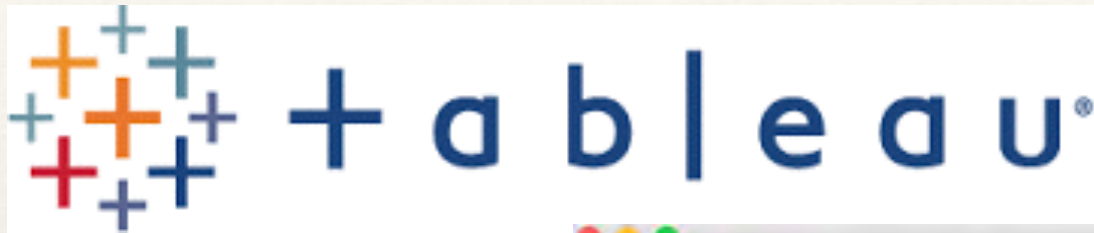# Basic Scatter Plots and Analysis
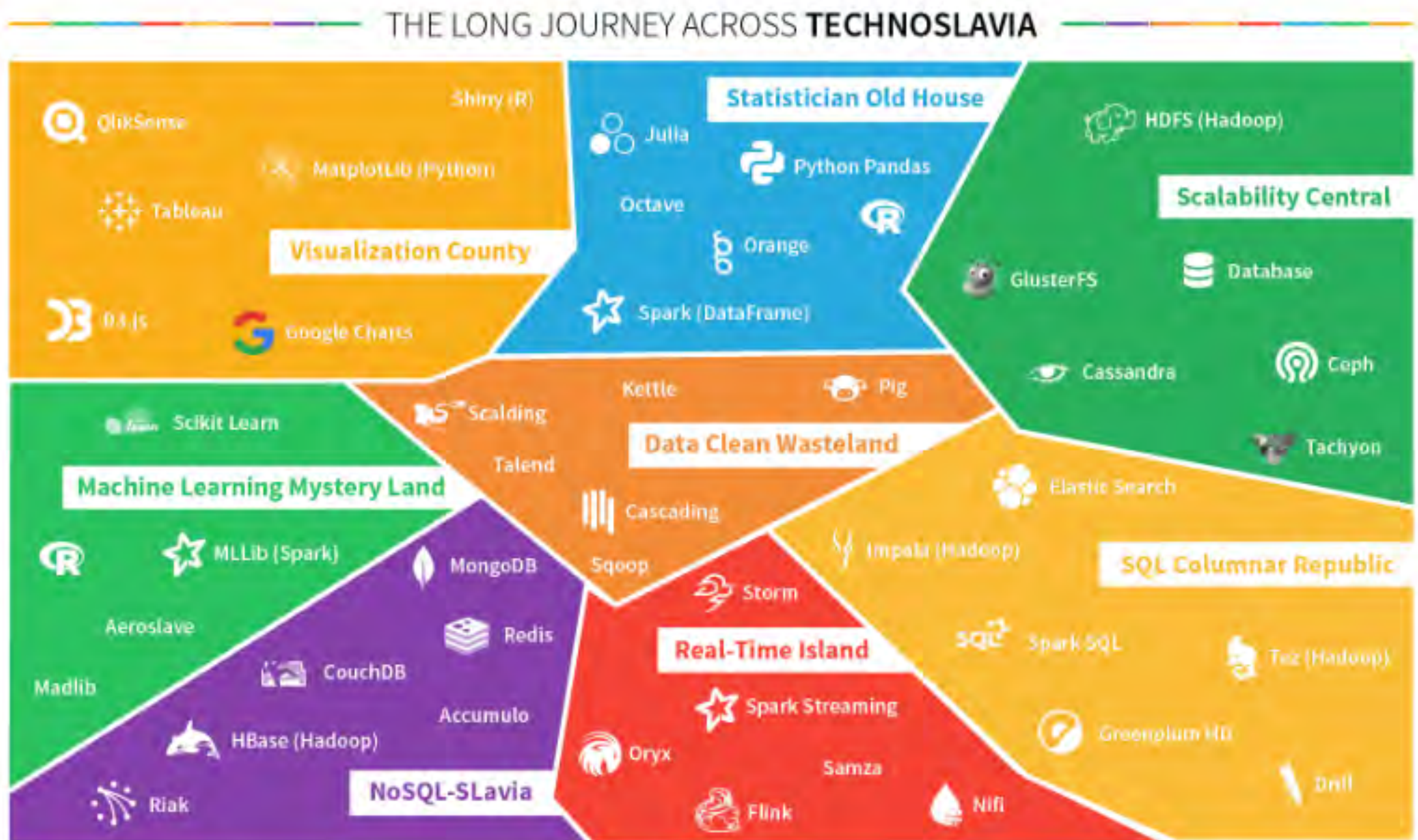
# Quick look at "small data"



Google Sheet  LibreOffice Calc  Apple Numbers

# Tableau

# Languages/Tools for Data Science



THE LONG JOURNEY ACROSS **TECHNOSLAVIA**

Shiny (R)

QlikSense

MatplotLib (Python)

Tableau

**Visualization County**

D3.js

Google Charts

**Statistician Old House**

Julia

Python Pandas

Octave

R

Orange

Spark (DataFrame)

HDFS (Hadoop)

**Scalability Central**

GlusterFS

Database

Cassandra

Ceph

Scalding

Kettle

Pig

Talend

**Data Clean Wasteland**

Cascading

Scikit Learn

**Machine Learning Mystery Land**

MLLib (Spark)

Sqoop

Impala (Hadoop)

Tachyon

Elastic Search

**SQL Columnar Republic**

Storm

Spark SQL

Tez (Hadoop)

Aeroslave

MongoDB

Redis

**Real-Time Island**

Madlib

CouchDB

Accumulo

Spark Streaming

Greenplum HD

HBase (Hadoop)

Oryx

Samza

Drill

Riak

**NoSQL-SLavia**

Flink

Nifi

# Programatic Analysis

# Python vs R: Data Science WAR

# R vs Python

## History

### R

**Creators**

Ross Ihaka and Robert Gentleman

**Release Year**

1995

**Must Knows**

1. R is an implementation of S programming language (Bell Labs).

2. R's design and evolution is handled by the R-core group and R foundation.

3. R's software environment was written primarily in C, Fortran and R.

### Python

**Creator**

Guido Van Rossum

**Release Year**

1991

**Must Knows**

1. Python was inspired by C, Modula-3, and particularly ABC.

2. Python gets its name from the "Monty Python's Flying Circus" comedy series.

3. Python Software Foundation (PSF) takes care of Python's advances.

# R vs Python

## Purpose

R focuses on better, user friendly data analysis, statistics and graphical models.

Python emphasizes productivity and code readability.

## Used By?

R has been used primarily in academics and research. However, R is rapidly expanding into the enterprise market.

*"The closer you are to statistics, research and data science, the more you might prefer R."*

Python is used by programmers that want to delve into data analysis or apply statistical techniques, and by developers that turn to data science.

*"The closer you are to working in an engineering environment, the more you might prefer Python."*

grossberg@cs.ccny.cuny.edu

# R vs Python

## Usability

| R | Python |
|---|---|
| Statistical models can be written with only a few lines. | Coding and debugging is easier to do in Python, mainly because of the "nice" syntax. |
| There are R stylesheets but not everyone uses them. | The indentation of the code affects its meaning. |
| The same piece of functionality can be written in several ways in R. | Any piece of functionality is always written the same way in Python. |

## Flexibility

| R | Python |
|---|---|
| It is easy to use complex formulas in R. All kinds of statistical tests and models are readily available and easily used. | Python is flexible for doing something novel that has never been done before. Developers can also use it for scripting a website or other applications. |

grossberg@cs.ccny.cuny.edu

## Ease of Learning

R has a steep learning curve at start. Once you know the basics, you can easily learn advanced stuff.

R is not hard for experienced programmers.

*Check out DataCamp's interactive exercises and tutorials.*

Python's focus on readability and simplicity makes that its learning curve is relatively low and gradual.

Python is considered a good language for starting programmers.

*Try using the book "Learn Python The Hard Way" and its accompanying site with videos and exercises.*

## Code Repositories

CRAN stands for the Comprehensive R Archive Network: it is a huge repository of R packages to which users can easily contribute.

Packages are collections of R functions, data, and compiled code. They can be installed in R with one line.

PyPi is the Python Package Index: it is a repository of Python software, consisting of libraries. Users can contribute to Pypi, but it is a bit complicated in practice.

Watch out with dependencies and installing Python libraries!

*"I don't see Python [...] building up a huge code repository comparable to CRAN. [R has] a gigantic head start, [and] [...] statistics simply is not Python's central mission;"*
*- Norm Matloff, professor of computer science*

14

# R vs. Python



**Usage**

R is mainly used when the data analysis tasks require standalone computing or analysis on individual servers.

Python is generally used when the data analysis tasks need to be integrated with web apps or if statistics code needs to be incorporated into a production database.
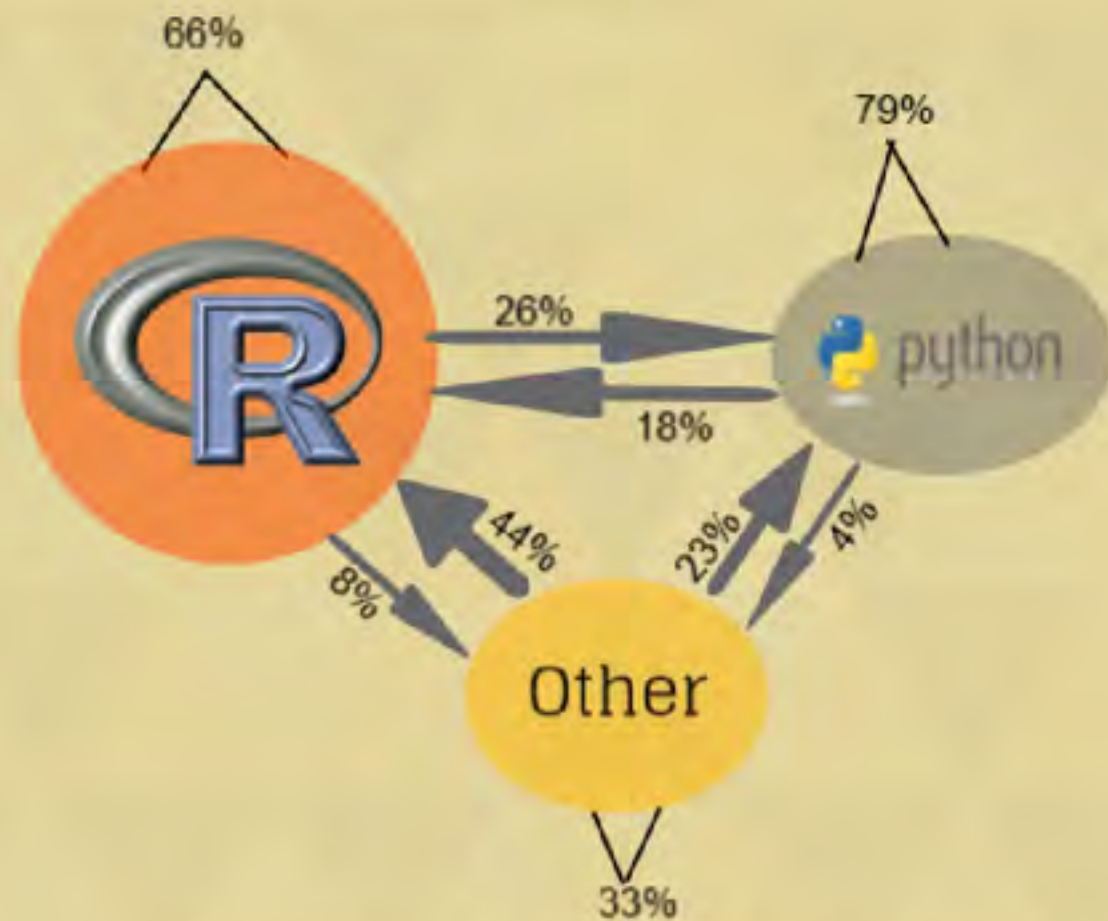
**Task**

For exploratory work, R is easier for beginners. Statistical models can be written with a few lines of code.

As a full-fledged programming language, Python is a good tool to implement algorithms for production use.

grossberg@cs.ccny.cuny.edu

# Both Winners



Switching Between R and Python?

Number of people switching between R and Python in 2013 *

66%

79%

26%

18%

44%

23%

4%

8%

Other

33%

*Percentages on the arrows are relative to the base

"My current strategy is to leverage the best of both worlds — do early stage data analysis in R, then switch to Python when it's time to get serious, be a team player, and ship some real code and data products."

• • •

"I use R to conduct statistical tests, graph data, and inspect large data sets. If I actually have to write an algorithm, I prefer Python..."

• • •

"I'd rather do math in a general-purpose language than try to do general-purpose programming in a math language."

# Jobs need both



Jobs and Salary?

O'Reilly 2014 Data Science Salary Survey

Average Annual Salaries In The Range Of:

Python = US$110, 000 to US$125,000 = R

R and Python job trends

Job Trends from Indeed.com

— R !"R D" !"A R" !"H R" !"R N"  !toys !kids !" R Walgreen" !walmart !"HVAC R" !"R Bard"  and (
— Python and ("big data" or "statistical analysis" or "data mining" or "data analytics" or "mach

# R: IDE R-Studio

# Some R resources

# Python

Don't Use  Python IDLE (it is not good)

# Unix Environment



Terminal

Shell (Bash)

21

# Some Must-Know Bash Commands

- ls -l -a

- find . -name "*.py"

- | [aka pipe]

- whereis

- cp

- rm

- mv

- cat

- touch

- mkdir

- > and >>

- less

- chmod

- vi [vim]

22

# Totally Useless Bash Commands

* Cowsay

* locate

* fortune

# Python Shell (in Bash in Terminal)



```
1. python3.5
mamilla:~ michael$ python
Python 3.5.2 |Anaconda 4.1.1 (x86_64)| (default, Jul  2 2016, 17:52:12)
[GCC 4.2.1 Compatible Apple LLVM 4.2 (clang-425.0.28)] on darwin
Type "help", "copyright", "credits" or "license" for more information.
>>> print("Mary had a little lamb.")
Mary had a little lamb.
>>>
```

Not so useful. No history.

# Create Edit Python File in vim (atom)



```
mamilla:~ michael$ vim foo.py
```



```
1 print( "Hello Cruel World!")
~
~
~
~
~
~
```

INSERT    [No Name] +    <    |    utf-8    |    no ft    100%    LN    1:29
-- INSERT --

# run using "python"

```
1. bash

mamilla:~ michael$ python foo.py
Hello Cruel World!
mamilla:~ michael$ 
```

26

# Better Shell Ipython

```
mamilla:~ michael$ ipython
Python 3.5.2 |Anaconda 4.1.1 (x86_64)| (default, Jul  2 2016, 17:52:12)
Type "copyright", "credits" or "license" for more information.

IPython 4.2.0 -- An enhanced Interactive Python.
?         -> Introduction and overview of IPython's features.
%quickref -> Quick reference.
help      -> Python's own help system.
object?   -> Details about 'object', use 'object??' for extra details.

In [1]:
```
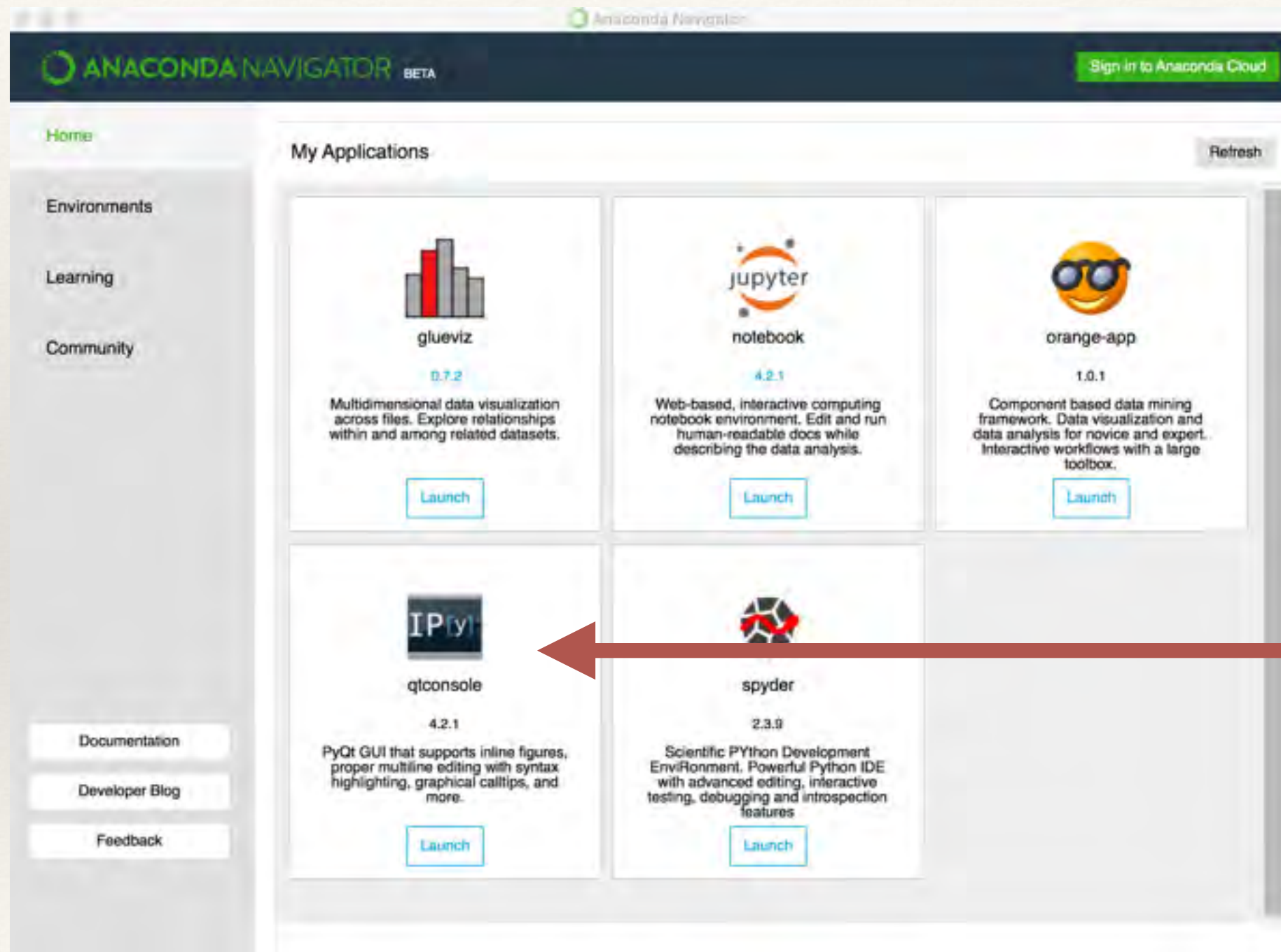
Lots of features (including history)
Launch from terminal OR launch from anaconda

27

# Ipython via Anaconda



Simple Ipython Shell

# Basic Ipython

Access to shell:
ls
pwd
!ls
History:
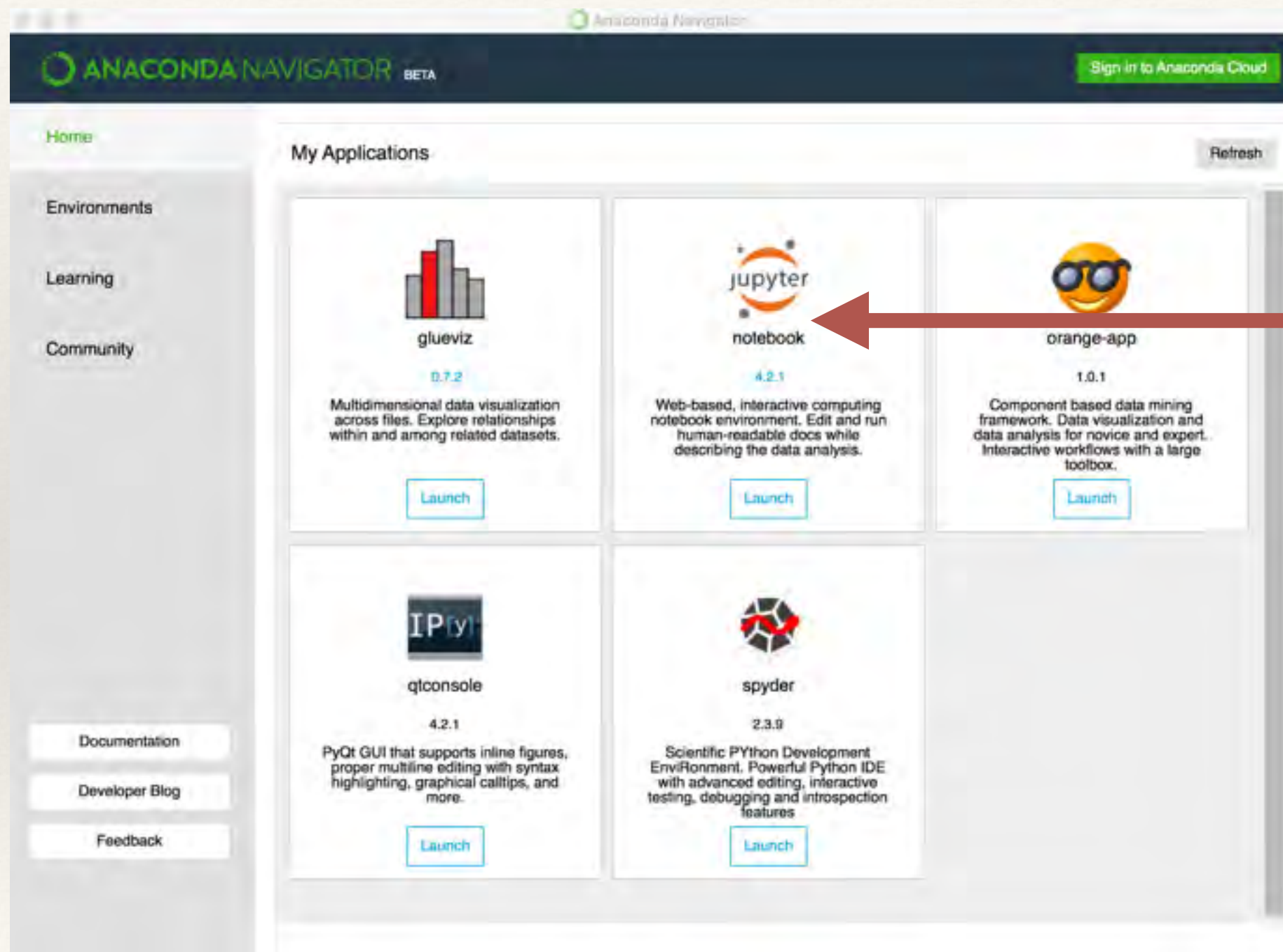[1] each command
control-p / -n

Help:
print? [gets help]
Magic:
%lsmagic
%automagic

29

# Ipython

❖ Ipython a Python Shell (we'll talk about Jupyter later)
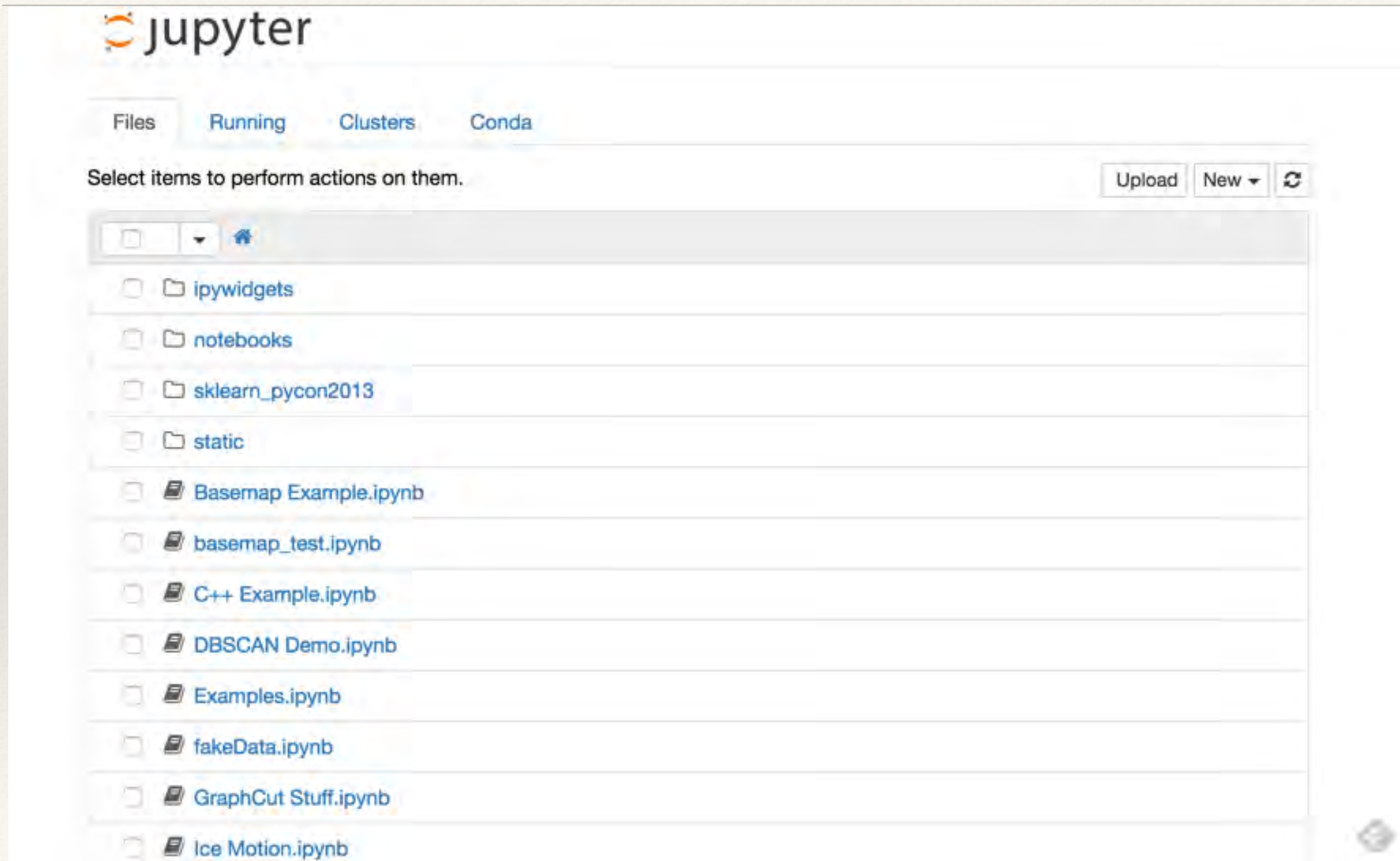
# Jupyter Notebook: from anaconda



Run Notebook Server

# Jupyter Notebook: from shell

```
mamilla:ipythonnb michael$ jupyter notebook
[W 12:43:33.949 NotebookApp] Unrecognized JSON config file version, assuming versio
n 1
[I 12:43:34.595 NotebookApp] [nb_conda_kernels] enabled, 2 kernels found
[I 12:43:35.193 NotebookApp] ✓ nbpresent HTML export ENABLED
[W 12:43:35.193 NotebookApp] ✗ nbpresent PDF export DISABLED: No module named 'nbbr
owserpdf'
[I 12:43:35.236 NotebookApp] [nb_anacondacloud] enabled
[I 12:43:35.240 NotebookApp] [nb_conda] enabled
[I 12:43:35.247 NotebookApp] Serving notebooks from local directory: /Users/michael
/ipythonnb
[I 12:43:35.247 NotebookApp] 0 active kernels
[I 12:43:35.247 NotebookApp] The Jupyter Notebook is running at: http://localhost:8
888/
[I 12:43:35.247 NotebookApp] Use Control-C to stop this server and shut down all ke
rnels (twice to skip confirmation).
[W 12:43:36.714 NotebookApp] /Users/michael/ipythonnb/DV59969S16 doesn't exist
[W 12:43:36.715 NotebookApp] /Users/michael/ipythonnb/DV83060S16 doesn't exist
```

# Jupyter Notebooks: Directory

grossberg@cs.ccny.cuny.edu

# Jupyter Notebooks: Launch



Start
Ipython
Notebook

# Jupyter notebook with ipython kernel



Perfect for exploration experiments and tutorials

# Basic DS Python Stack (know well)

grossberg@cs.ccny.cuny.edu

# Data Science From Scratch

❖ Chapter 2: Python Crash Course