

# Python for Data Analysis and Scientific Computing

X433.3 (2 semester units in COMPSCI)

Instructor Alexander I. Iliev, Ph.D.

# Course Content Outline

- **Python for Data Analysis and Scientific Computing®**

- Python - pros and cons
- Installing the environment with core packages
- Python modules, packages and scientific blocks
- Working with the shell, IPython and the editor

HW1 (5pts)

- **Basic language specifics 1/2**

- Basic arithmetic operations, assignment operators, data types, containers
- Iterative programming (if/elif/else)
- Conditional expressions
- Recursion programming (for/continue/while/break)
- Functions: definition, return values, local vs. global variables

- **Basic language specifics 2/2**

- Classes / Functions (cont.): objects, methods, passing by value and reference
- Scripts, modules, packages
- I/O interaction with files
- Standard library
- Exceptions

- **NumPy 1/3**

- Why NumPy?
- Data type objects
- NumPy arrays
- Indexing and slicing of arrays

HW2 (5pts)

- **Matplotlib**

- What is Matplotlib?
- Basic plotting
- Tools: title, labels, legend, axis, points, subplots, etc.
- Advanced plotting: scatter, pie, bar, 3D plots, etc.

HW3 (5pts)

# Course Content Outline

- **NumPy 2/3**
  - Array operations
  - Reductions
  - Broadcasting
  - Array: shaping, reshaping, flattening, resizing, dimension changing
  - Data sorting
- **Midterm / Project proposal due (30pts)**
- **NumPy 3/3**
  - Type casting
  - Masking data
  - Organizing arrays
  - Loading data files
  - Dealing with polynomials
  - Good coding practices
- **Scipy 1/2**
  - What is Scipy?
  - Working with files
  - Algebraic operations
  - The Fast Fourier Transform
  - Signal Processing
- **HW4 (5pts)**
- **Scipy 2/2**
  - Interpolation
  - Statistics
  - Optimization
  - Multithreading and Multiprocessing
- **Project**
  - Project Presentation
- **Final Project (40pts)**

# Course Content Outline

- **Methods of Instruction**

Lectures and in-class discussions will be the main tools of instruction. Homework problems will help students absorb the material and get to practice in their free time. Since this is a more specialized course that naturally delves into specific topics a midterm exam will be held to assess the performance of the students midway through the course.

- **Credit Requirements**

Students must complete all homework assignments and pass all exams. They also must receive a passing grade on the final exam to receive a passing grade in the course.

- **Course Grade Weighting**

1. In-class participation: 10%
2. Homework: 20% - Wednesdays - use the weekend to prepare
3. Midterm 30% - Monday - use the weekend to prepare
4. Final Project 40% - Wednesday

- **Optional Reading for the Course**

Title: Python for Data Analysis

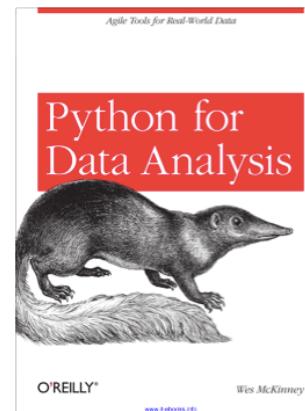
ISBN-10: 1449319793 and ISBN-13: 978-1449319793

Author(s): Wes McKinney

Publisher: O'Reilly Media

Edition Number: 1 edition

Publication Date: November 1, 2012



# Course Content Outline

- **Python for Data Analysis and Scientific Computing®**

- Python - pros and cons
- Installing the environment with core packages
- Python modules, packages and scientific blocks
- Working with the shell, IPython and the editor

HW1

- **Basic language specifics 1/2**

- Basic arithmetic operations, assignment operators, data types, containers
- Iterative programming (if/elif/else)
- Conditional expressions
- Recursion programming (for/continue/while/break)
- Functions: definition, return values, local vs. global variables

- **Basic language specifics 2/2**

- Functions (cont.): objects, methods, passing by value and reference
- Scripts, modules, packages
- I/O interaction with files
- Standard library
- Exceptions

- **NumPy 1/3**

- Why NumPy?
- Data type objects
- NumPy arrays
- Indexing and slicing of arrays

- **Matplotlib**

- What is Matplotlib?
- Basic plotting
- Tools: title, labels, legend, axis, points, subplots, etc.
- Advanced plotting: scatter, pie, bar, 3D plots, etc.

# Data

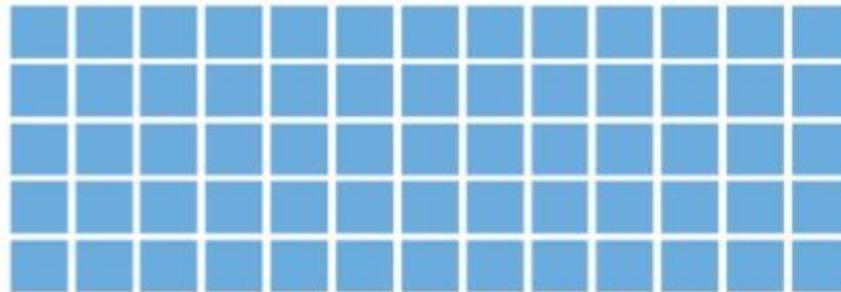
- The world's digital content and media is **growing rapidly** at a never stopping rate



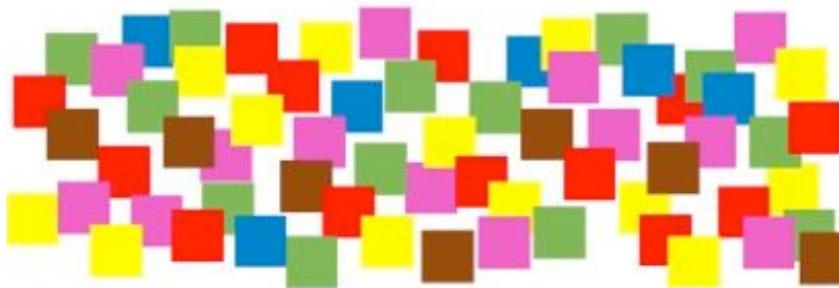
*90% of the data in the world today has been created in the last two years alone, at 2.5 quintillion bytes of data a day*

*report from IBM Marketing Cloud in 2017*

# Data



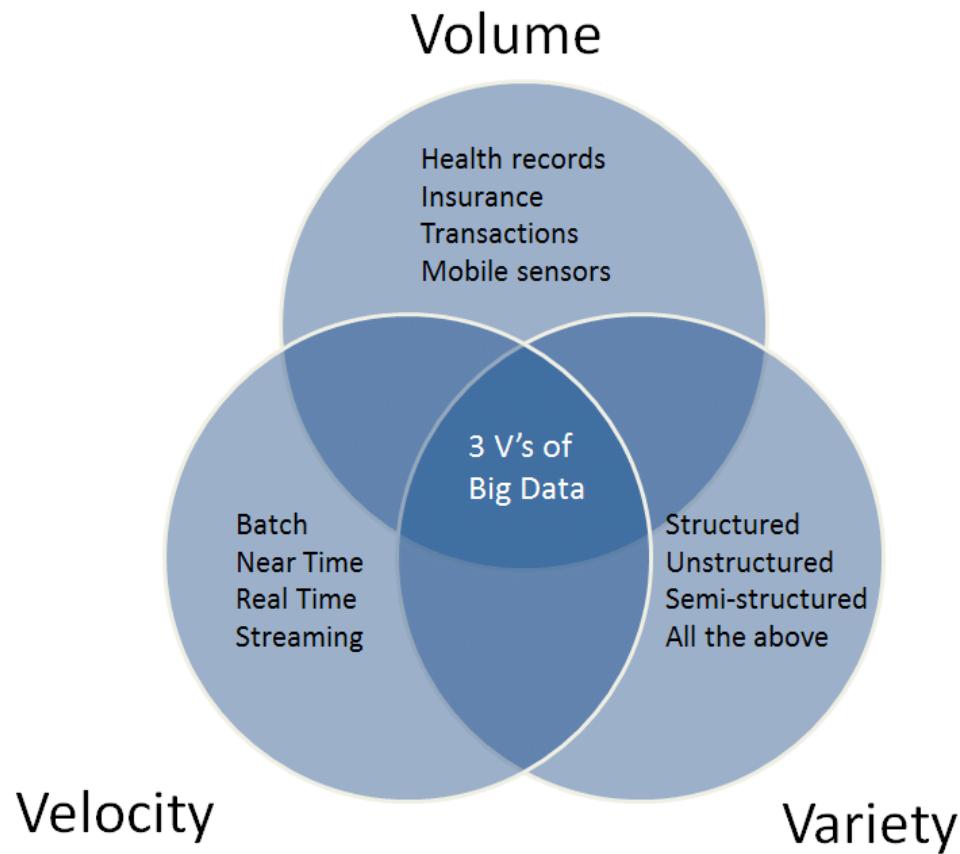
**Structured Data**



**Unstructured Data**

Graphical representations illustrate the difference between  
structured and unstructured data

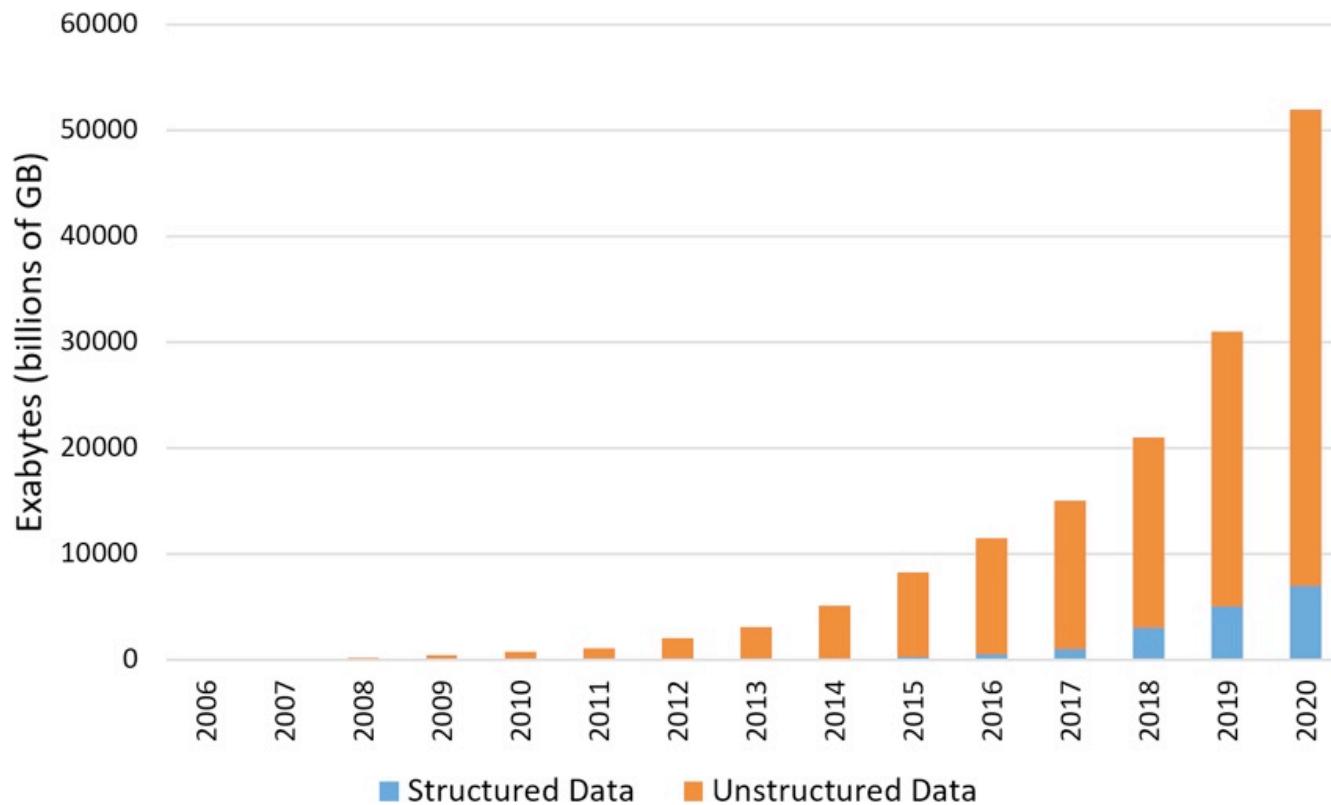
# Big Data



Variability, Veracity, and Value

# Data

## The Cambrian Explosion...of Data



# Data

Multiplying Factor	SI Prefix	Scientific Notation	Name
1 000 000 000 000 000 000 000 000	Yotta (Y)	$10^{24}$	1 septillion
1 000 000 000 000 000 000 000	Zetta (Z)	$10^{21}$	1 sextillion
1 000 000 000 000 000 000	Exa (E)	$10^{18}$	1 quintillion
1 000 000 000 000 000	Peta (P)	$10^{15}$	1 quadrillion
1 000 000 000 000	Tera (T)	$10^{12}$	1 trillion
1 000 000 000	Giga (G)	$10^9$	1 billion
1 000 000	Mega (M)	$10^6$	1 million
1 000	kilo (k)	$10^3$	1 thousand
0 001	milli (m)	$10^{-3}$	1 thousandth
0 000 001	micro (u)	$10^{-6}$	1 millionth
0 000 000 001	nano (n)	$10^{-9}$	1 billionth
0 000 000 000 001	pico (p)	$10^{-12}$	1 trillionth
0 000 000 000 000 001	femto (f)	$10^{-15}$	1 quadrillionth
0 000 000 000 000 000 001	atto (a)	$10^{-18}$	1 quintillionth
0 000 000 000 000 000 000 001	zepto (z)	$10^{-21}$	1 sextillionth
0 000 000 000 000 000 000 000 001	yocto (y)	$10^{-24}$	1 septillionth

Metric prefixes defined at the 19th General Conference on Weights and Measures in 1991

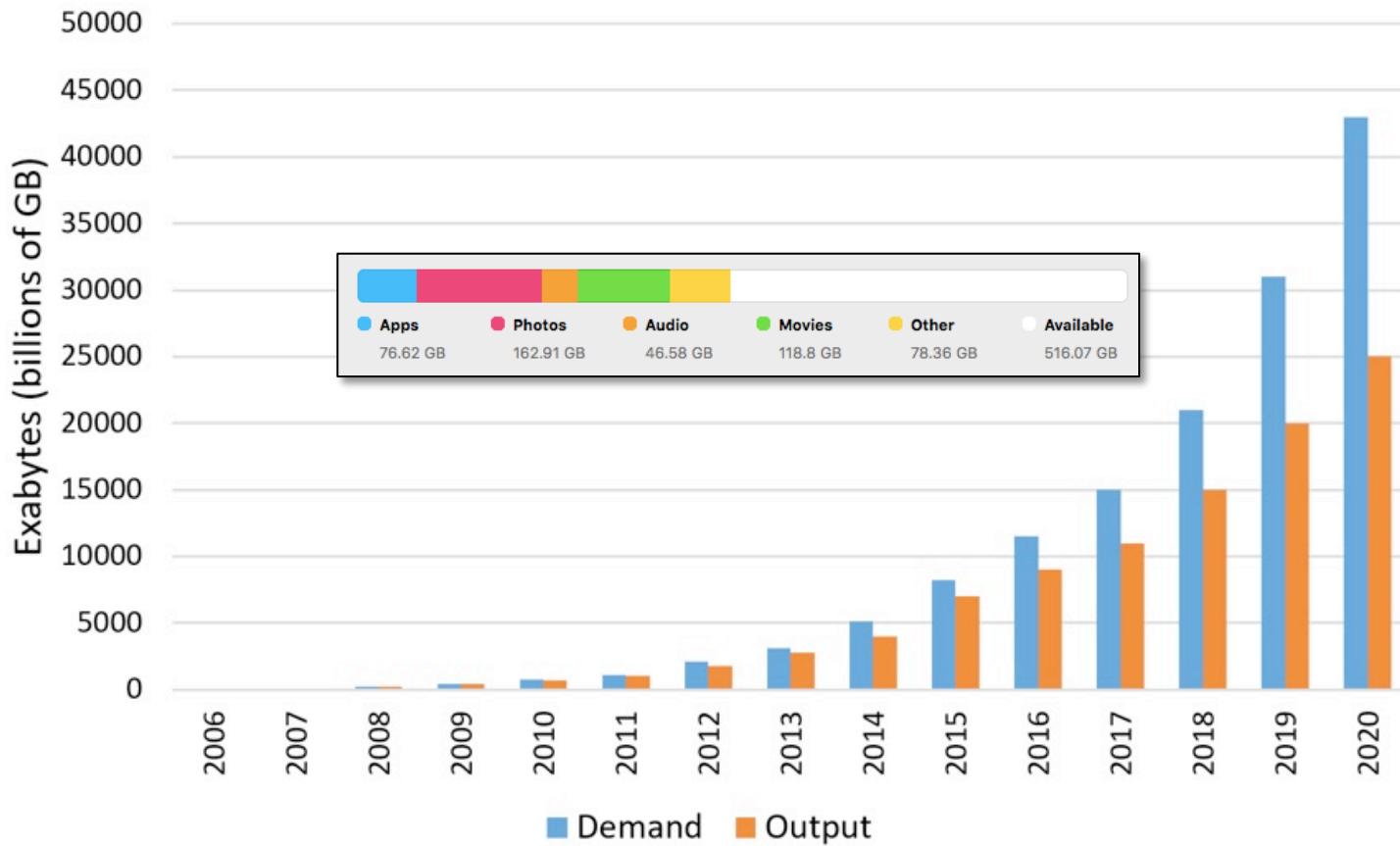
# Data

Multiplying Factor	SI Prefix	Scientific Notation	Name
1 208 925 819 614 629 174 706 176	Yottabytes	$2^{80}$	1 septillion
1 180 591 620 717 411 303 424	Zettabytes	$2^{70}$	1 sextillion
1 152 921 504 606 846 976	Exabytes	$2^{60}$	1 quintillion
1 125 899 906 842 624	Petabytes	$2^{50}$	1 quadrillion
1 099 511 627 776	Terabytes	$2^{40}$	1 trillion
1 073 741 824	Gigabytes	$2^{30}$	1 billion
1 048 576	Megabytes	$2^{20}$	1 million
1 024	kilobytes	$2^{10}$	1 thousand

Examples of prefixes used to measure digital data with a binary system

# Data

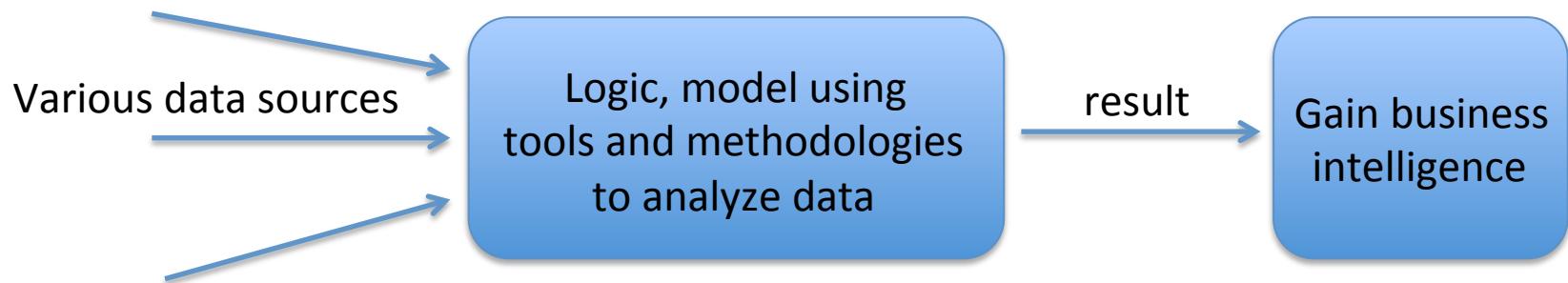
## Storage Supply & Demand



Storage supply and demand growth over two decades

# Why using Data Analytics?

- To solve complex scientific or business tasks
- Using complex data from various sources
- Using tools and methods to extract meaningful scientific or business insights



# Data Analytics vs Data Analysis

- **Data analytics** is a broader field, which includes data analysis as part of its core
- **Data analytics** defines the methodology used behind the analysis
- **Methodology** means to comprehend the cognitive processes that are used in order to understand any given problem and explore data in logical fashion
- **Data analytics** also includes:
  - data creation, **extraction**, and **collection**;
  - **transformation** using specific tool sets and techniques;
  - different **methods for analysis**;
  - communication of results.

# The Problem



# The Problem

- Example:

Facebook



- Over **2 billion** active users worldwide generate enormous amounts of user data
- Every **60 seconds** there are:
  - **136,000 photos** are uploaded
  - **510,000 comments** are posted
  - **283,000 status updates** are made
- lots of unstructured data (audio, photos, videos)

# Data Analysis

- What is it?

Data analysis, is a process of inspecting, cleansing, transforming, and modeling data with the goal of discovering useful information, suggesting conclusions, and supporting decision-making.

source: Wikipedia

# Data Analysis

- How does it work?
  - First user **data is collected**
  - Then **user behavior** is studied
  - Data is packaged for **advertisers**
  - Enhance useful functionality (ex. shopping) and **gather even more data from users** (habits, preferences, etc.)

# Why using Data Analytics?

- Data can be very noisy and hard to use

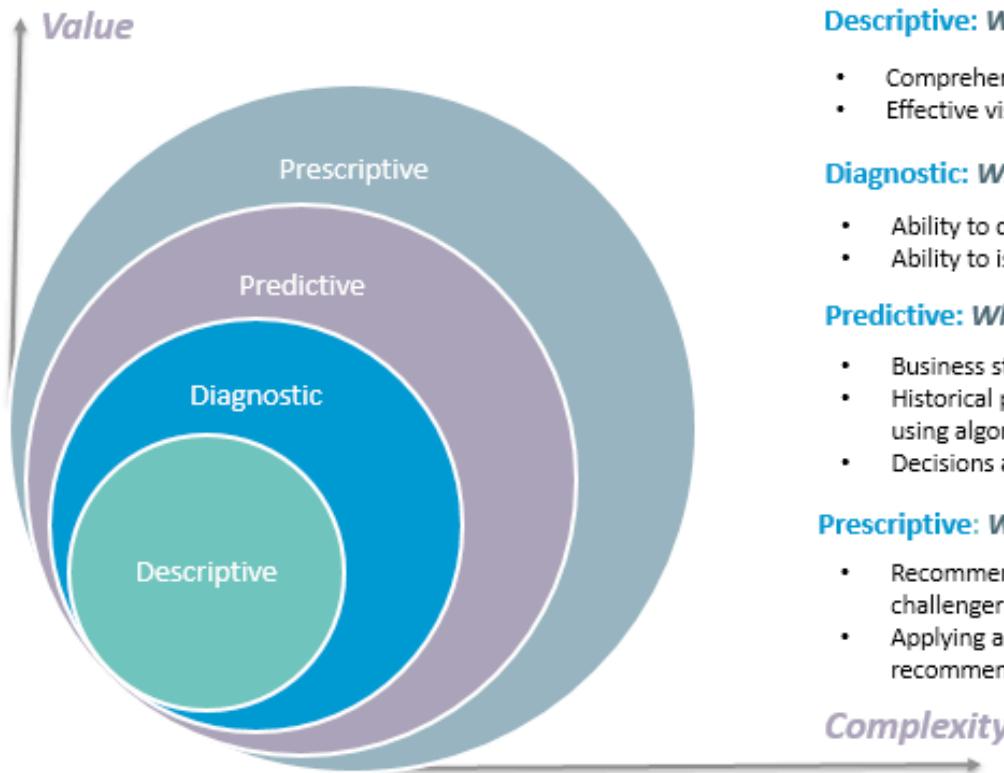


# Types of Data Analytics

- Descriptive data analytics studies data based on previous samples  
(what happened?)
- Diagnostic data analytics studies the reasons and isolates the problem  
(why it happened?)
- Predictive data analytics predicts the possible future outcome using statistical modeling  
(what might happen?)
- Prescriptive data analytics makes logical decisions to suggest best course of action  
(what to do?)

# Types of Data Analytics

## 4 types of Data Analytics



### What is the data telling you?

#### **Descriptive:** *What's happening in my business?*

- Comprehensive, accurate and live data
- Effective visualisation

#### **Diagnostic:** *Why is it happening?*

- Ability to drill down to the root-cause
- Ability to isolate all confounding information

#### **Predictive:** *What's likely to happen?*

- Business strategies have remained fairly consistent over time
- Historical patterns being used to predict specific outcomes using algorithms
- Decisions are automated using algorithms and technology

#### **Prescriptive:** *What do I need to do?*

- Recommended actions and strategies based on champion / challenger testing strategy outcomes
- Applying advanced analytical techniques to make specific recommendations

# Planning your strategy

- Problem analysis to mitigate risks
- Create your action plan
- Effort estimation
- Visualize and organize the results
- Creating business logic
- Report the results to leadership using analysis for:
  - various approach
  - different cost
  - benefits for each case

# More about the data...

- Is the data available?
- How do we access it?
- How big is the data?
- What type(s) of data is it?
- How good is the source of data?
- What is the scope (local vs global)?
- How frequently is collected?
- ... more



# More about the data...

- Plan ahead how to work with the data:
  - Collect it
  - Store it (is possible)
  - Study the dataset
  - Prepare and clean it
  - Analyze it
  - Organize and visualize it
  - Make conclusions about the data
  - Analyze small pieces more often than wait and process the entire dataset. It just might be too late!

# More about the data...

- Create and **keep data dictionary** containing:
  - databases
  - tables
  - fields
  - keys
  - schemas
  - data types
  - scope (local vs. global)
  - how it was loaded (using web services or ingested)
  - frequency (how often was the data collected?)

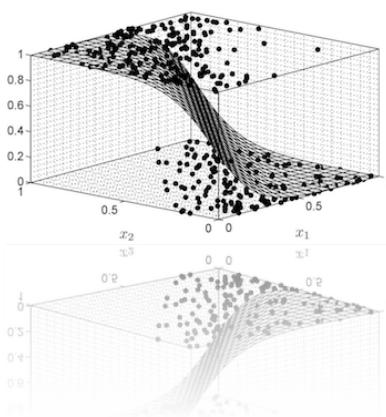
# Case study

- Example 1:

data source



statistical modeling

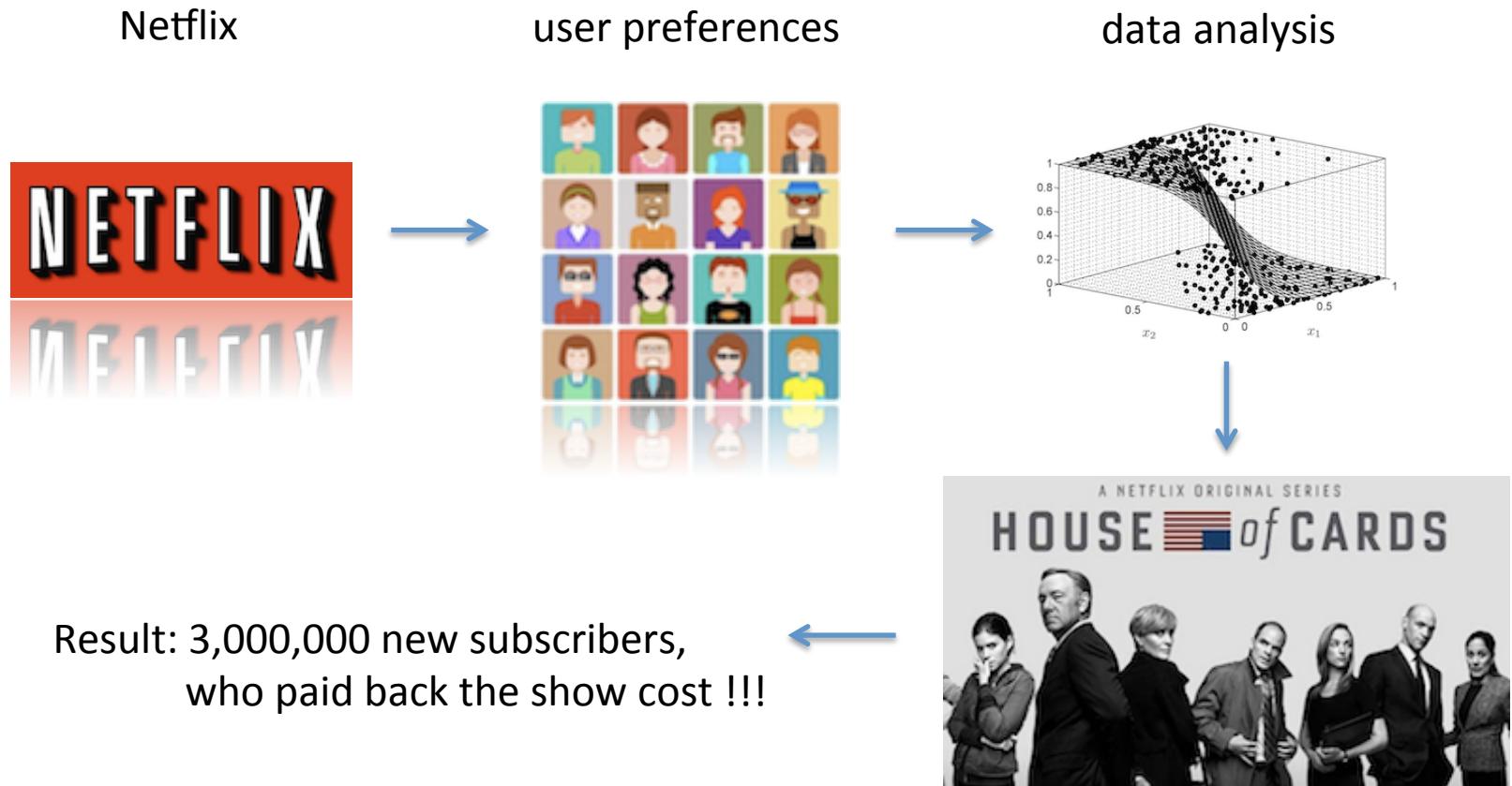


tailored offering



# Case study

- Example 2:



# The Excel power user

- Using Microsoft Excel spreadsheets containing data for analysis

## Pros

- Most business will use this method because **it is easy**
- It is **affordable**
- It is **versatile** as can be applied quickly

## Cons

- It provides a very **limited toolset** for analyzing and visualizing data
- **Not for large and complex datasets**

# Python for Data Analysis and Scientific Computing

- Python implementation started in 1989 by Guido Van Rossum
- Python intended to be the successor of the ABC language
- Python 2.0 released in late 2000
- Python 3.0 released in late 2008 and was mostly incompatible with its predecessor
- EOL for Python 2.7 was initially set at 2015, but postponed to 2020
- Google announced to work on a transcompiler for Python 2.7 – Jan.2017

# Python for Data Analysis and Scientific Computing

- Python comparison to Java:

- Dynamic vs static typing
- Braces vs indentation
- Speed and portability
- Interpreted vs compiled

# Python for Data Analysis and Scientific Computing

- Python with NumPy and SciPy – why using it?
  - Because you can develop **proof-of-concept** solutions fairly easily and free
  - **Open-source** software that is widely spread
  - **Fast** development time
  - **Many scientific libraries**
  - Mobile computing for **smartphone development** by using Kivy for Android and iOS
  - **Web-site incorporation** by using Django

# Python for Data Analysis and Scientific Computing

- Python – pros and cons:

## Pros

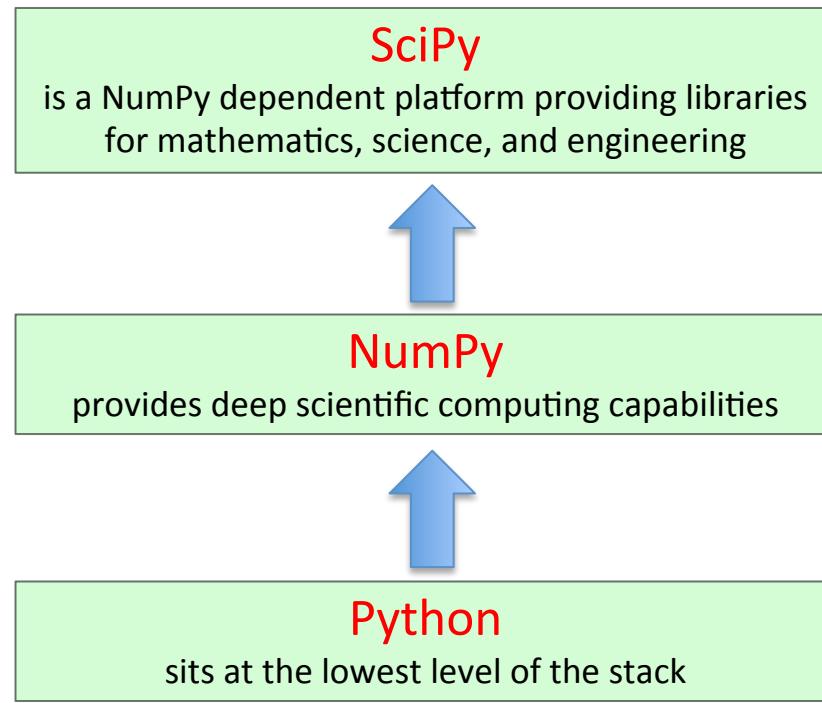
- High-level object-oriented language that is **easy to learn**
- **Good readability** of the code
- **Quick development time**
- **Rich scientific libraries** and many others like web server management, serial port access, etc.
- **Easy to share** and communicate results
- **Easy to read**, manipulate and process data
- A single environment for many tasks. No need to learn another language
- **Free and open-source** software that is widely spread

## Cons

- **Speed of execution** since it is higher level language
- **Mobile computing** for smartphone development (unless you use Kivy: Linux, Windows, OS X, Android and iOS)
- It is **not included in Web browsers**, because it is hard to secure
- Design flaws because it is **dynamically typed language** and some errors show at runtime such as syntactic errors (mistyping variable names)
- Python can be both **interpreted (.py)** and **compiled (.pyc)**
- **Documentation** in the past wasn't at the level of PHP or Java, but this has changed
- **Not all algorithms** found in other specialized platforms **or toolboxes** (Matlab) **were available in the past**

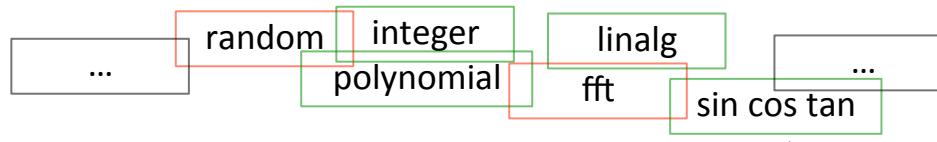
# Python for Data Analysis and Scientific Computing

- Basic stack using NumPy and SciPy for Python:



# Python for Data Analysis and Scientific Computing

- Python + NumPy
  - ... is like Matlab
  - NumPy provides the code for any fundamental data structure:



- There two basic data structures in NumPy:

## NDArray

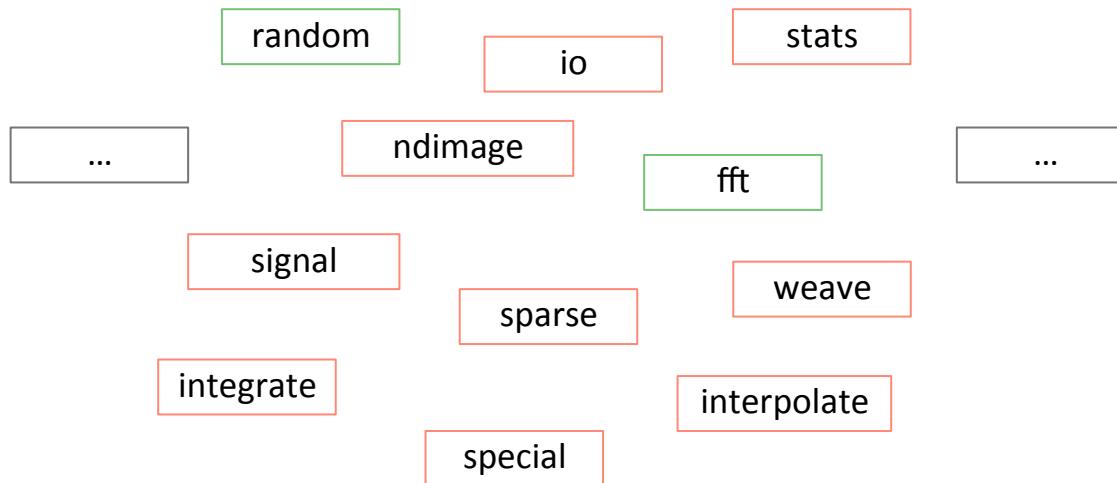
- is an N-dimensional array type
- describes a collection of “items” of the same type
- the items are indexed using N integers

## UFunc

- provides fast operations on arrays
- operates element-by-element on arrays
- example are: *sin, cos, tan, +, -, etc.*

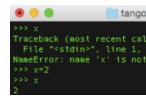
# Python for Data Analysis and Scientific Computing

- Python + NumPy + SciPy
  - ... is like Matlab + toolboxes
  - SciPy provides many specialized packages, which makes it very competitive and powerful platform



# Python for Data Analysis and Scientific Computing

- Choosing your environment:
  - Using the **shell** (primarily for Linux and OSX users):
    - Obtaining Python with core packages:  
<https://www.python.org/downloads/>
    - Obtaining the NumPy & SciPy libraries:  
<http://www.scipy.org/scipylib/download.html>
  - GUI environments using an interactive editor and debuggers for Python:
    - IPython – a GUI environment for interactive Python use with shells (terminal- and Qt-based):  
<http://www.ipython.org/>
    - Python(x,y) – free scientific and engineering environment for numerical computations, data analysis and data visualization:
    - Canopy – Python based environment providing core scientific analytic and scientific Python packages, for scientific development and visualization:  
<https://store.enthought.com/>
    - IEP – Interactive Editor for Python, is an interactive cross-platform Python IDE that is simple and efficient. IEP consists of an **editor**, a **shell**, and a set of **tools** to help the developer  
<http://www.iep-project.org/>
    - Pyzo – a free open-source environment based on Python:  
<http://www.pyzo.org/downloads.html>
    - ... other

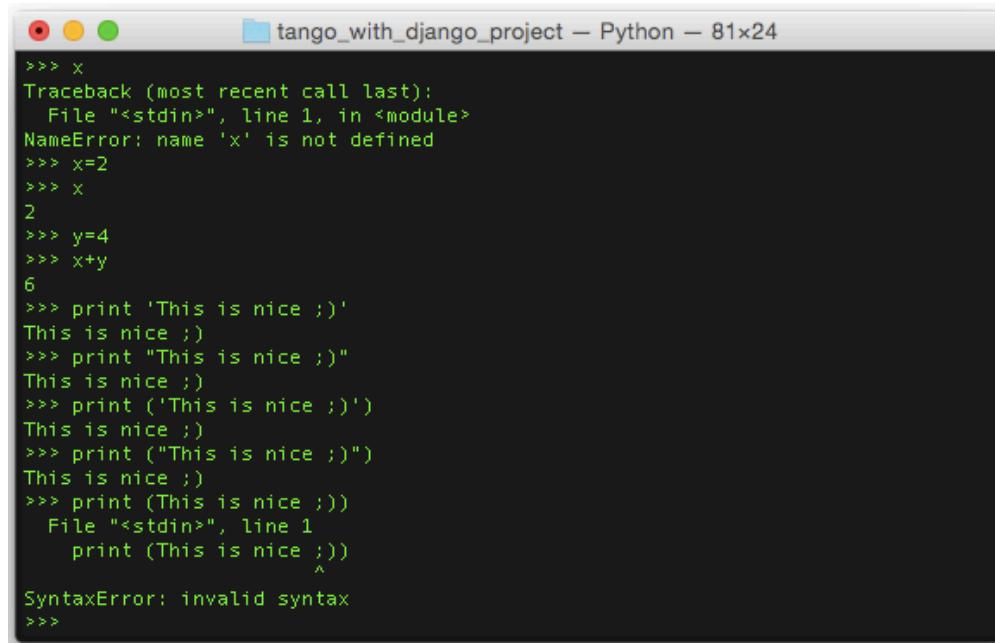


# Python for Data Analysis and Scientific Computing

- Working with the shell

- The shell:

- present on **OSX, Unix** and **Linux systems**
    - provides a **quick and easy to use** plain text environment for swift execution of simple commands
    - it is **flexible** and **lightweight**
    - it has a Unix/Linux look and feel and is **not used only for Python** and its packages and modules



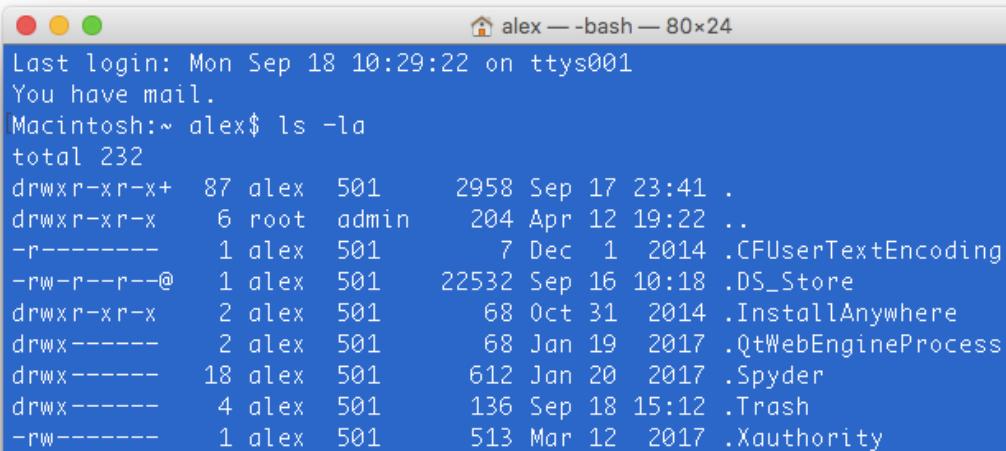
The screenshot shows a terminal window with the title "tango\_with\_django\_project — Python — 81x24". The window displays the following Python code and its execution:

```
>>> x
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
NameError: name 'x' is not defined
>>> x=2
>>> x
2
>>> y=4
>>> x+y
6
>>> print 'This is nice ;)'
This is nice ;)
>>> print "This is nice ;)"
This is nice ;)
>>> print ('This is nice ;)')
This is nice ;)
>>> print ("This is nice ;)")
This is nice ;)
>>> print (This is nice ;))
  File "<stdin>", line 1
    print (This is nice ;))
                           ^
SyntaxError: invalid syntax
>>>
```

# Python for Data Analysis and Scientific Computing

- Using the shell (linux and OSX users):
  - Using ‘pip’: a tool for installing Python packages from the [Python Package Index](#)
    - In your terminal type:  
\$ sudo easy\_install pip, or  
\$ sudo pip install --upgrade pip
    - Then you can install python and some packages:  
\$ pip install python  
\$ pip install numpy  
\$ pip install scipy  
\$ pip install matplotlib  
\$ pip install django

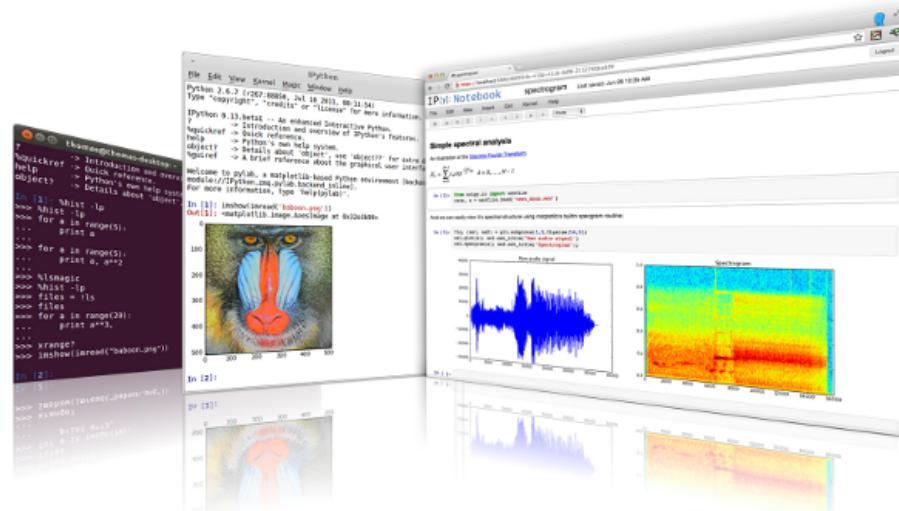
*\* a side note on permissions and ownership:*



```
Last login: Mon Sep 18 10:29:22 on ttys001
You have mail.
Macintosh:~ alex$ ls -la
total 232
drwxr-xr-x+  87 alex  501      2958 Sep 17 23:41 .
drwxr-xr-x+   6 root   admin     204 Apr 12 19:22 ..
-r-----    1 alex  501       7 Dec  1 2014 .CFUserTextEncoding
-rw-r--r--@   1 alex  501     22532 Sep 16 10:18 .DS_Store
drwxr-xr-x+   2 alex  501      68 Oct 31 2014 .InstallAnywhere
drwx-----    2 alex  501      68 Jan 19 2017 .QtWebEngineProcess
drwx-----   18 alex  501     612 Jan 20 2017 .Spyder
drwx-----    4 alex  501     136 Sep 18 15:12 .Trash
-rw-----    1 alex  501     513 Mar 12 2017 .Xauthority
```

# Python for Data Analysis and Scientific Computing

- Working with IPython
  - IPython
    - Is a GUI environment for interactive Python use with shells (terminal- and Qt-based) and these key features:
      - **Interpreters** that can be embedded into users' custom projects
      - It has a **build-in web-based** interactive notebook environment with complete set of shell features and support for embedded figures, animations and **multimedia**
      - Use of **GUI toolkits** and support for interactive data **visualization**
      - A high-performance **library for parallel computing** portable to multicore systems, clusters and **cloud computing**



# Python for Data Analysis and Scientific Computing

- Working with IPython and the editor
  - IPython and the editor
    - History, persistent across sessions
    - Caching of output results during a session
    - Extensible tab completion, with support by default for completion of python variables and keywords, filenames and function keywords
    - Comprehensive system of ‘magic’ commands for controlling the environment and performing many tasks related either to IPython or the operating system
    - A rich configuration system with easy switching between different setup
    - System shell access with user-extensible alias setup
    - Session logging and reloading
    - Code is easy to embed and use in other Python programs and GUIs
    - Embedded debugger

# Python for Data Analysis and Scientific Computing

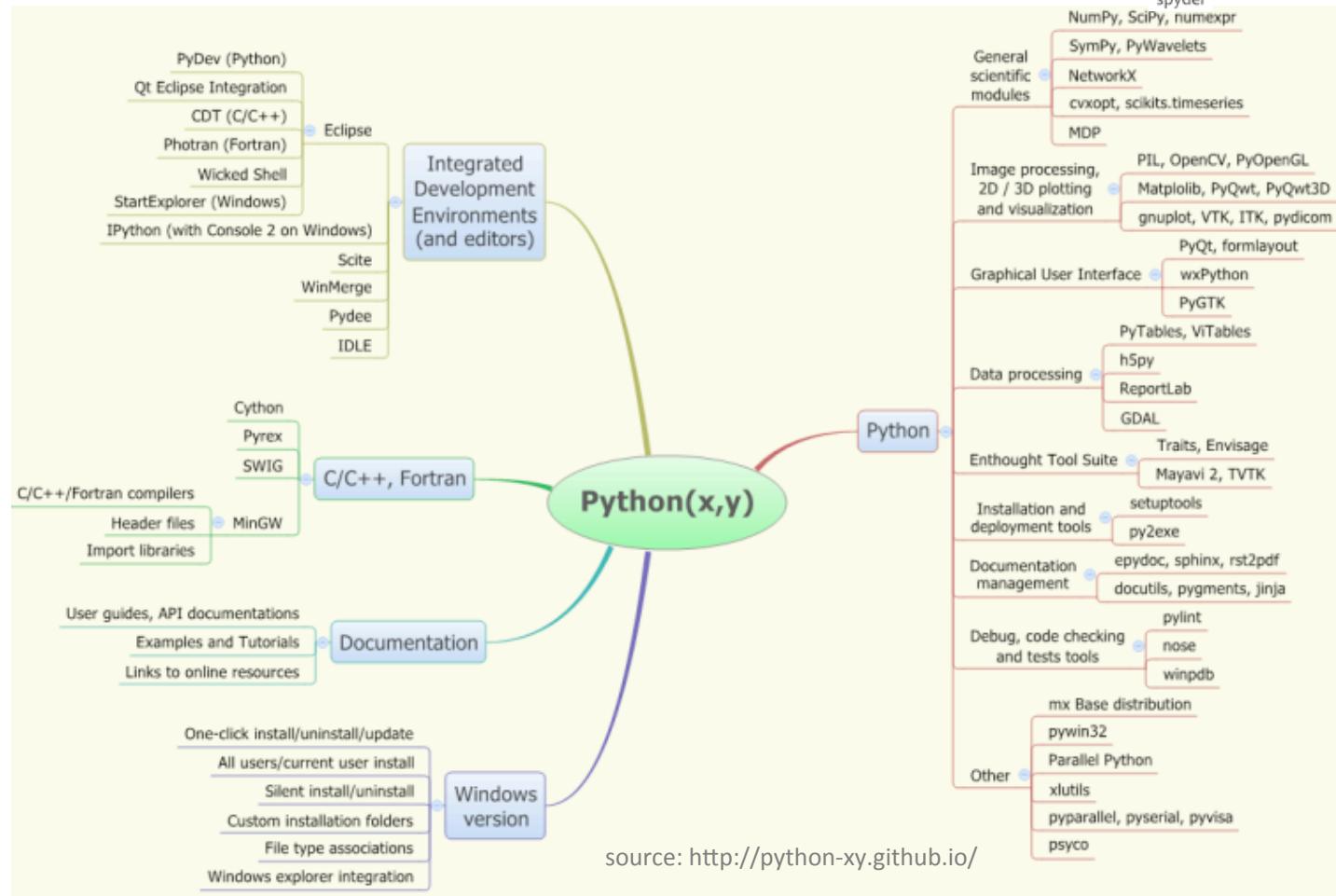
- Working with alternative environments:

- Python(x,y) – a free scientific-oriented Python Distribution based on Qt and Spyder –



It offers:

- object-oriented prog.
- MATLAB-like syntax
- parallel computing
- 2D & 3D plotting
- any level coding



source: <http://python-xy.github.io/>

# Python for Data Analysis and Scientific Computing

- Working with alternative environments:
  - Canopy – scientific and analytic Python package distribution. **has a free version, but limited**

It offers:

- object-oriented prog.
- MATLAB-like syntax
- many packages
- graphical debugger is **NOT** free

The screenshot shows the Enthought Canopy product page with the following details:

Category	Canopy Express	Canopy with Python Essentials	Canopy with Python Foundations All Access	Canopy Triple Play
Price	FREE	\$199 / year	\$678 / year	\$1,125 / year
Contact Us	Contact us for discounts for 5+ users			
Download	Free Download	Add to Cart	Add to Cart	Add to Cart
Python Packages	100+ Core	300+ Extended		
See Included Python package details <a href="#">+</a>				
Integrated Analysis Environment				
Graphical Debugger		●	●	●
Integrated IPython	●	●	●	●
Advanced Code Editor	●	●	●	●
Application Development Platform	●	●	●	●

# Python for Data Analysis and Scientific Computing

- Working with alternative environments:
  - Canopy – scientific and analytic Python package distribution. **has a free version, but limited**

Package Manager:

- Packages are installed through Package Manager

The screenshot shows the Enthought Canopy Package Manager interface. At the top, there's a navigation bar with 'Refresh' and 'Not logged in.' On the right, a search bar contains the text 'seaborn'. Below the search bar, a table lists packages. The table has three columns: 'Package Name', 'Latest Available Version', and a header row for 'Installed' (0/106), 'Available' (1/494), 'Updates' (0/0), and 'History'. The 'Available' column is highlighted in green. The table shows one package, 'seaborn', with the latest available version being '0.7.0-6'. At the bottom of the page, there's a detailed view for the 'seaborn' package. It shows the package name 'seaborn', a note 'No summary available.', and two status boxes: 'Installed:' (Currently not installed) and 'Available on store: (enthought/free)' with a dropdown menu showing '0.7.0-6'. A blue button labeled 'Install v0.7.0-6' is located to the right of the dropdown.

Installed	Available	Updates	History	Package Name	Latest Available Version
0/106	1/494	0/0		seaborn	0.7.0-6

**seaborn**

No summary available.

Installed: Currently not installed.

Available on store: (enthought/free)

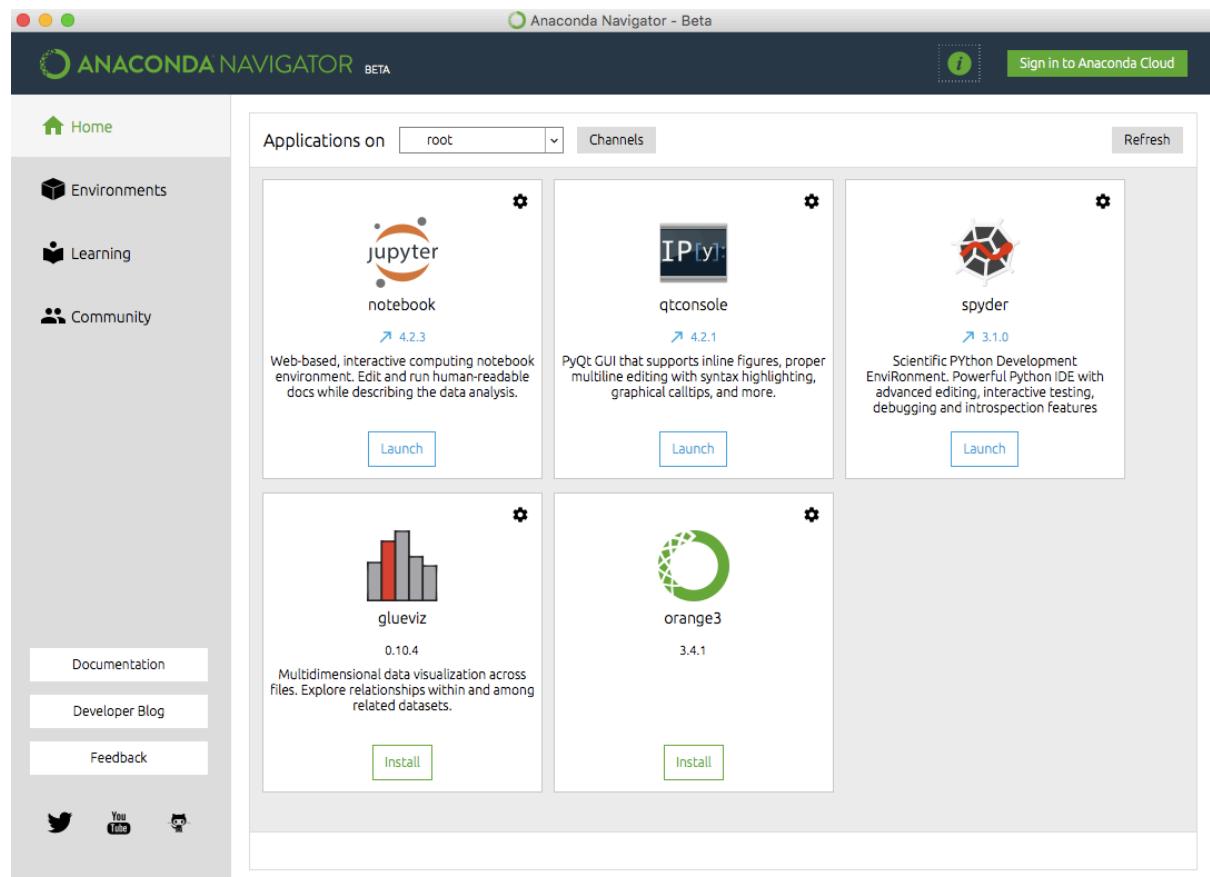
0.7.0-6

Install v0.7.0-6

No description available.

# Python for Data Analysis and Scientific Computing

- Working with alternative environments:
  - Anaconda Navigator



# Python for Data Analysis and Scientific Computing

- Working with alternative environments:
  - Anaconda Navigator

Available are:

- Jupyter Notebook
- IP[y] qtconsole
- Spyder
- glueviz

- **Jupyter**: is **also available in Canopy**: it is a **web-based** interactive computing **notebook environment**. Easy to read and describe the data analysis
- **IP[y] qtconsole**: is a PyQt based GUI working with various figures and plots among other editing features
- **Spyder**: is **for scientific python development** with interactive testing env **and debugging features**. It can be installed independently in some platforms such as Linux
- **glueviz**: n-dimensional data visualization using multiple files. View the relationships among multiple datasets

# Python for Data Analysis and Scientific Computing

- Working with alternative environments:
  - Canopy – scientific and analytic Python package distribution.

The screenshot shows a Jupyter Notebook interface titled "Titanic Berkeley Project V2 - final". The notebook has a toolbar with File, Edit, View, Insert, Cell, Kernel, and Help options. A status bar at the bottom right indicates "Python 2". The main area displays two code cells:

**Setting up Python Environment needed**

```
In [*]: #Importing all the libraries needed
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline
import brewer2mpl
import seaborn as sns
from scipy import stats
from scipy import linalg

/Users/alex/Library/Enthought/Canopy_64bit/User/lib/python2.7/site-packages/matplotlib/font_manager.py:273: UserWarning: Matplotlib is building the font cache using fc-list. This may take a moment.
warnings.warn('Matplotlib is building the font cache using fc-list. This may take a moment.')

In [15]: # Setting up better defaults for matplotlib
from matplotlib import rcParams

rcParams['axes.facecolor'] = 'white'
rcParams['font.size'] = 18

#colorbrewer2 Dark2 qualitative color table
dark2_colors = brewer2mpl.get_map('Dark2', 'Qualitative', 7).mpl_colors
```

Annotations on the left side of the interface:

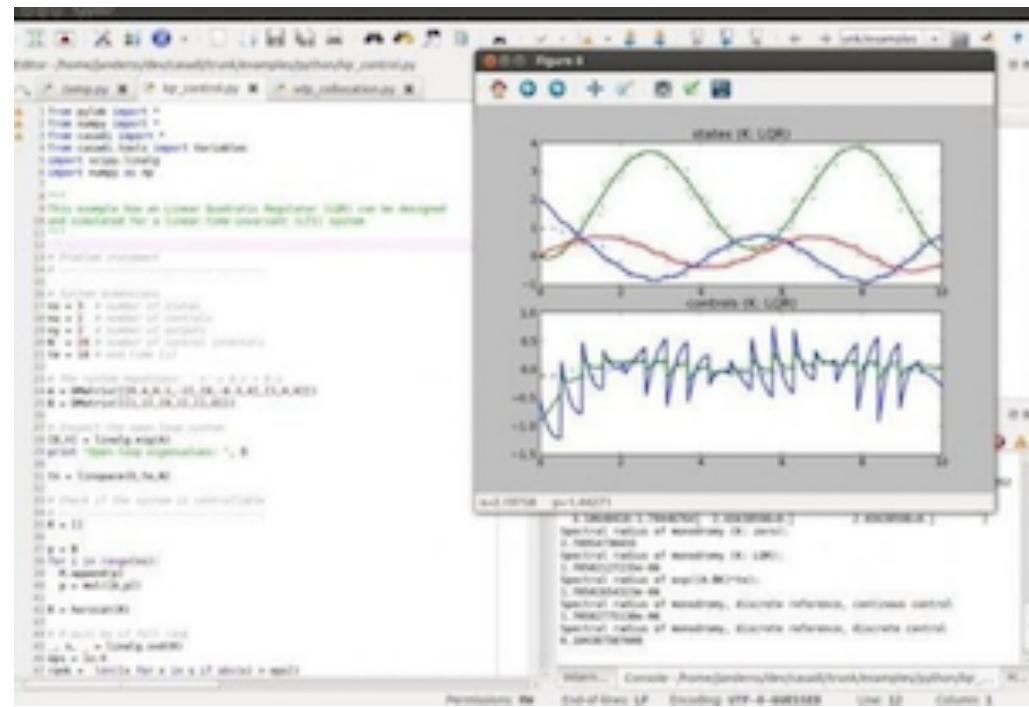
- [\*] the line is in a process of executing ←
- Execute your code sequentially using Shift+Enter

# Python for Data Analysis and Scientific Computing

- Working with alternative environments:

- Spider ...

- Spider is a free Open-source

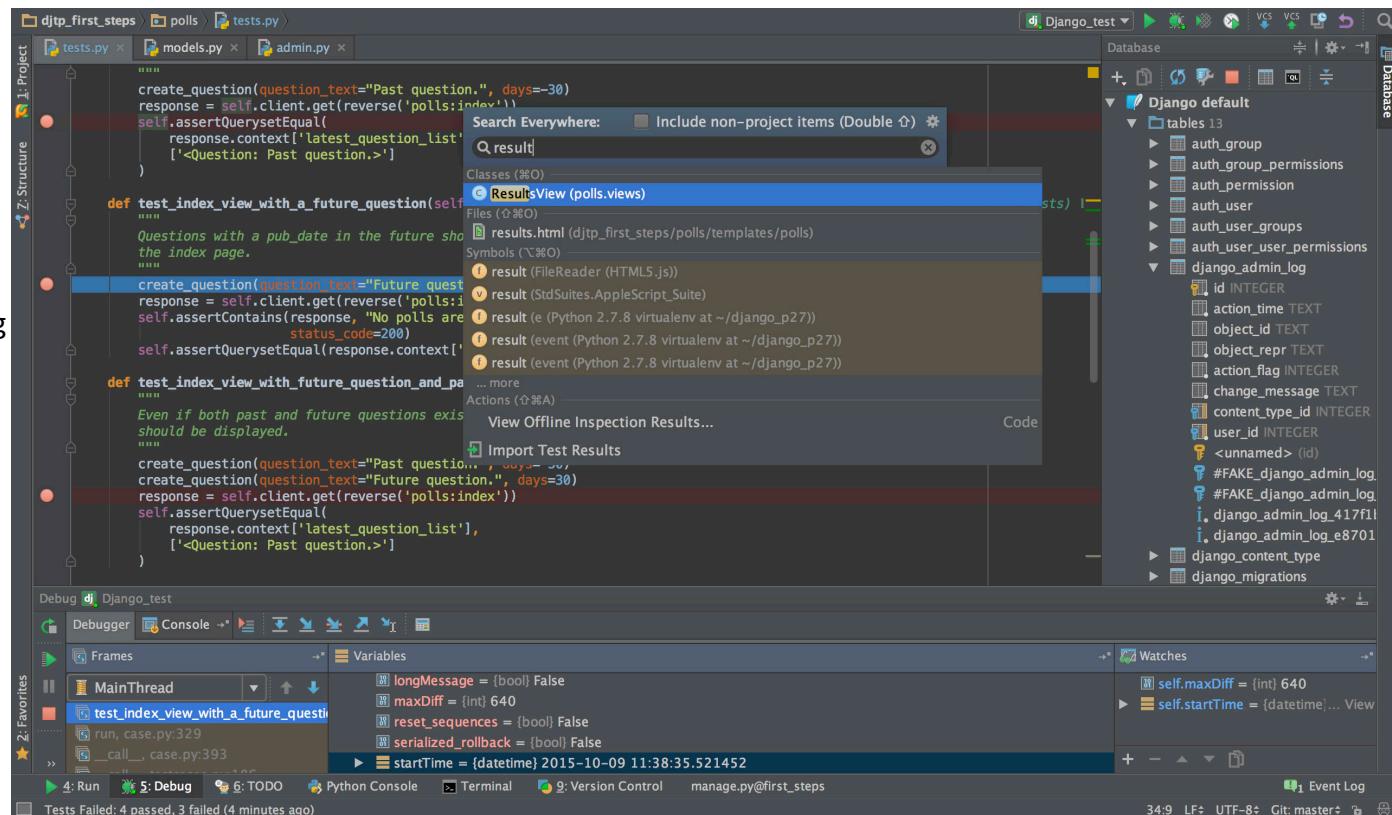


# Python for Data Analysis and Scientific Computing

- Working with alternative environments:
  - Pycharm – scientific and analytic Python package distribution. **has a free version, but limited**

It offers:

- object-oriented prog.
- save time by auto completing the routines
- on-the-fly error checking
- focus on bigger picture
- many features
- many packages



source: <https://www.jetbrains.com/pycharm/>

# Python for Data Analysis and Scientific Computing

- Working with alternative environments:

- Eric ...

- Eric*** is a free Open-source



The screenshot shows the Eric development environment interface. At the top is a menu bar with options like File, Edit, View, Debug, Unitest, Project, Refactoring, Extras, Settings, Window, Bookmarks, and Help. Below the menu is a toolbar with various icons for file operations, search, and navigation. The main workspace contains several panes: a "Project Browser" pane on the left showing a "Name" dropdown and a "VCS Status" section; a central "Code Editor" pane displaying Python code with syntax highlighting; a "Log Viewer" pane at the bottom left; a "Task Viewer" pane at the bottom right; and a "Template Viewer" pane on the right side. The status bar at the bottom shows the file path "File: /home/ethan/Desktop/test.py" and the current line and position "Line: 1 Pos: 0".

```
#!/usr/bin/python
# -*- coding: utf-8 -*-

"""
Triple quoted nonsense.
b = """Triple quoted with an escaped quote '."""
c = """Triple single quotes. Take that, Internet!"""
d = """This is a double quote. """
# In case the editor gets confused, this line will be helpful.
def foobar():
    return c

if __name__ == '__main__':
    foobar()
```

# Python for Data Analysis and Scientific Computing

- Working with alternative environments:

- Pyzo ... what?

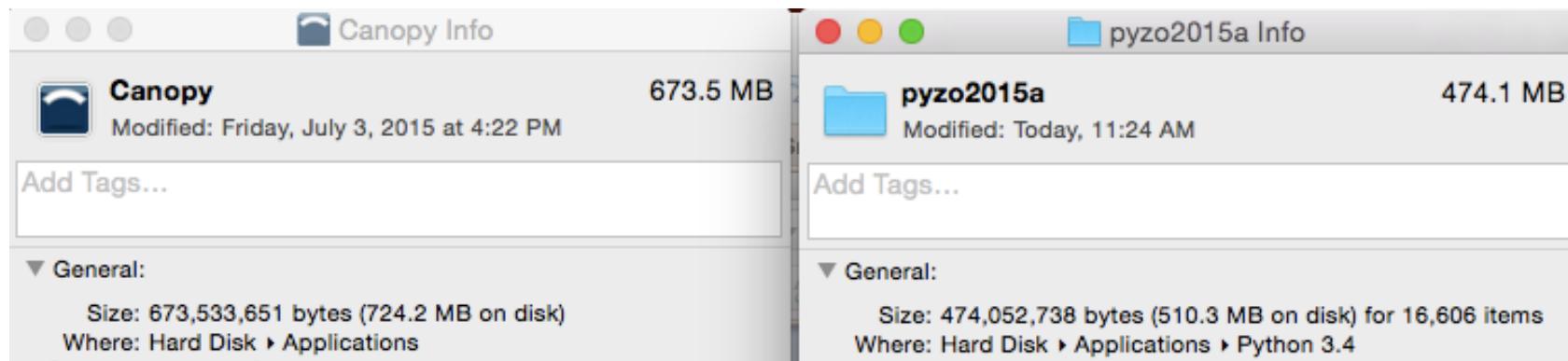
- **Pyzo** is a **free Open-source** BSD-licensed
    - It is an easy to use **GUI** computing environment **based on Python 3**
    - It is built to make Python's potential for **scientific and engineering computing** to be easily accessible
    - Pyzo is combined with IEP to make a powerful easy to use comprehensive environment for scientific computing and visualization
    - Pyzo is a **comprehensive development environment** pre-built with:
      - many scientific packages including **NumPy**, **SciPy** and **Matplotlib**
      - uses **IEP as the front-end IDE**
      - **IPython** and a **notebook**
    - Complete packages for Pyzo are given here: <http://www.pyzo.org/packages.html>



PYZO: PYTHON TO THE PEOPLE

# Python for Data Analysis and Scientific Computing

- Working with alternative environments:
  - Pyzo takes less space than Canopy (without Canopy debugger installed)



# Python for Data Analysis and Scientific Computing

- Python packages for data analysis:
  - **NumPy**: provides fast precompiled functions for numerical routines, large n-dimensional arrays, and matrices, as well as high-level mathematical functions to operate on these arrays
  - **SciPy**: provides scientific, technical, and deep analytical computing functionality and contains modules for statistics, optimization, linear algebra, integration, interpolation, and plenty more
  - **Matplotlib**: is a plotting library used by NumPy that provides the ability to embed plots into applications using GUI toolkits like wxPython, Qt, or GTK+
  - **Pandas**: offers data structures and manipulation and analysis of numerical tables

# Python for Data Analysis and Scientific Computing

- Python packages for data analysis:
  - **Scikit-learn**: contains various algorithms for classification, regression and clustering such as: naive Bayes, SVM, logistic regression, random forests, gradient boosting, k-means and DBSCAN. It is designed to use the scientific libraries NumPy and SciPy
  - **Theano**: with Theano you can define, evaluate, and optimize mathematical expressions applied to n-dimensional arrays
  - **Statsmodels**: provides statistical models, statistical tests, plotting functions, and descriptive statistics, while working with different types of data
  - **Other**: Mipy, NLTK, Keras, PyBrain, Hebel, nolearn, Orange

# Python for Data Analysis and Scientific Computing

- Working with Pyzo



**Download**  
**About Python**  
**About Pyzo**  
**Learn**  
**FAQ**

## Packages

Pyzo comes with a range of packages for general functionality, science, development, testing, UI design, etc. ... Click on any of the links to go to the projects website where you will find more information such as tutorials, documentation, and mailing lists.

The [Scipy Stack](#):

- [numpy](#) - using arrays in Python
- [scipy](#) - core scientific functionality
- [matplotlib](#) - 2D plotting library
- [pandas](#) - data structures and analysis
- [sympy](#) - symbolic math
- [nose](#) - software testing
- [IPython](#) - advanced interactive shell

Further scientific packages:

- [pyopengl](#) - talk to your GPU
- [visvis](#) - 3D plotting and visualization
- [skimage](#) - 2D and 3D image processing
- [sklearn](#) - machine learning
- [imageio](#) - read and write image data

# Python for Data Analysis and Scientific Computing

- Working with Pyzo
  - Useful magic commands:

Linux-like:

- pwd
- ls
- mkdir
- rmdir
- whos
- dir()
- cd
- history

Python IEP/Pyzo specific:

- edit *file*
- del(parameter)
- reset
- reset -f
- run *file*
- time
- colors (nocolor,linux,lightbg)
- dhist



PYZO: PYTHON TO THE PEOPLE

Tip: Using a magic command followed by '?' will give you help for it. Example: 'cd?'

magic comprehensive list: <http://ipython.org/ipython-doc/3/interactive/magics.html>

# Python for Data Analysis and Scientific Computing

- How to check your packages? .. (In your interactive editor - I use Pyzo)

```
In [3]: import pip  
In [4]: pip  
Out[4]: <module 'pip' from '/Applications/Python 3.4/pyzo2015a/lib/python3.4/site-packages/pip/__init__.py'>  
In [5]: pip list  
Out[5]:  
You are using pip version 6.0.8, however version 7.0.1 is available.  
.... snip snip snip ...  
Django 1.8.1  
.... snip snip snip ...  
matplotlib 1.4.2  
.... snip snip snip ...  
numpy 1.9.1  
.... snip snip snip ...  
scipy 0.15.1      ->      one can use: "scipy.__version__" to check the package version (scipy in this case)
```

or

```
In [5]: pip.get_installed_distributions()  
Out[5]:  
.... snip snip snip ...  
Django 1.8.1 (/Applications/Python 3.4/pyzo2015a/lib/python3.4/site-packages),  
.... snip snip snip ...  
matplotlib 1.4.2 (/Applications/Python 3.4/pyzo2015a/lib/python3.4/site-packages),  
.... snip snip snip ...  
numpy 1.9.1 (/Applications/Python 3.4/pyzo2015a/lib/python3.4/site-packages),  
.... snip snip snip ...  
scipy 0.15.1 (/Applications/Python 3.4/pyzo2015a/lib/python3.4/site-packages),
```

- Using:

```
In [6]: help()
```

# Python for Data Analysis and Scientific Computing

- Python modules, packages and scientific blocks
  - A **module** is a Python file with .py extension containing definitions and statements
  - Modules implement a set of functions and can be imported into other modules using the import command:

Example:

Lets create our first Python program called – *lecture1.py*:

```
# Let's make a simple addition until certain criteria is not met
def add(a,operation):           # the criteria to be met are: 'a' and 'operation'
    if operation == 'add':
        x, y = 2, 4
        while x < a:
            print(x,y),
            x = x + y
```

Possible usage:

```
# importing a module called 'lecture1'
import lecture1      -> must be in the same working folder
lecture1.add(24, 'add')
```

Tips:

1. to check what's loaded in memory simply look at the workspace
2. if you enter an infinite loop hit 'ctrl-c' to break the loop and stop the execution
3. try debugging your program using breakpoints

# Python for Data Analysis and Scientific Computing

- Python modules, packages and scientific packages
  - A *module* in Python is a source file that can expose classes, functions and global variables
  - A *package* in Python is a directory of one or several Python modules (source .py files)

Example:

# 1. this is how we import a *package*:

```
import math
```

Usage:

```
math.acos(0.3454453)
```

# 2. this is how we import a *module* from a package:

```
import acos from math ... ( or: from math import acos)
```

Usage:

```
acos(0.3454453)
```

Tips:

1. Observe what's loaded in workspace in the two cases above
2. you can delete a particular item from workspace by using '*del(item)*'
3. to clear the workspace use '*reset*'
4. The built-in function '*dir()*' is used to find what's currently defined and shown in workspace

# Python for Data Analysis and Scientific Computing

- Python modules, packages and scientific blocks
  - 'sys': includes information about the modules available to us. It includes built-in and imported modules
  - **Imported modules:** usually Python programs are a combination of several modules imported with a main application  
They can be seen in the following way:

```
import sys, textwrap
names = sorted(sys.modules.keys())
name_text = ' || '.join(names)
print(textwrap.fill(name_text))
```
  - **Built-in modules:** the interpreter can be built and compiled with some C modules, that do not appear in the imported modules list since they were included in the package  
They can be seen in the following way:

```
import sys
for module_name in sys.builtin_module_names:
    print(module_name)
```
  - **The search path:** the search path for modules is managed as a list saved in **sys.path**. The default path list consist of the script directory used to start the application and the current working directory  
**sys.path**

# Python for Data Analysis and Scientific Computing

Magic commands in **IPython**:

List of important **Magic** commands in **IPython**:

%automagic	magic [-brief, ]
ls	lsmagic
pwd	ed, edit
cd	load
mv	run
dir	time
mkdir	writefile
rm, rmdir	env
who, whos	system
alias, unalias	config
hist, history	save
dhist	logoff logon logstart logstate logstop

# Python for Data Analysis and Scientific Computing

HW assignment 1:

Make sure you have chosen, installed and configured your Python-NumPy-SciPy environment