

# Scientific Programming in Python

Inteligencia Artificial en los Sistemas de Control Autónomo  
Máster Universitario en Ingeniería Industrial

Departamento de Automática

## Objectives

1. Introduce some Python tools for scientific programming.
2. Motivate the need of efficient matrix manipulation.
3. Handle matrices and dataframes in Python.
4. Basic data visualization with Python.

## Bibliography

Jake VanderPlas. Python Data Science Handbook. Chapters 1, 2, 3 and 4. O'Reilly. (Link).

# Table of Contents

## 1. Overview

- Data Science
- The data scientist toolkit
- Anaconda
- Python IDEs for Data Science

## 2. Basics

- Magic commands
- Pasting code blocks
- Running external code
- Input and output history
- iPython shell commands
- Automagic

## 3. NumPy

- Understanding Data Types in Python
- Introduction
- Matrix creation
- NumPy data types
- NumPy array attributes
- Accessing single elements
- Accessing subarrays
- Reshaping of arrays
- Concatenation of arrays

- Splitting of arrays
- Universal functions
- Aggregations
- Broadcasting
- Structured arrays

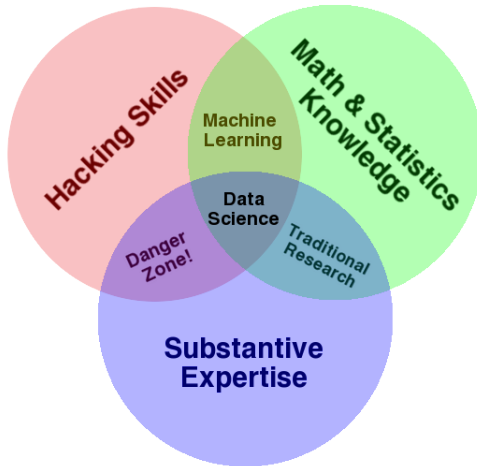
## 4. Pandas

- Introduction
- The Pandas **Series** object
- The Pandas **DataFrame** object
- Constructing **DataFrame** objects
- Data indexing and selection
- Operating on data
- Missing data
- Combining datasets: `pd.concat()`
- Combining datasets: `pd.merge()`
- Aggregation in Pandas
- Grouping in Pandas

## 5. Visualization

- Visualization examples
- Introduction to Matplotlib
- Introduction to Seaborn
- Seaborn: Distribution

# Data Science























## Pasting code blocks: %paste and %cpaste

- %paste: Paste one time
- %%cpaste: Paste several times

```
In [25]: %cpaste
Pasting code; enter '--' alone on the line
to stop or use Ctrl-D.
:         def donothing(x):
            return x:
:--
```

```
In [20]: %paste
def donothing(x):
    return x

## -- End pasted text --
```











iPython  
Automagic

## Problems with some shell commands

```
In [23]: !pwd
/repositorios/pythonCourse
In [24]: !cd ..
In [25]: !pwd
/repositorios/pythonCourse
```

Some magic commands here to help

- %cd,%ls,%mkdir,%pwd,
- ...

Those magics are regularly used ...

- ... so common that % is no longer required (automagic)
- Working with iPython is almost like working with a Unix-like shell

## Automagic commands

cat, cp, env, ls, man, mkdir, more,  
mb, pwd, rm and rmdir

## NumPy

# Understanding Data Types in Python (I)

```
/* C code */
int result = 0;
for(int i=
```

- Data types must be declared
- Data types cannot change

## Dynamic typing

```
# Python code
result = 0
for
```

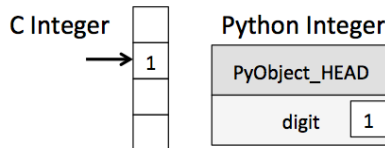
## NumPy

## Understanding Data Types in Python (II)

## Dynamic typing must be implemented somewhere ...

Python 3.4 source code

```
struct _longobject {
    long ob_refcnt;
    PyTypeObject *ob_type;
    size_t ob_size;
    long ob_digit[1];
};
```

























# NumPy

## Splitting of arrays

Three methods to split arrays

- `np.split()`
- `np.vsplit()`
- `np.hsplit()`

`np.split()`

```
In [1]: x = [1, 2, 3, 99, 99, 3, 2, 1]
In [2]: x1, x2, x3 = np.split(x, [3, 5])
In [3]: print(x1, x2, x3)
[1 2 3] [99 99] [3 2 1]
```

`np.vstack()`

```
In [1]: grid = np.arange(16).reshape((4, 4))
In [2]: print(grid)
[[ 0  1  2  3]
 [ 4  5  6  7]
 [ 8  9 10 11]
 [12 13 14 15]]
In [3]: upper, lower = np.vsplit(grid, [2])
In [4]: print(upper)
[[0 1 2 3]
 [4 5 6 7]]
```













































# Pandas

## Introduction

A data science workflow needs more features

- Label columns and rows
- Missing data
- Operations on groups
- Data input

Pandas implements all those features, and more

- Built on NumPy's ndarray

Pandas provides two main objects

- Series
- DataFrame

### Convention

```
import numpy as np
import pandas as pd
```

























# Pandas

## Missing data (I)

NumPy supports missing data in floating-point data

- Specific value defined by IEEE
- Available as `np.nan`

Pandas supports missing data through two mechanisms

- `None` object, interpreted as NaN (Not a Number)
- `np.nan`: for floating-point data
- Almost automatic NaN handling (types upcast)

```
>>> pd.Series([1, np.nan, 2, None])
0      1.0
1      NaN
2      2.0
3      NaN
dtype: float64
```





















# Pandas

## Aggregation in Pandas (I)

The first step in data analysis is summarization

- First contact with data
- Insight to the dataset

Aggregation methods

- Applied to columns

AGGREGATION	DESCRIPTION
<code>count()</code>	Total number of items
<code>first(), last()</code>	First and last item
<code>mean(), median()</code>	Mean and median
<code>min(), max()</code>	Minimum and maximum
<code>std(), var()</code>	Standard dev. and variance
<code>mad()</code>	Mean absolute deviation
<code>prod()</code>	Product of all items
<code>sum()</code>	Sum of all items
<code>describe()</code>	Data summary

```
>>> import seaborn as sns
>>> planets = sns.load_dataset('planets')
>>> planets.head()
```

		method	number	orbital_period	mass	distance	year
0	Radial Velocity	1	269.300	7.10	77.40	2006	
1	Radial Velocity	1	874.774	2.21	56.95	2008	
2	Radial Velocity	1	763.000	2.60	19.84	2011	
3	Radial Velocity	1	326.030	19.40	110.62	2007	
4	Radial Velocity	1	516.220	10.50	119.47	2009	

```
>>> planets.dropna().describe()
```

	number	orbital_period	mass	distance	year
count	498.00	498.000000	498.00	498.0000	498.000
mean	1.73	835.778671	2.50	52.0682	2007.377
std	1.17	1469.128259	3.63	46.5960	4.167
min	1.00	1.328300	0.00	1.3500	1989.000
25%	1.00	38.272250	0.21	24.4975	2005.000
50%	1.00	357.000000	1.24	39.9400	2009.000
75%	2.00	999.600000	2.86	59.3325	2011.000
max	6.00	17337.500000	25.00	354.0000	2014.000

```
>>> planets.mean()
```

number	1.785507
orbital_period	2002.917596
mass	2.638161
distance	264.069282
year	2009.070531
dtype:	float64











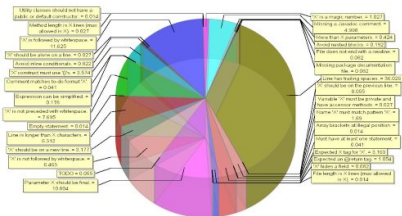




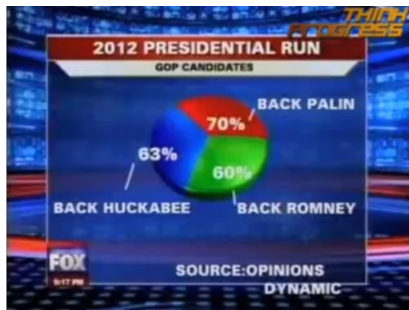


# Visualization

## Bad visualization examples (I)



(Source)



(Source)













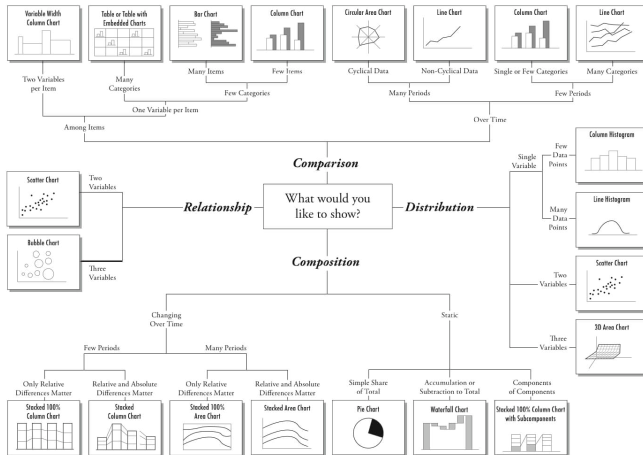


# Visualization

## Motivation (III)

### Chart Suggestions—A Thought-Starter

www.ExtremePresentation.com  
© 2009 A. Abela — a.x.abela@gmail.com



(Source)













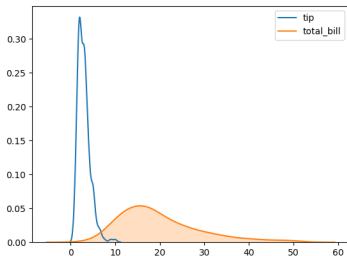






















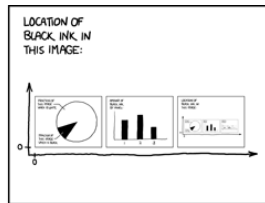
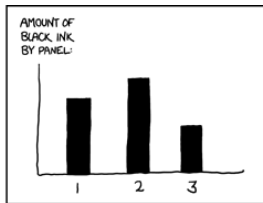
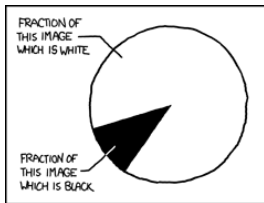












(Source)