| Title of Work. | No. of Sentences. | Average Number of Words per Senten | | | | |
|---|---|---|---|---|---|---|
| | | First 100. | Second 100. | Third 100. | Fourth 100. | Fifth 100. |
| r Buergergeneral ' ......... | 500 | 6.7 | 5.1 | 4.5 | 3.9 | 4.7 |
| etz v. Berlichingen ' * ...... | 2500 | 8.7 | 9.2 | 8.7 | 7.8 | 8.1 |
| tters to Frau v. Stein ' .... | 500 | 13.5 | 12.5 | 11.3 | 11.2 | 12.4 |
| ust ': First Part ......... | 500 | 14.6 | 15.2 | 13.6 | 13.0 | 12.0 |
| ust ': Second Part ........ | 500 | 16.3 | 17.2 | 15.5 | 12.3 | 16.5 |
| inecke Fuchs ' ............ | 500 | 18.5 | 16.5 | 16.9 | 14.8 | 16.6 |
| e Leiden d. j. Werthers ' ... | 500 | 20.7 | 20.9 | 19.1 | 22.7 | 18.2 |
| liaenische Reise: Rom ' ... | 500 | 23.1 | 23.7 | 22.7 | 21.5 | 22.7 |

# Introduction to Data Analysis with Python

Prof. David Losada, Universidade de Santiago de Compostela
**Email**: david.losada@usc.es
Marcos Fernández, Investigador predotural, University of Santiago de Compostela
**Email**: marcosfernandez.pichel@usc.es

# Why Python?

- Python is a modern, open source, **object-oriented programming** language, created by a Dutch programmer, Guido van Rossum, in 1991.
- Officially, it is an **interpreted scripting language**. There are several implementations, being the most used the one implemented in C.
- It offers the power and flexibility of lower level (i.e. compiled) languages, without the steep learning curve, and without most of the associated debugging pitfalls.
- The language is very **clean and readable**, and it is available for almost every modern computing platform.
- It is widely used: Google, DropBox, Deutsche Borse, NASA, etc.

# Why Python?

In Data Science, **programming is mostly about discovering solutions by writing code and using a language to express them.**

Python offers a number of advantages to data scientists:

- Powerful and easy to use.
- Anything that can be coded in C, FORTRAN, or Java can be done in Python, almost always in fewer lines of code, and with fewer debugging headaches.
- Its standard library is extremely rich, including modules for string manipulation, regular expressions, file compression, mathematics, scientific programming, profiling and debugging.
- Unnecessary language constructs, such as END statements and brackets are absent, making the code efficient and easy to read.
- Python is object-oriented, which allows data structures to be abstracted in a natural way.

# Why Python?

- It is interactive
  - Python may be run interactively on the command line. Commands may entered serially followed by the Return key.
  - This is useful for exploratory programming.
- It is extensible
  - Python is often referred to as a "glue" language, meaning that it is a useful in a mixed-language environment. C or FORTRAN code can be compiled directly into Python programs.
  - Sometimes, it can be slow relative to compiled languages. In many cases this performance deficit is due to a short loop of code that runs thousands or millions of times. Such bottlenecks may be easily removed by coding a function in C or Cython that can be compiled into a Python module.

# Why Python?

- There is a vast body of Python third party modules:
  - **NumPy**: Numerical Python (NumPy) is a set of extensions that provides the ability to specify and manipulate array data structures.
  - **SciPy**: An open source library of scientific tools for Python, SciPy supplements the NumPy module. SciPy includes modules for graphics and plotting, optimization, integration, special functions, signal and image processing, genetic algorithms, ODE solvers, and others.
  - **Matplotlib**: Matplotlib is a python 2D plotting library which produces publication-quality figures.
  - **pandas**: A module that provides high-performance, easy-to-use data structures and data analysis tools. In particular, the DataFrame class is useful for spreadsheet-like representation and mannipulation of data.
  - **IPython**: An enhanced Python shell, designed to increase the efficiency and usability of coding, testing and debugging Python.

# Why Python?

- Strong points:
  - Viable alternative to SAS/Matlab/R/...
  - Uptake in the financial sector.
  - Upcoming generation of programmers with Python as a first language.
  - Python communities in HPC/Scientific Computing.

- Weak points
  - No (public) great success story.
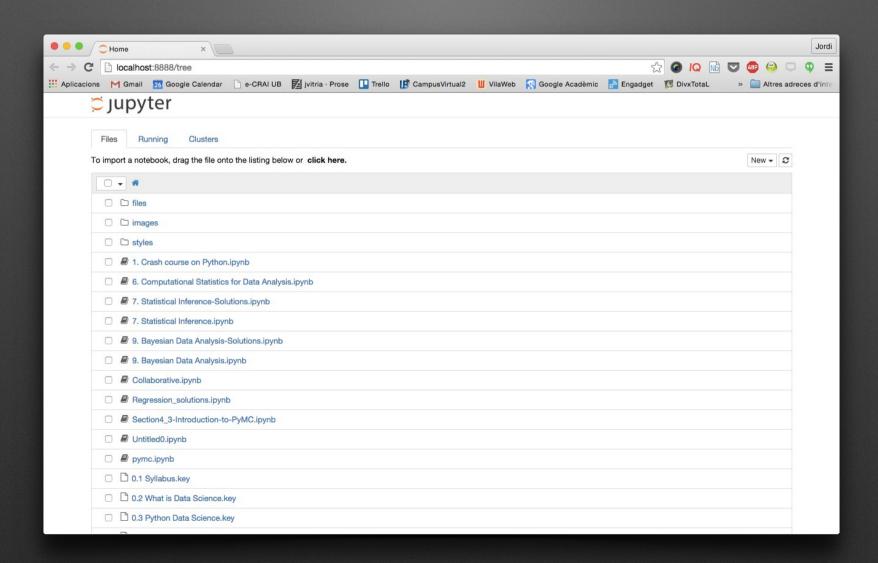  - Weak support.

# What is IPython?

I is for interactive

In this class, we will use **IPython notebooks** for teaching and programming assignments.

An IPython notebook lets you write and execute Python code in your web browser.

IPython notebooks make it very easy to tinker with code and execute it in bits and pieces; for this reason IPython notebooks are widely used in scientific computing.

# What is IPython?

I is for interactive

# What is IPython?

I is for interactive