

HERRAMIENTAS PARA CIENCIA DE DATOS

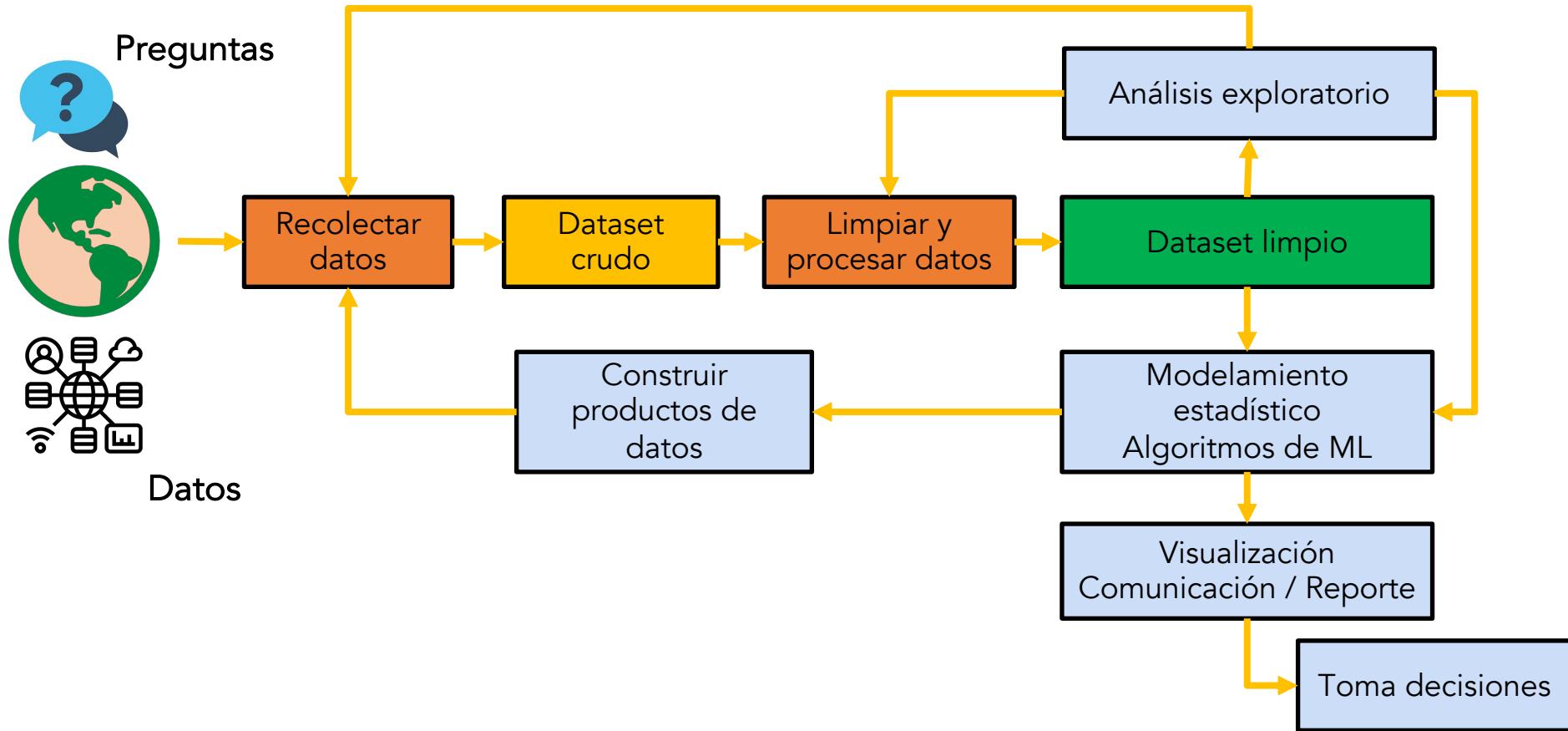
CLASE 2

AGENDA DE HOY

- Repaso Clase 1
- Herramientas computacionales para CD
- Demostración en Jupyter Notebook
- Actividad 1

¿Dudas respecto al
programa?

EL PROCESO DE CIENCIA DE DATOS



DATOS vs. INFORMACIÓN

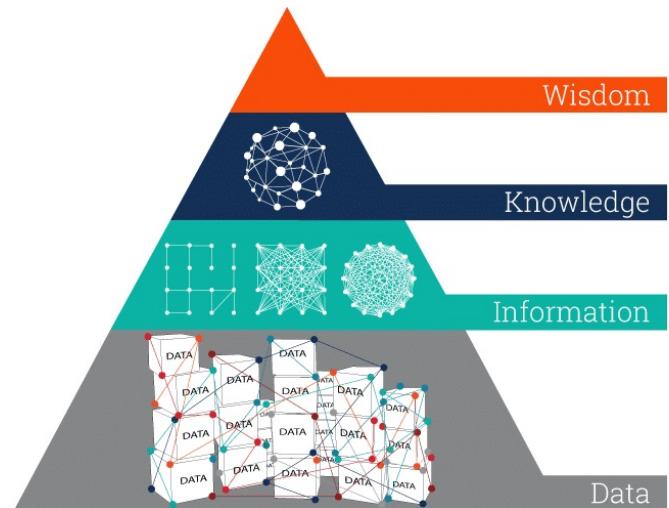
Datos → abstracciones o mediciones obtenidas del mundo real.

Información → datos que han sido procesados, estructurados o contextualizados de manera que son significativos para las personas.

Conocimiento → información que ha sido interpretada y comprendida por una persona, de manera que puede actuar a partir de ella.

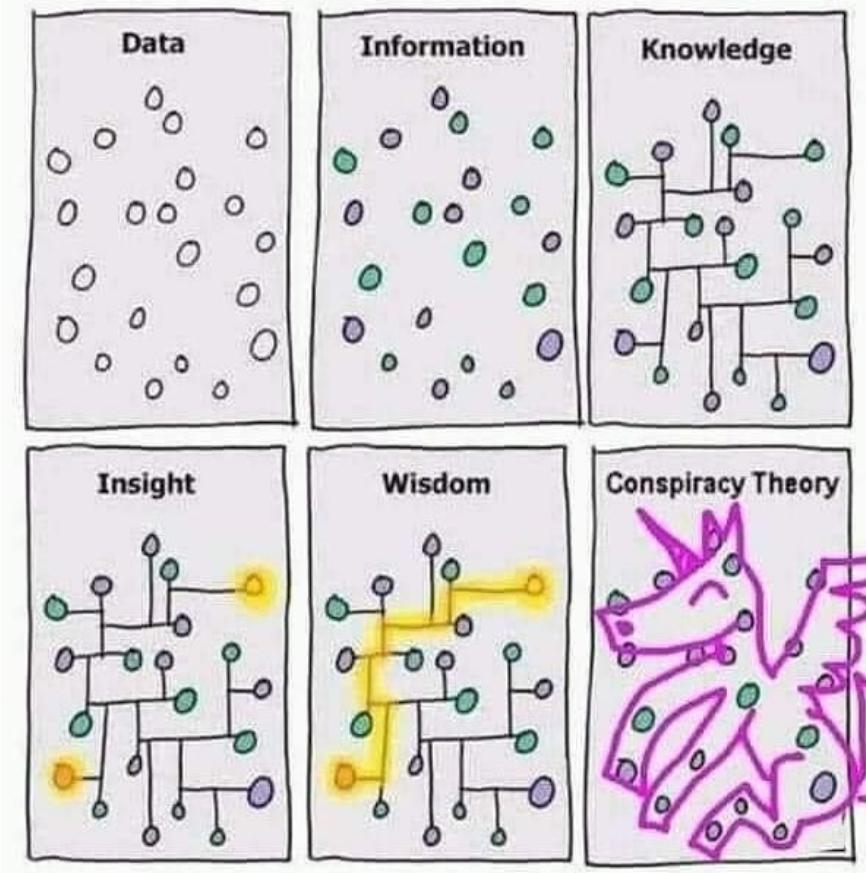
Sabiduría → actuar de forma apropiada a partir del conocimiento

Pirámide DIKW



<https://www.ontotext.com/knowledgehub/fundamentals/dikw-pyramid/>

DATOS vs. INFORMACIÓN



"The combination of some data and an aching desire for an answer does not ensure that a reasonable answer can be extracted from a given body of data"

-John Tukey

"La combinación de algunos datos y un intenso deseo de obtener una respuesta no asegura que sea posible extraer una respuesta razonable de un conjunto dado de datos"

Herramientas Computacionales para Ciencia de Datos

REPRODUCIBILIDAD CIENTÍFICA

Reproducibilidad:

- Capacidad de reproducir procedimientos científicos (experimentos, análisis, resultados).
- Clave para el desarrollo, colaboración y transparencia científica y profesional

¿Por qué es importante?

- Para repetir un análisis con distintos parámetros.
- Para escalar un análisis → entrega a equipo de producción
- Para crear confianza en los resultados
- Para compartir el conocimiento
- Para ahorrar tiempo y stress 😊

nature <https://www.nature.com/articles/533452a>

Explore content ▾ Journal information ▾ Publish with us ▾

nature > news feature > article

Published: 25 May 2016

1,500 scientists lift the lid on reproducibility

Monya Baker

Nature 533, 452–454 (2016) | Cite this article
12k Accesses | 1301 Citations | 3886 Altmetric | Metrics

This article has been updated

Survey sheds light on the ‘crisis’ rocking research.

More than 70% of researchers have tried and failed to reproduce another scientist's experiments, and more than half have failed to reproduce their own experiments. Those are some of the telling figures that emerged from *Nature*'s survey of 1,576 researchers who took a brief online questionnaire on reproducibility in research.

Digital innovators from Alibaba to Google to Facebook to Netflix to Microsoft to Amazon have already actively embraced large-scale, rapid experimentation as integral to their innovation efforts. But every organization seeking authentic insight from experiment needs to be wary of the problems that have made scientific research unreliable.

<https://towardsdatascience.com/why-your-company-needs-reproducible-research-d4a08f978d39>

REPRODUCIBILIDAD CIENTÍFICA

¿Cómo lo hacemos?

- Disponibilización de **datos**.
- Automatización de procesos, sin pasos “manuales”.
- Desarrollando procesos computacionales transparentes.
- Con suficiente **detailed** para reproducir todo el análisis.
- Usando modelos y herramientas eficientes de **reusabilidad**.
 - Software libre
 - Paquetes de código abierto
 - “Notebooks” computacionales
 - Repositorios y versionamiento

nature <https://www.nature.com/articles/533452a>

Explore content ▾ Journal information ▾ Publish with us ▾

nature > news feature > article

Published: 25 May 2016

1,500 scientists lift the lid on reproducibility

Monya Baker

Nature 533, 452–454 (2016) | Cite this article

12k Accesses | 1301 Citations | 3886 Altmetric | Metrics

 This article has been updated

Survey sheds light on the ‘crisis’ rocking research.

More than 70% of researchers have tried and failed to reproduce another scientist’s experiments, and more than half have failed to reproduce their own experiments. Those are some of the telling figures that emerged from *Nature*’s survey of 1,576 researchers who took a brief online questionnaire on reproducibility in research.

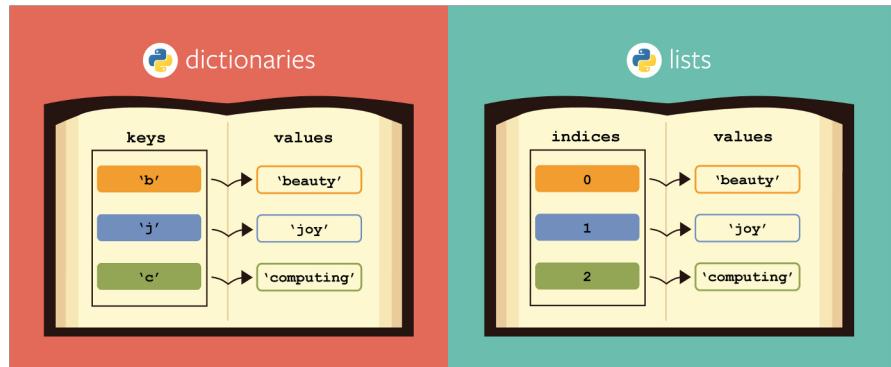
Digital innovators from Alibaba to Google to Facebook to Netflix to Microsoft to Amazon have already actively embraced large-scale, rapid experimentation as integral to their innovation efforts. But every organization seeking authentic insight from experiment needs to be wary of the problems that have made scientific research unreliable.

<https://towardsdatascience.com/why-your-company-needs-reproducible-research-d4a08f978d39>

- Lenguaje **interpretado** de programación: los comandos o instrucciones se ejecutan paso por paso, sin compilación previa.
- Gran y activa **comunidad** de computación científica y ciencia de datos desarrollada en torno a Python.
- Facilidad para **integrar** códigos desarrollados en C, C++, FORTRAN, etc. (más eficientes para algunos algoritmos de cálculo)
- Integración de **múltiples librerías** con distintas funcionalidades y usos (numpy, sciPy, pandas, scikit-learn, etc).
- Menos eficiente que otros lenguajes para computación paralela y concurrencias.

```
>>> import this  
The Zen of Python, by Tim Peters
```

Beautiful is better than ugly.
Explicit is better than implicit.
Simple is better than complex.
Complex is better than complicated.
Flat is better than nested.
Sparse is better than dense.
Readability counts.
Special cases aren't special enough to break the rules.
Although practicality beats purity.
Errors should never pass silently.
Unless explicitly silenced.
In the face of ambiguity, refuse the temptation to guess.
There should be one-- and preferably only one --obvious way to do it.
Although that way may not be obvious at first unless you're Dutch.
Now is better than never.
Although never is often better than *right* now.
If the implementation is hard to explain, it's a bad idea.
If the implementation is easy to explain, it may be a good idea.
Namespaces are one honking great idea -- let's do more of those!



PYTHON

Variables

You can name a variable anything as long as it obeys the following rules:

1. It can be only one word.
2. It can use only letters, numbers, and the underscore (`_`) character.
3. It can't begin with a number.
4. Variable name starting with an underscore (`_`) are considered as "unuseful".

Data Types

Data Type	Examples
Integers	-2, -1, 0, 1, 2, 3, 4, 5
Floating-point numbers	-1.25, -1.0, --0.5, 0.0, 0.5, 1.0, 1.25
Strings	'a', 'aa', 'aaa', 'Hello!', '11 cats'

Comparison Operators

Operator	Meaning
<code>==</code>	Equal to
<code>!=</code>	Not equal to
<code><</code>	Less than
<code>></code>	Greater Than
<code><=</code>	Less than or Equal to
<code>>=</code>	Greater than or Equal to

Functions

```
>>> def hello(name):
...     print('Hello {}'.format(name))
...
>>> hello('Alice')
Hello Alice
Hello Bob
```

if Statements

```
if name == 'Alice':
    print('Hi, Alice.')
```

else Statements

```
name = 'Bob'
if name == 'Alice':
    print('Hi, Alice.')
else:
    print('Hello, stranger.')
```

elif Statements

```
name = 'Bob'
age = 5
if name == 'Alice':
    print('Hi, Alice.')
elif age < 12:
    print('You are not Alice, kiddo.')
```

Dictionaries and Structuring Data

Example Dictionary:

```
myCat = {'size': 'fat', 'color': 'gray', 'disposition': 'loud'}
```

The keys(), values(), and items() Methods

values():

```
>>> spam = {'color': 'red', 'age': 42}
>>> for v in spam.values():
...     print(v)
red
42
```

keys():

```
>>> for k in spam.keys():
...     print(k)
color
age
```

Lists

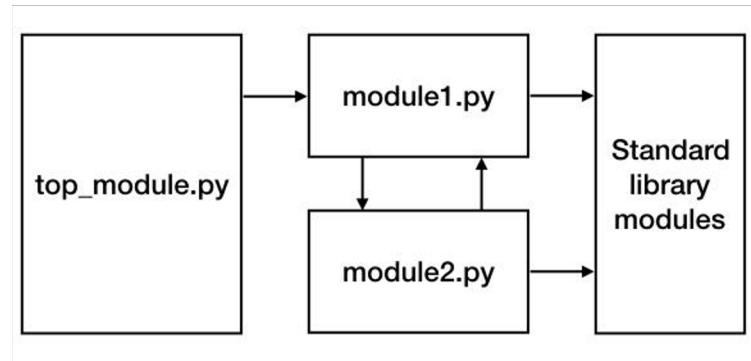
```
>>> spam = ['cat', 'bat', 'rat', 'elephant']
>>> spam
['cat', 'bat', 'rat', 'elephant']
```

Getting Individual Values in a List with Indexes

```
>>> spam = ['cat', 'bat', 'rat', 'elephant']
>>> spam[0]
'cat'
```

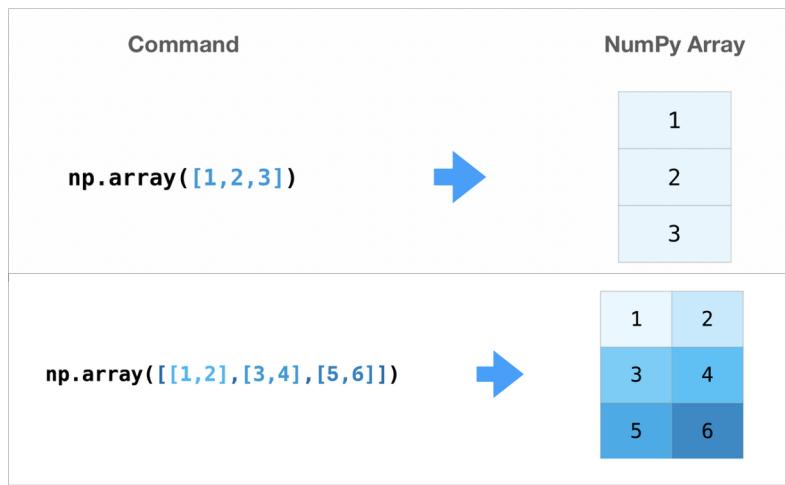
LIBRERÍAS EN PYTHON

- **Script:** archivo de Python que es ejecutado directamente para hacer “algo”.
- **Módulo:** archivo Python que es importado por scripts u otros módulos. Define clases, funciones y variables a ser usados por otros módulos.
- **Paquete:** colección de módulos relacionados que proveen una cierta funcionalidad.
- **Librería:** colección de códigos reutilizable. Puede estar formadas por muchos módulos o paquetes.
- **Python Package Index (PyPi):** repositorio de software en lenguaje Python.
 - Permite encontrar e instalar software desarrollado y compartido por la comunidad Python.



NumPy <https://numpy.org/>

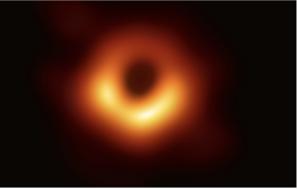
- Provee las estructuras, algoritmos y funciones para la mayoría de las aplicaciones científicas relacionadas con datos numéricos.
- **Arreglos** multi-dimensionales (ndarray): eficientes para almacenar y manipular data.
- Funciones para realizar **operaciones matemáticas** entre **arreglos** ➔ evita iteraciones
 - Álgebra lineal, transformadas de Fourier, números aleatorios, etc.
- **Interoperable**: estructura de data primaria para muchas otras herramientas de Python.
- **Eficiente**: basada en códigos en C bien optimizados.
- Herramientas para leer y escribir datos/arreglos al disco/archivos.



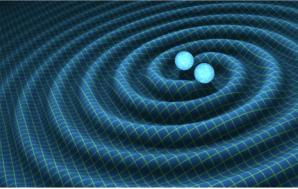
data	data[0,1]	data[1:3]	data[0:2,0]
0 1 2	0 1 2	0 1 2	0 1 2
1 3 4	1 3 4	1 3 4	1 3 4
2 5 6	2 5 6	2 5 6	2 5 6

CASE STUDIES

FIRST IMAGE OF A BLACK HOLE

A black hole image showing a bright, glowing orange and yellow ring against a dark background.

DETECTION OF GRAVITATIONAL WAVES

A visualization of gravitational waves as concentric blue and green ripples emanating from two blue spheres.

SPORTS ANALYTICS

A red cricket ball resting on a green grassy field.

How NumPy, together with libraries like SciPy and Matplotlib that depend on NumPy, enabled the Event Horizon Telescope to produce the first ever image of a black hole

In 1916, Albert Einstein predicted gravitational waves; 100 years later their existence was confirmed by LIGO scientists using NumPy.

Cricket Analytics is changing the game by improving player and team performance through statistical modelling and predictive analytics. NumPy enables many of these analyses.

POSE ESTIMATION USING DEEP LEARNING

A cheetah in a grassy field with small colored dots (red, green, blue) placed on its body at various joints to track its movement.

DeepLabCut uses NumPy for accelerating scientific studies that involve observing animal behavior for better understanding of motor control, across species and timescales.

NumPy



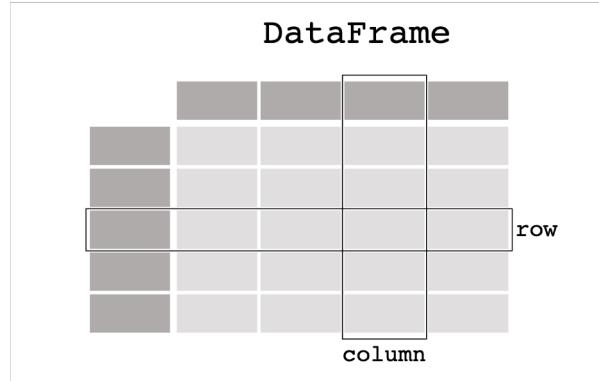
The fundamental package for scientific computing with Python

PANDAS

<https://pandas.pydata.org/>

<https://pandas.pydata.org/>

- Provee las estructuras y funciones de alto nivel para trabajar con **datos estructurados y tablas** en forma rápida, fácil y expresiva.
 - **DataFrame:** estructura de datos tabular, orientada a columnas, con etiquetas de filas y columnas.
 - **Series:** objeto tipo arreglo 1-D, con etiqueta.
-
- Combina la eficiencia de cálculo de NumPy, con capacidades de manipulación de datos típicas de hojas de cálculo (e.g. Excel) y bases de datos relacionales.
 - Provee funcionalidades de indexación para reformatear (reshape), seccionar (slice), agregar, filtrar, seleccionar subconjuntos de datos.



Each column in a **DataFrame** is a **Series**

Series

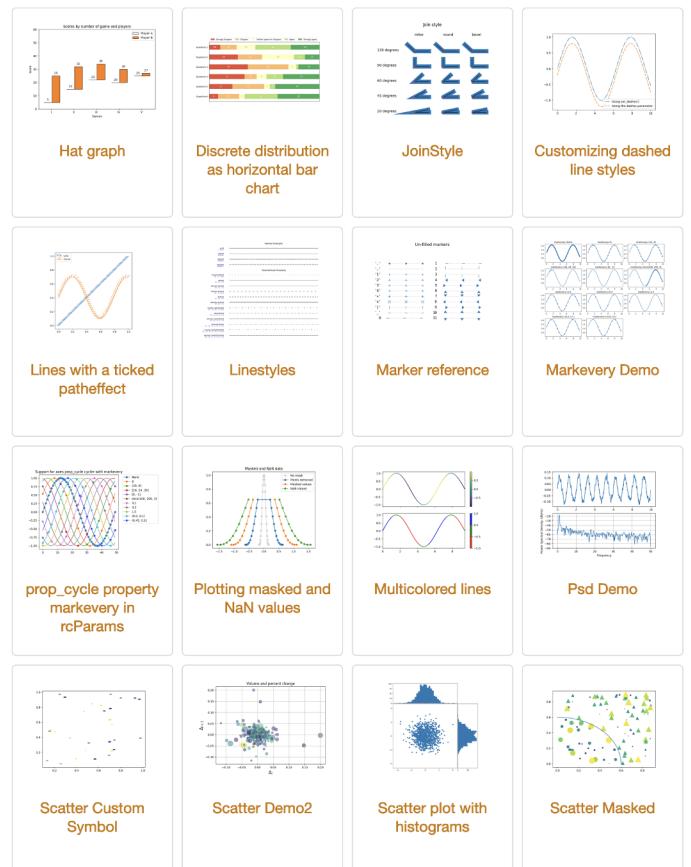


MATPLOTLIB <https://matplotlib.org>



- Librería de Python más popular para producción de gráficos y otras visualizaciones 2D/3D estáticas, animadas e interactivas.
- Permite controlar y personalizar todos los detalles de una figura.
- Existen varios paquetes que extienden y construyen sobre las funcionalidades de Matplotlib, incluyendo algunas interfaces de ploteo de alto nivel:
 - Seaborn: interfaz para generación de gráficos estadísticos atractivos e informativos. Ayuda a explorar y entender la data.
 - Cartopy: paquete diseñado para el procesamiento de data geoespacial, para producir mapas y otros análisis espaciales
- Galerías:
 - <https://matplotlib.org/stable/gallery/index.html>
 - <https://www.python-graph-gallery.com/>
- Cheatsheets:

<https://github.com/matplotlib/cheatsheets#cheatsheets>





Version 3.2

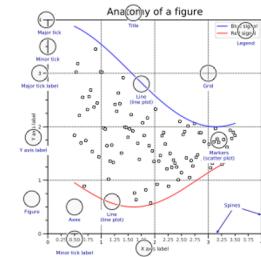
Quick start

```
import numpy as np
import matplotlib as mpl
import matplotlib.pyplot as plt

X = np.linspace(0, 2*np.pi, 100)
Y = np.cos(X)

fig, ax = plt.subplots()
ax.plot(X,Y,color='C1')
fig.savefig("figure.pdf")
fig.show()
```

Anatomy of a figure



Subplots layout

```
subplots([cols,rows,...])
fig, axs = plt.subplots(3,3)

G = gridspec([cols,rows,...])
ax = G[0,:]

ax.inset_axes(extent)

d=make_axes_locatable(ax)
ax=d.new_horizontal('10%')
```

Getting help

- matplotlib.org
- github.com/matplotlib/matplotlib/issues
- discourse.matplotlib.org
- stackoverflow.com/questions/tagged/matplotlib
- [gitter.im/matplotlib](https://gitter.im/matplotlib/lounge)
- twitter.com/matplotlib
- Matplotlib users mailing list

Basic plots

```
plot([X],Y,[fmt],...)
X, Y, fmt, color, marker, linestyle

scatter(X,Y,...)
X, Y, [s]izes, [c]olors, markers, cmap

bar[h](x,height,...)
x, height, width, bottom, align, color

imshow(Z,[cmap],...)
Z, cmap, interpolation, extent, origin

contourf([X],[Y],Z,...)
X, Y, Z, levels, colors, extent, origin

quiver([X],[Y],U,V,...)
X, Y, U, V, C, units, angles
```

```
pie(x,[explode],...)
Z, explode, labels, colors, radius

text(x,y,text,...)
x, y, text, va, ha, size, weight, transform

fill_[between](x)[...]
X, Y1, Y2, color, where
```

Advanced plots

```
step(X,Y,[fmt],...)
X, Y, fmt, color, marker, where

boxplot(X,...)
X, notch, sym, bootstrap, widths

errorbar(X,Y,xerr,yerr,...)
X, Y, xerr, yerr, fmt

hist(X, bins, ...)
X, bins, range, density, weights

violinplot(D,...)
D, positions, widths, vert

barbs([X],[Y], U, V, ...)
X, Y, U, V, C, length, pivot, sizes

eventplot(positions,...)
positions, orientation, lineoffsets

hexbin(X,Y,C,...)
X, Y, C, gridsize, bins

xcorr(X,Y,...)
X, Y, normed, detrend
```

Scales

linear	any values	log
z	z	$values > 0$

Projections

```
subplot(...,projection=p)
p='polar'
p='3d'
p=Orthographic()
from cartopy.crs import Cartographic
```

Lines

linestyle or ls

 capstyle or dash_capstyle

Markers

																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																											<img alt="triangle-right-down-left-up-right-down-down
--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	---

SCIPY <https://www.scipy.org/>

- Ecosistema de software de código abierto basado en Python para matemáticas, ciencia e ingeniería.
- La *librería SciPy* provee una colección de paquetes enfocados en distintos problemas o requerimientos propios de computación científica:
 - Integración (`scipy.integrate`)
 - Álgebra lineal (`scipy.linalg`)
 - Optimización (`scipy.optimize`)
 - Análisis de señales (`scipy.signal`)
 - Estadísticas (`scipy.stats`)
 - Interpolación (`scipy.interpolate`)
- Usa arreglos de numpy como estructura básica de datos.

SciPy.org



SciPy (pronounced "Sigh Pie") is a Python-based ecosystem of open-source software for mathematics, science, and engineering. In particular, these are some of the core packages:

 NumPy Base N-dimensional array package	 SciPy library Fundamental library for scientific computing	 Matplotlib Comprehensive 2-D plotting
 IPython Enhanced interactive console	 SymPy Symbolic mathematics	 pandas Data structures & analysis
 Large parts of the SciPy ecosystem (including all six projects above) are fiscally sponsored by NUMFOCUS		

SCIKIT LEARN <https://scikit-learn.org/stable/>



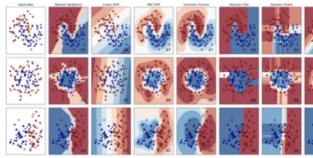
- Construido sobre NumPy, SciPy y matplotlib.
- Provee herramientas simples y eficientes para **aprendizaje de máquinas**:
 - Preprocesamiento
 - Clasificación
 - Regresión
 - Clustering
 - Reducción de dimensionalidad
 - Selección de modelos

Classification

Identifying which category an object belongs to.

Applications: Spam detection, image recognition.

Algorithms: SVM, nearest neighbors, random forest, and more...



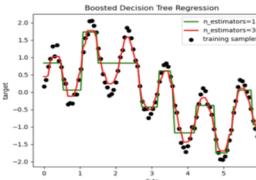
Examples

Regression

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, Stock prices.

Algorithms: SVR, nearest neighbors, random forest, and more...



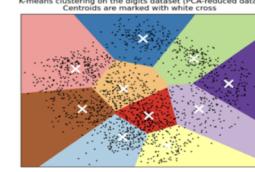
Examples

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: k-Means, spectral clustering, mean-shift, and more...



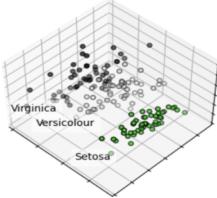
Examples

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Algorithms: k-Means, feature selection, non-negative matrix factorization, and more...



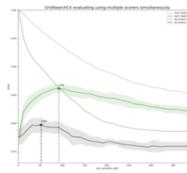
Examples

Model selection

Comparing, validating and choosing parameters and models.

Applications: Improved accuracy via parameter tuning

Algorithms: grid search, cross validation, metrics, and more...



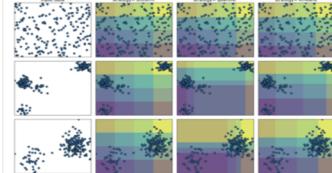
Examples

Preprocessing

Feature extraction and normalization.

Applications: Transforming input data such as text for use with machine learning algorithms.

Algorithms: preprocessing, feature extraction, and more...



Examples

GEOPANDAS



<https://geopandas.org/>

- Librería que facilita el trabajo con datos geoespaciales en Python.
- Extiende los tipos de datos usados por pandas para permitir manipulación y operaciones espaciales con datos geométricos.
 - DataFrame → GeoDataFrame
 - Series → GeoSeries
- Se basa en otras librerías pre-existentes:
 - **fiona** → acceso a formatos de datos geoespaciales
 - **shapely** → operaciones geométricas

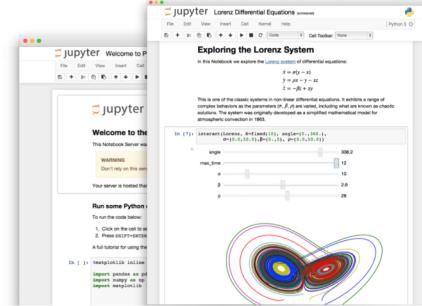
[27]: data.head()

[27]:

	GROUP	CLASS	geometry	area
0	64	32421	POLYGON ((379394.248 6689991.936, 379389.790 6...	76.027392
1	64	32421	POLYGON ((378980.811 6689359.377, 378983.401 6...	2652.054186
2	64	32421	POLYGON ((378804.766 6689256.471, 378817.107 6...	3185.649995
3	64	32421	POLYGON ((379229.695 6685025.111, 379233.366 6...	13075.165279
4	64	32421	POLYGON ((379825.199 6685096.247, 379829.651 6...	3980.682621

JUPYTER NOTEBOOKS

- Aplicación web open source para crear y compartir documentos (**.ipynb**) con:
 - Código ejecutable
 - Resultados del código
 - Texto narrativo e imágenes
- Soporta múltiples lenguajes de programación.
- Notebooks pueden ser exportados a múltiples formatos estáticos.
- Instalación: <https://jupyter.org/install.html>
- Documentación:
https://jupyter-notebook.readthedocs.io/en/stable/examples/Notebook/examples_index.html



The Jupyter Notebook is an open-source web application that allows you to create and share documents that contain live code, equations, visualizations and narrative text. Uses include: data cleaning and transformation, numerical simulation, statistical modeling, data visualization, machine learning, and much more.

[Try it in your browser](#) [Install the Notebook](#)

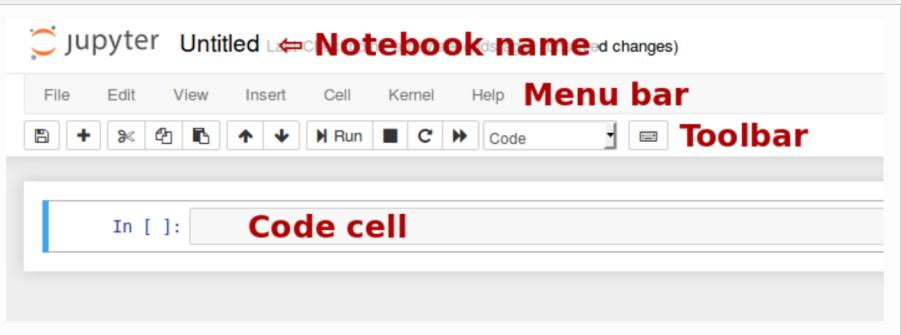
</> Language of choice Share notebooks Interactive output Big data integration

Jupyter supports over 40 programming languages, including Python, R, Julia, and Scala.

Notebooks can be shared with others using email, Dropbox, GitHub and the [Jupyter Notebook Viewer](#).

Your code can produce rich, interactive output: HTML, images, videos, LaTeX, and custom MIME types.

Leverage big data tools, such as Apache Spark, from Python, R and Scala. Explore that same data with pandas, scikit-learn, ggplot2, TensorFlow.



jupyter Untitled Notebook [name](#) [changes](#)

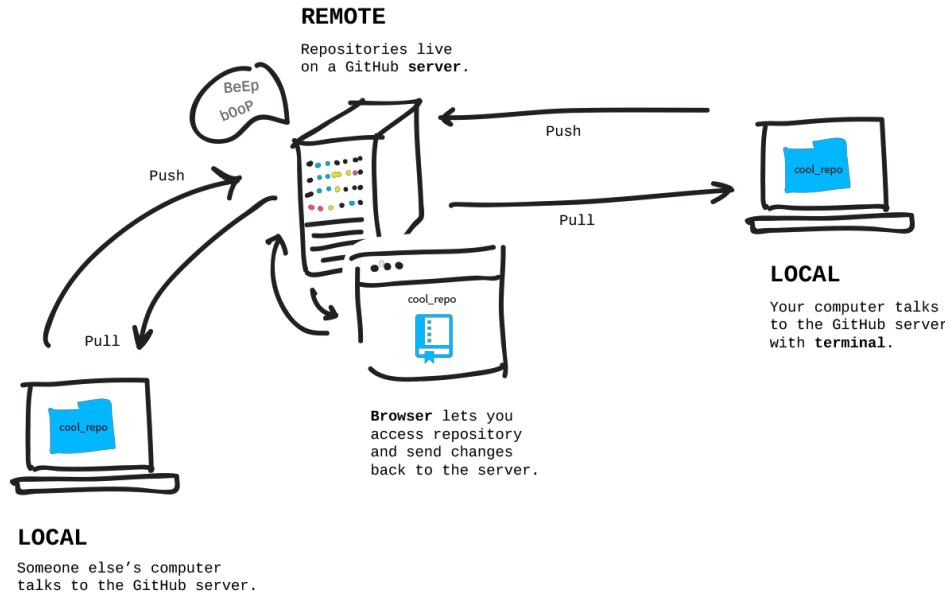
File Edit View Insert Cell Kernel Help **Menu bar**

Toolbar

In []: **Code cell**

GITHUB

- GitHub es una plataforma de hospedaje de código para el **control de versiones** y la **colaboración**.
- Permite que personas trabajen juntos en proyectos desde donde sea.
- Cada proyecto es un **repositorio**, que todos los **archivos** de tu proyecto y el **historial de revisiones** de cada uno de ellos.
 - Puedes debatir y administrar el trabajo de tu proyecto dentro del repositorio.



<https://docs.github.com/es/get-started/quickstart/hello-world>

GITHUB y PORTFOLIOS

How to Build a Data Science Portfolio



Michael Galarnyk Jul 8, 2018 · 18 min read



A portfolio is one way to show people you are that data science unicorn.

<https://towardsdatascience.com/how-to-build-a-data-science-portfolio-5f566517c79c>

These data science portfolios will awe and inspire you (mid-2020 edition)

Use these to improve your own data science portfolio, learn new skills or discover new, interesting projects.

<https://towardsdatascience.com/these-data-science-portfolios-will-awe-and-inspire-you-mid-2020-edition-728e1021f60>

The Promise of Portfolios: Training Modern Data Scientists

by Deborah Nolan and Sara Stoudt

Published on Jul 30, 2021

last released
3 weeks ago

ABSTRACT

Data scientists rely on many technical skills and the ability to reason about data to solve problems. As educators grapple with how to prepare students in this new field, they are faced with identifying both what a student must know and what a student should be able to do by the end of their data science education, and also how to collect evidence of those abilities.

We present a way to unite and coordinate individual efforts toward training well-rounded data scientists: a data science portfolio that highlights strong communication. This structuring of classroom assignments provides a way to evaluate students' mastery of material in each class and also allows for a student to build a professional portfolio that remains valuable after the class is over.

Data science portfolio pieces broadly include written and visual assignments that give students practice crafting data-driven arguments and narratives for a variety of audiences. This flexible nature of the portfolio gives students a way to demonstrate their abilities in an inclusive, 'choose your own adventure' way.

<https://hdsr.mitpress.mit.edu/pub/pm67plq9/release/2>

Git PUSH PULL

