

# **Curso Big Data**

**Taller de Empleo  
Camargo X**

**Biuse Casaponsa  
[deducedatasolutions.com](http://deducedatasolutions.com)**

“Big data” hace referencia a conjuntos de datos tan grandes y complejos como para que hagan falta aplicaciones informáticas no tradicionales de procesamiento de datos para tratarlos adecuadamente. *Wikipedia*



# 'DATA IS THE NEW OIL.'

From the beginning of recorded time until 2003, we created

**5 exabytes  
of data.**

In 2011 the same amount was created every two days.

By 2013, it's expected that the time will shrink to 10 minutes.

Every hour, we create enough Internet traffic to fill

**7 billion  
DVDs.**

Side by side, that's seven times the height of Everest.

Coined in 2006 by Clive Humby, a British data commercialization entrepreneur, this now famous phrase was embraced by the World Economic Forum in a 2011 report, which considered data to be an economic asset, like oil.

English is the dominant language of the web. But by 2014 it will be **Chinese**, if its current rate of increase continues.

Top languages used on the web (May 2011):



**247 billion  
EMAILS**  
are sent **every day**. (Up to 80% are spam.)

There are nearly as many bits of information in the digital universe as there are stars in our actual universe.

As of August 2012, there were just over **4 million** articles in the English Wikipedia.

There are **133 million BLOGS** on the web.

on the web.

Just as a study of activity on Twitter gave residents, family members, and journalists advance warning of details about the devastating earthquake and tsunami in Japan,

**high-frequency traders**, with the help of computer algorithms, use Big Data to follow trends and to act quickly on their findings.

These specialized algorithms

make split-second decisions to buy or sell a commodity.

New cable being laid under the Atlantic will shave

**5 milliseconds**

from the current 65 milliseconds it takes for trading firms to travel between New York and London.

With new fiber-optic cable,

the round-trip time between New York and London will be 59.6 milliseconds.

This 5-millisecond saving is worth many millions of dollars to the trading firms who use the cable (and who will pay millions to do so).

How they save 5 milliseconds:

The depth of the Atlantic Ocean varies from 1,000 feet to 13,000 feet. The new cable will lie on areas of the ocean floor that are up to 1,000 feet shallower than the current fastest cable. A different route, the new cable is shorter, meaning that the time it takes for messages to travel along it is shortened.

The new cable takes a shallower, therefore shorter route.

USA

50%  
of 5-year-old  
kids in the U.S. are given  
access to a smartphone.



# 'DATA IS THE NEW OIL.'

From the beginning of recorded time until 2003, we created

**5 exabytes** (5 billion gigabytes) of data.

In 2011 the same amount was created every two days.

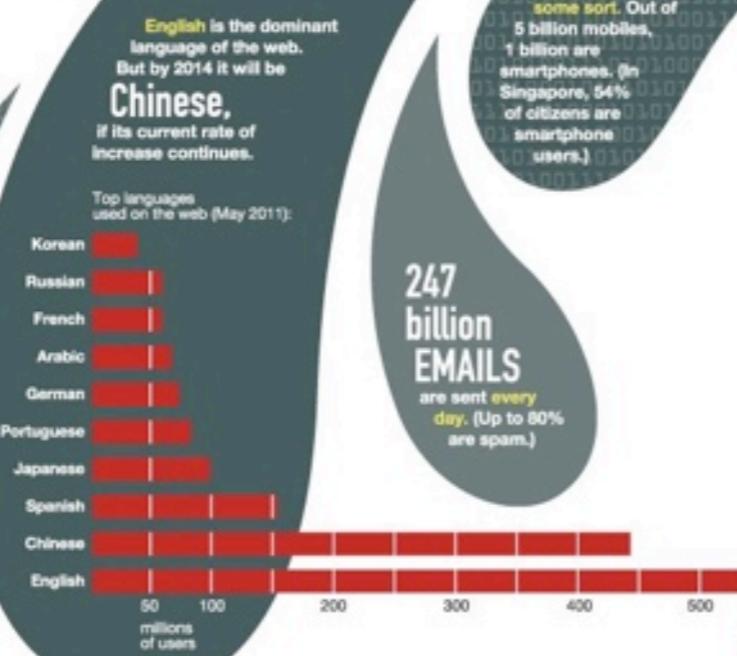
By 2013, it's expected that the time will shrink to 10 minutes.

Every hour, we create enough Internet traffic to fill  
**7 billion DVDs.**

Side by side, that's seven times the height of Everest.

Coined in 2006 by Clive Humby, a British data commerce entrepreneur, this now famous phrase was embraced by the World Economic Forum in a 2011 report, which considered data to be an economic asset, like oil.

## 1 exabyte = 1 000 millones de GB



**247 billion EMAILS** are sent every day. (Up to 80% are spam.)

80% of all humans own a mobile phone of some sort. Out of 5 billion mobiles, 1 billion are smartphones. (In Singapore, 54% of citizens are smartphone users.)

Just as a study of activity on Twitter gave residents, family members, and journalists advance warning of details about the devastating earthquake and tsunami in Japan, **high-frequency traders**, with the help of computer algorithms, use Big Data to follow trends and to act quickly on their findings.

These specialized algorithms make split-second decisions to buy or sell a commodity.

New cable being laid under the Atlantic will shave

**5 milliseconds**

from the current 65 milliseconds it takes for trading firms to travel between New York and London.

With new fiber-optic cable, the round-trip time between New York and London will be 59.6 milliseconds.

This 5-millisecond saving is worth many millions of dollars to the trading firms who use the cable (and who will pay millions to do so).

How they save 5 milliseconds:

The depth of the Atlantic Ocean varies from 1,000 feet to 13,000 feet. The new cable will lie on areas of the ocean floor that are up to 1,000 feet shallower than the current fastest cable. A different route, the new cable is shorter, meaning that the time it takes for messages to travel along it is shortened.

The new cable takes a shallower, therefore shorter route.

USA

50% of 5-year-old kids in the U.S. are given access to a smartphone.



# 'DATA IS THE NEW OIL.'

From the beginning of recorded time until 2003, we created

**5 exabytes** (5 billion gigabytes) of data.

In 2011 the same amount was created every two days.

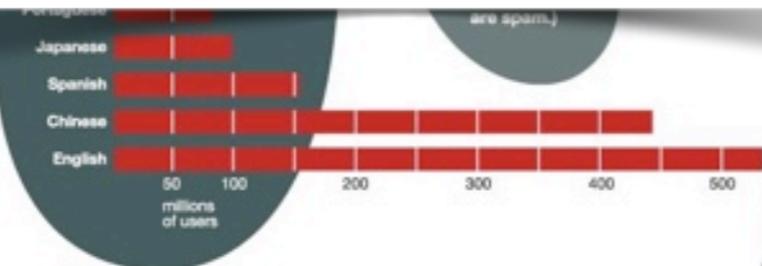
By 2013, it's expected that the time will shrink to 10 minutes.

Every hour, we create enough Internet traffic to fill  
**7 billion DVDs.**

Side by side, that's seven times the height of Everest.

Coined in 2006 by Clive Humby, a British data commerce entrepreneur, this now famous phrase was embraced by the World Economic Forum in a 2011 report, which considers data to be an economic asset, like oil.

## 1000 millones de GB = 2 millones de ordenadores



## 1 exabyte = 1 000 millones de GB

Just as a study of activity on Twitter gave residents, family members, and journalists advance warning of details about the devastating earthquake and tsunami in Japan,

**high-frequency traders,** with the help of computer algorithms, use Big Data to follow trends and to act quickly on their findings.

These specialized algorithms make split-second decisions to buy or sell a commodity.

New cable being laid under the Atlantic will shave

**5 milliseconds**

from the current 65 milliseconds it takes for traffic to travel between New York and London.

With new fiber-optic cable, the round-trip time between New York and London will be 59.6 milliseconds.

This 5-millisecond saving is worth many millions of dollars to the trading firms who use the cable (and who will pay millions to do so).

How they save 5 milliseconds

The depth of the Atlantic Ocean varies from 1,000 feet to 36,000 feet. The new cable will lie on areas of the ocean floor that are up to 1,000 feet shallower than the current fastest cable. A different route, the new cable is shorter, meaning that the time it takes for messages to travel along it is shortened.

The new cable takes a shallower, therefore shorter route.

50% of 5-year-old kids in the U.S. are given access to a smartphone.



# 'DATA IS THE NEW OIL.'

From the beginning of recorded time until 2003, we created

**5 exabytes** (5 billion gigabytes) of data.

In 2011 the same amount was created every two days.

By 2013, it's expected that the time will shrink to 10 minutes.

Every hour, we create enough Internet traffic to fill  
**7 billion DVDs.**

Side by side, that's seven times the height of Everest.

1 exabyte =  
1 000 millones de GB

1 000 millones de GB = 2 millones de ordenadores

generados 5 exabytes hasta 2003

Coined in 2006 by Clive Humby, a British data commerce entrepreneur, this now famous phrase was embraced by the World Economic Forum in a 2011 report, which considers data to be an economic asset, like oil.

Just as a study of activity on Twitter gave residents, family members, and journalists advance warning of details about the devastating earthquake and tsunami in Japan,

**high-frequency traders,** with the help of computer algorithms, use Big Data to follow trends and to act quickly on their findings.

These specialized algorithms make split-second decisions to buy or sell a commodity.

New cable being laid under the Atlantic will shave

**5 milliseconds**

from the current 65 milliseconds it takes for traffic to travel between New York and London.

With new fiber-optic cable,

the round-trip time between New York and London will be 59.6 milliseconds.

This 5-millisecond saving is worth many millions of dollars to the trading firms who use the cable (and who will pay millions to do so).

How they save 5 milliseconds

The depth of the Atlantic Ocean varies from 1,000 feet to 13,000 feet. The new cable will lie on areas of the ocean floor that are up to 1,000 feet shallower than the current fastest cable. A different route, the new cable is shorter, meaning that the time it takes for messages to travel along it is shortened.

The new cable takes a shallower, therefore shorter route.

USA

# "DATA IS THE NEW OIL."

From the beginning of recorded time until 2003, we created

**5 exabytes** (5 billion gigabytes) of data.

In 2011 the same amount was created every two days.

By 2013, it's expected that the time will shrink to 10 minutes.

Every hour, we create enough Internet traffic to fill  
**7 billion DVDs.**

Side by side, that's seven times the height of Everest.

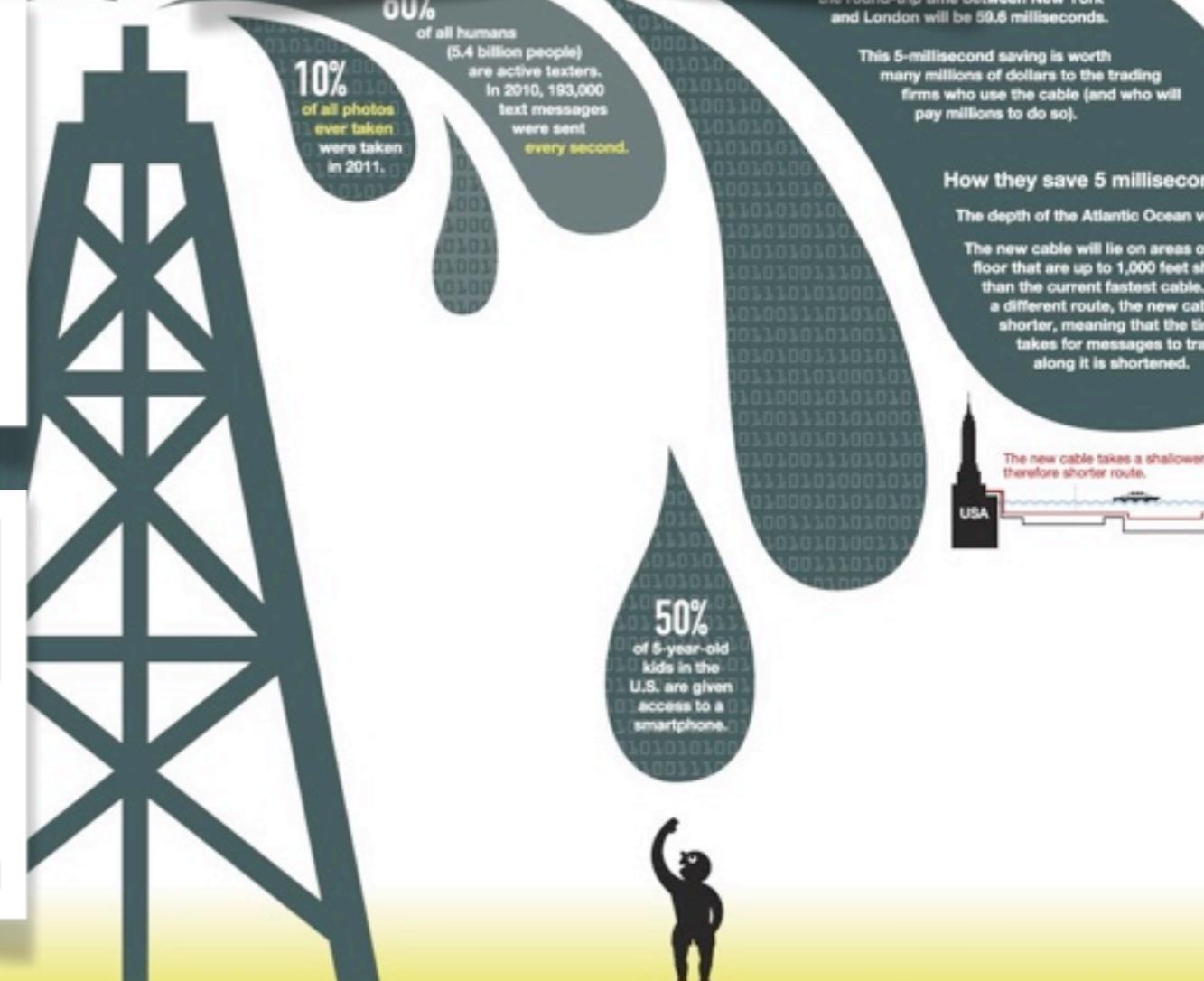
Coined in 2006 by Clive Humby, a British data commerce entrepreneur, this now famous phrase was embraced by the World Economic Forum in a 2011 report, which considers data to be an economic asset, like oil.

1000 millones de GB = 2 millones de ordenadores

generados 5 exabytes hasta 2003

1 exabyte =  
1 000 millones de GB

2011 cada 2 días se generaban 5 exabytes



**“DATA  
IS THE  
NEW  
OIL.”**

From the beginning of recorded time until 2003, we created

**5 exabytes** (5 billion gigabytes)  
of data.

In 2011 the same amount was created every two days.

By 2013, it's expected that the time will shrink to 10 minutes.

Every hour, we  
create enough  
Internet traffic  
to fill  
**7 billion**  
**DVDs**

Side by side, that's  
that's seven times  
the height of

generados 5  
exabytes hasta  
2003

1000 millones de  
GB = 2 millones  
de ordenadores

1 exabyte =  
1000 millones de  
GB

2011 cada 2 días  
se generaban 5  
exabytes

2025 se espera  
que generemos  
más de 400  
exabytes al día

# 2019 *This Is What Happens In An Internet Minute*



**Todo lo que  
hacemos en  
internet se  
guarda  
(o se puede  
guardar)**

# 2019 *This Is What Happens In An Internet Minute*



**Todo lo que  
hacemos en  
internet se  
guarda  
o se puede  
guardar)**

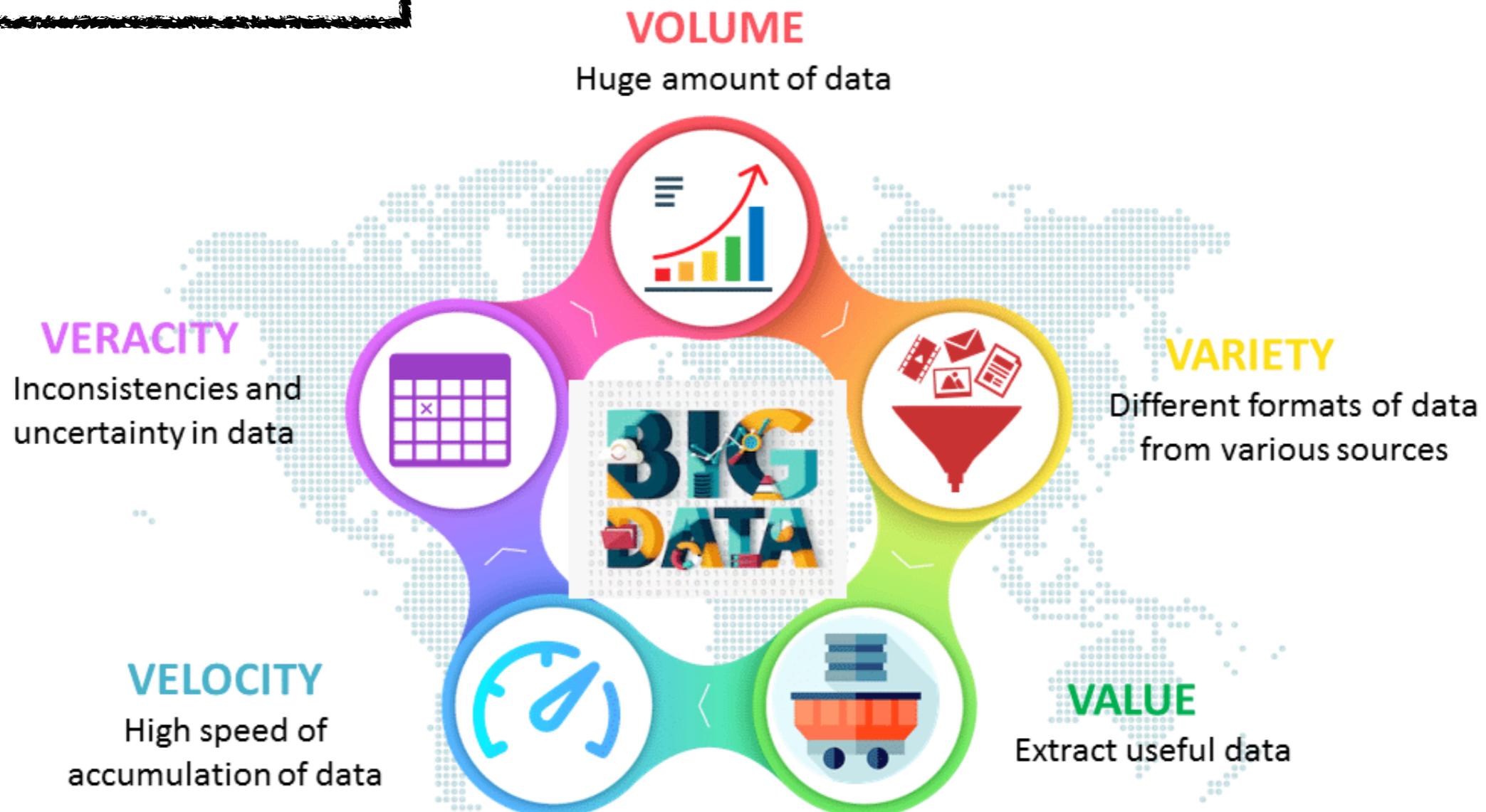
- Cada vez recogemos y almacenamos más datos pero “*solo unos pocos % se analizan correctamente*” (Cunef)
- Ej. “las empresas que utilizan sus datos para tomar decisiones tienen 23 veces más posibilidades de conseguir nuevos clientes y 6 veces menos de perderlos” (McKinsey Business Technology Office)
- También tiene interés en otros campos no económicos: predicción **de crisis climáticas**, gestión de **catástrofes, seguridad**, ciudades **inteligentes**, internet de las cosas

- Cada vez recogemos y almacenamos más datos pero “solo unos pocos % se analizan correctamente” (Cunef)
- Ej. “las empresas que utilizan sus datos para tomar decisiones tienen 23 veces más posibilidades de conseguir nuevos clientes y 6 veces menos de perderlos” (McKinsey Business Technology Office)
- También tiene interés en otros campos no económicos: predicción **de crisis climáticas**, gestión de **catastrofes, seguridad**, ciudades **inteligentes**, internet de las cosas

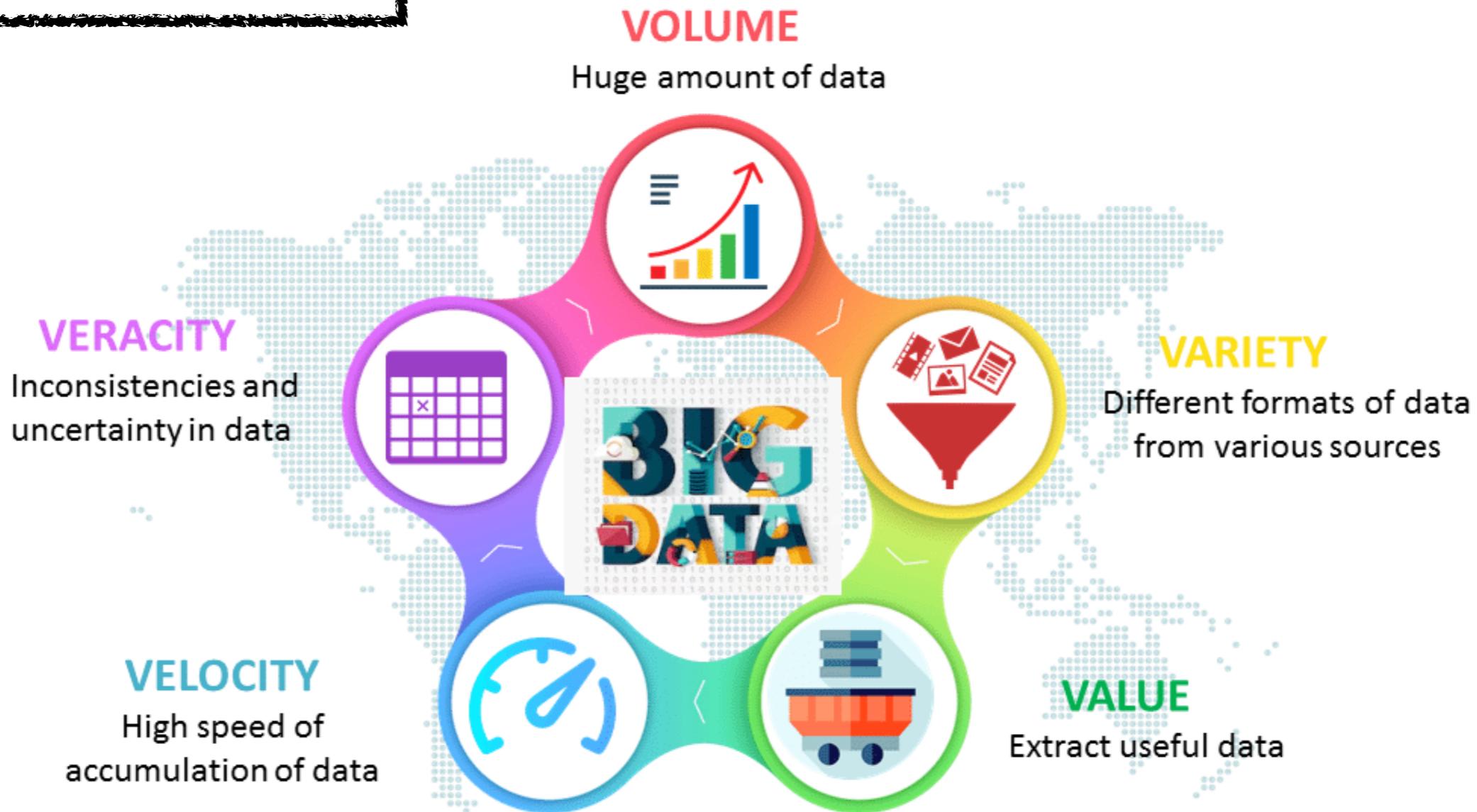


La Casa Blanca ha invertido más de 200,000 millones en proyectos de Big Data

# 5V BIG DATA

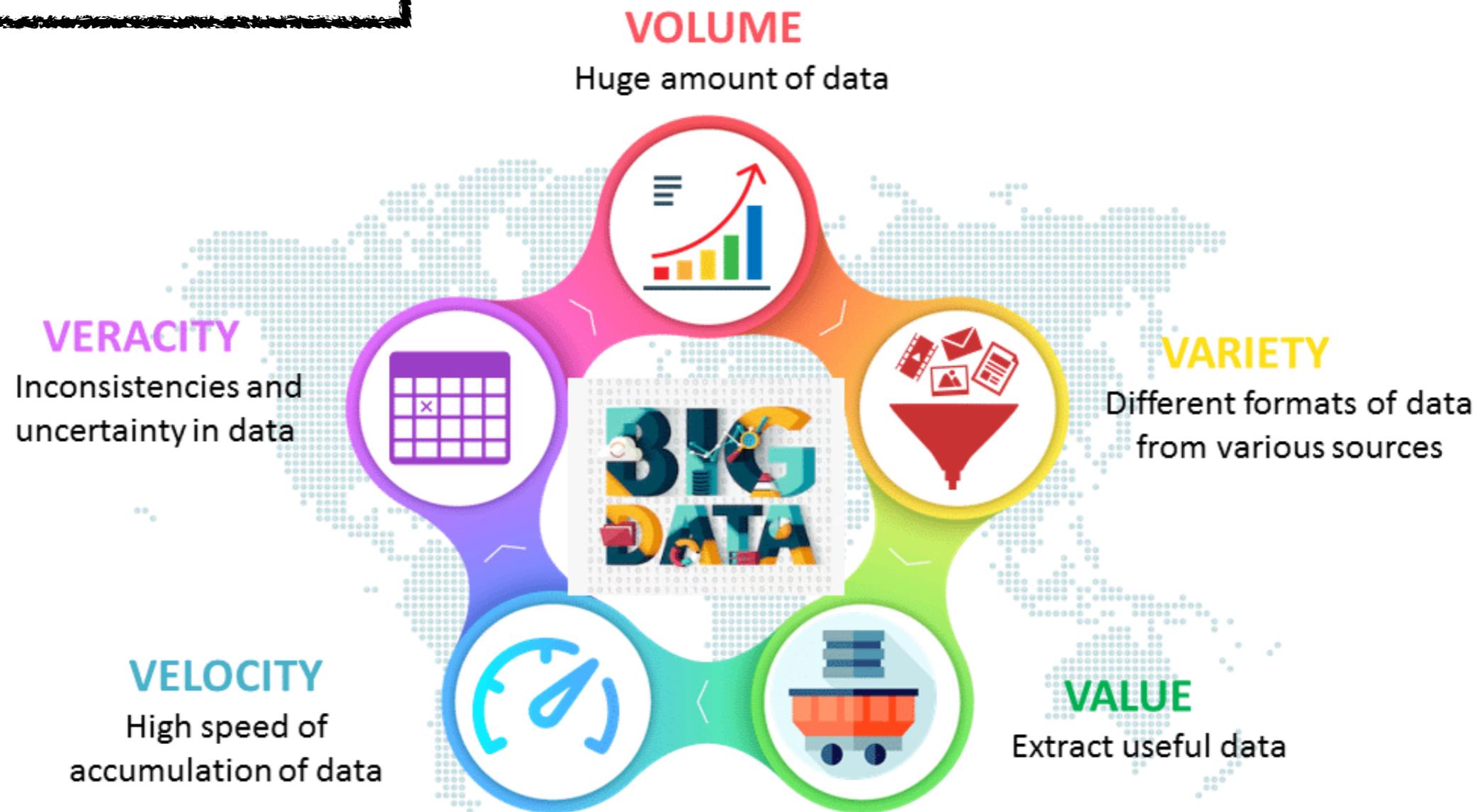


# 5V BIG DATA



Informáticos,  
ingernieros de  
sistemas,...

# 5V BIG DATA



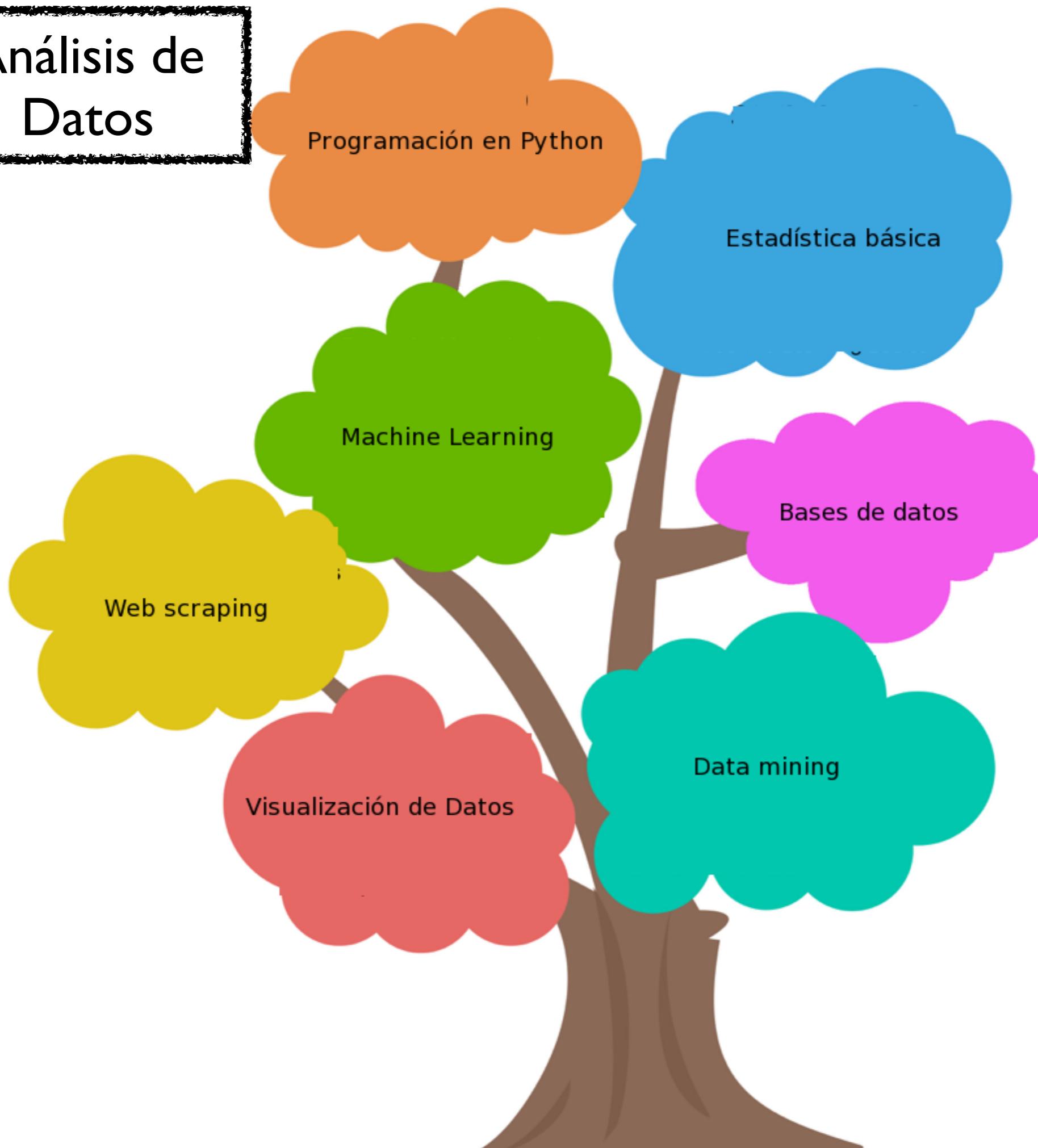
Informáticos,  
ingernieros de  
sistemas,...

Analistas de Datos  
(matemáticos,  
físicos,...)

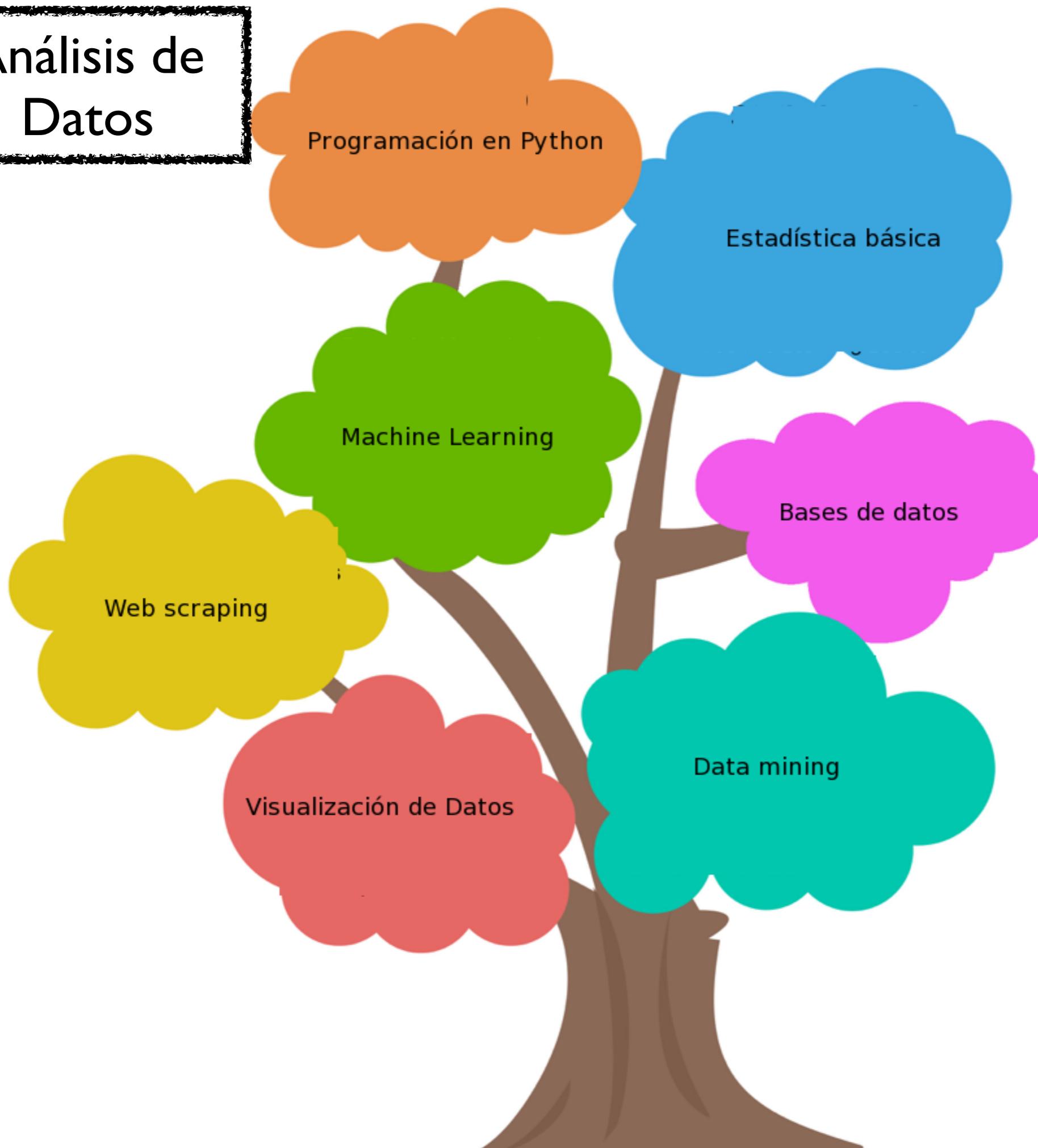
# Análisis de Datos



# Análisis de Datos



# Análisis de Datos



# Análisis de Datos

**Biuse.** Programación y redes neuronales



Web scraping

**Steven.** Visualización



Programación en Python

Visualización de Datos

Estadística básica

Machine Learning

Bases de datos

Data mining

**Marcos.** Estadística



**Carlos.** Bases de datos y automatización de sistemas

**Diego.** Bases de datos, web scraping, machine learning



# Let's talk about you



Hemos asumido:

**Al empezar:**

\* nivel de matemáticas bajo

\* nivel programación zero

\* cultura digital **alta**

# Let's talk about you

Hemos asumido:

## Al empezar:

- \* nivel de matemáticas bajo
- \* nivel programación zero
- \* cultura digital **alta**

## Al acabar:

- \* nivel de matemáticas un poco más alto
- \* conocimientos de programación
- \* cultura big data
- \* Soltura en herramientas para interaccionar con datos

# Programación en Python

# Programación

**Programar :** Hablar con cualquier dispositivo (móvil, ordenador,...) en un leguaje que el ordenador entienda (java, C++, python,...)

**Órdenes sencillas.** Ej. suma, ve al link, etc

**O más complejas:** Secuencia de órdenes.  
**script, código, programa**

# Programación

**Programar :** Hablar con ordenador ha estado programado (móvil, ordenador,...) en un ordenador entienda (java, C-Google,...)

**TODO** lo que hacemos con el ordenador ha sido programado por alguien. Apps, E-mail, búsquedas...

Y hacen **exactamente** eso.

**Órdenes sencillas.** Ej. suma, ve al link, etc

**O más complejas:** Secuencia de órdenes.  
**script, código, programa**

# Programación

**Un código no tiene porque hacer cálculos:**

ej.1

1. Abre un archivo
2. Busca la palabra árbol
3. Cambiala por cerezo
4. Guarda archivo

ej.2

1. Si presiono flecha mueve bruja
2. Mueve bruja hasta que presione flecha
3. Si bruja toca dragón, bruja muere



4. Si toca enter bruja hace hechizo

# Programación

**Un código no tiene porque hacer cálculos:**

ej.1

1. Abre un archivo
2. Busca la palabra árbol
3. Cambiala por cerezo
4. Guarda archivo

ej.2

1. Si presiono flecha mueve bruja
2. Mueve bruja hasta que presione flecha
3. Si bruja toca dragón, bruja muere



4. Si toca enter bruja hace hechizo

Scratch

# Programación en Python



## **Python :**

- Sencillo
- Mucha información On-line
- Muchos paquetes ya programados para usar
- Interactúa bien con otros programs (excel, VLC, etc...)
- Gratis !!

## **¿Quién usa python?:**

- Google
- NASA
- Youtube
- IFCA
- .....



Google Cloud Platform



# Programación en Python



## **Python :**

- funciones básicas 3+2, type(a), print()
- y paquetes numpy, scipy, matplotlib, ...

## **Anaconda :**

- Plataforma de python que te permite instalar paquetes fácilmente
- No es estrictamente necesario



# Programación en Python



## Documento de texto .py :

codigo.py

```
# programa python
x=3
y=4
print('la suma es',x+y)
```

```
>> python codigo.py
>> la suma es 7
```

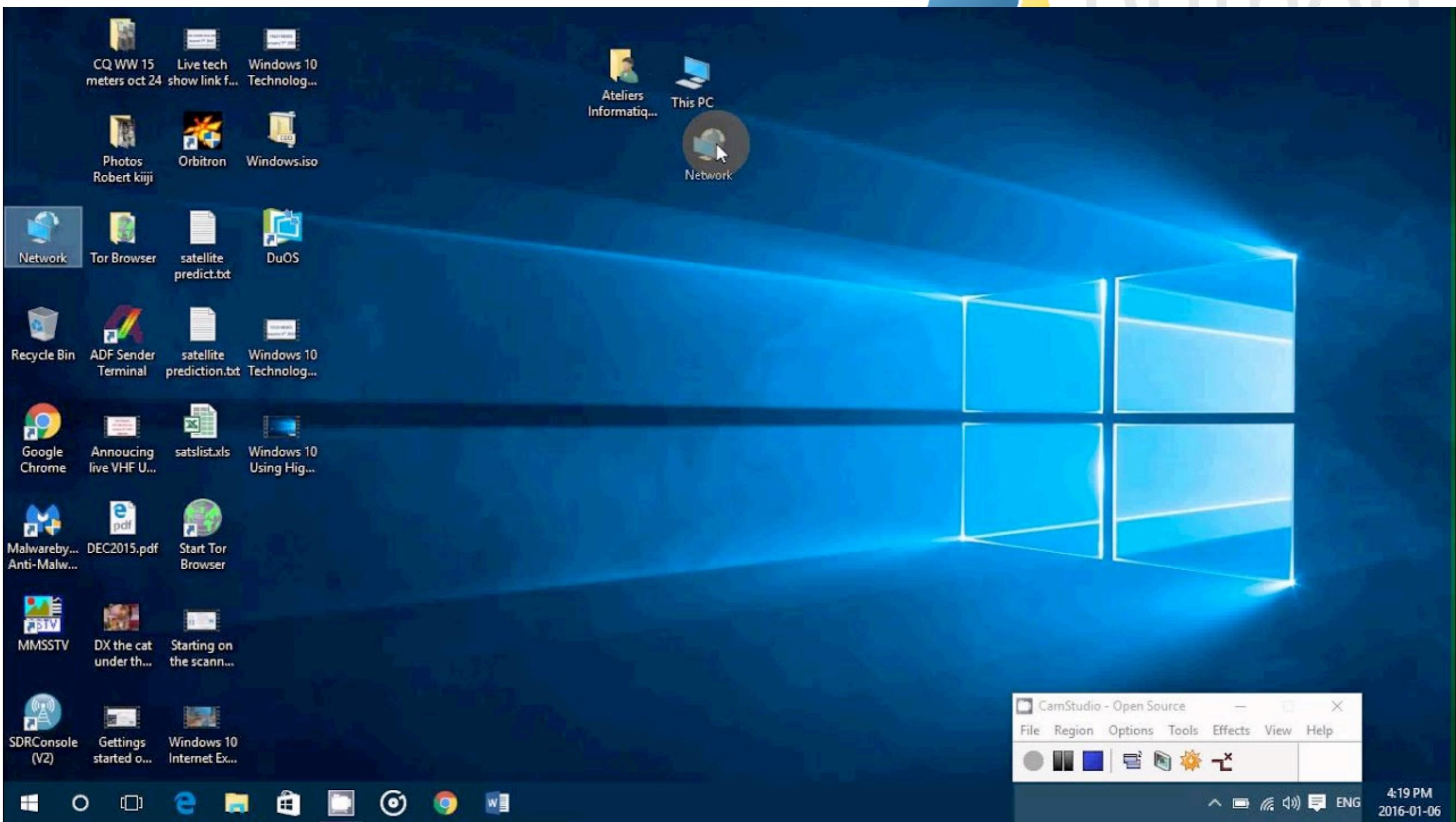
## Jupyter Notebook:

A screenshot of a Jupyter Notebook interface. At the top, there's a menu bar with File, Edit, View, Insert, Cell, Kernel, Widgets, and Help. Below the menu is a toolbar with various icons for file operations like new, open, save, and run. The main area shows a code cell labeled "In [1]:" containing the Python code from the previous slide. The output of the cell, "la suma es 7", is displayed below the code. The notebook has a clean, modern design with a light gray background and white cells.

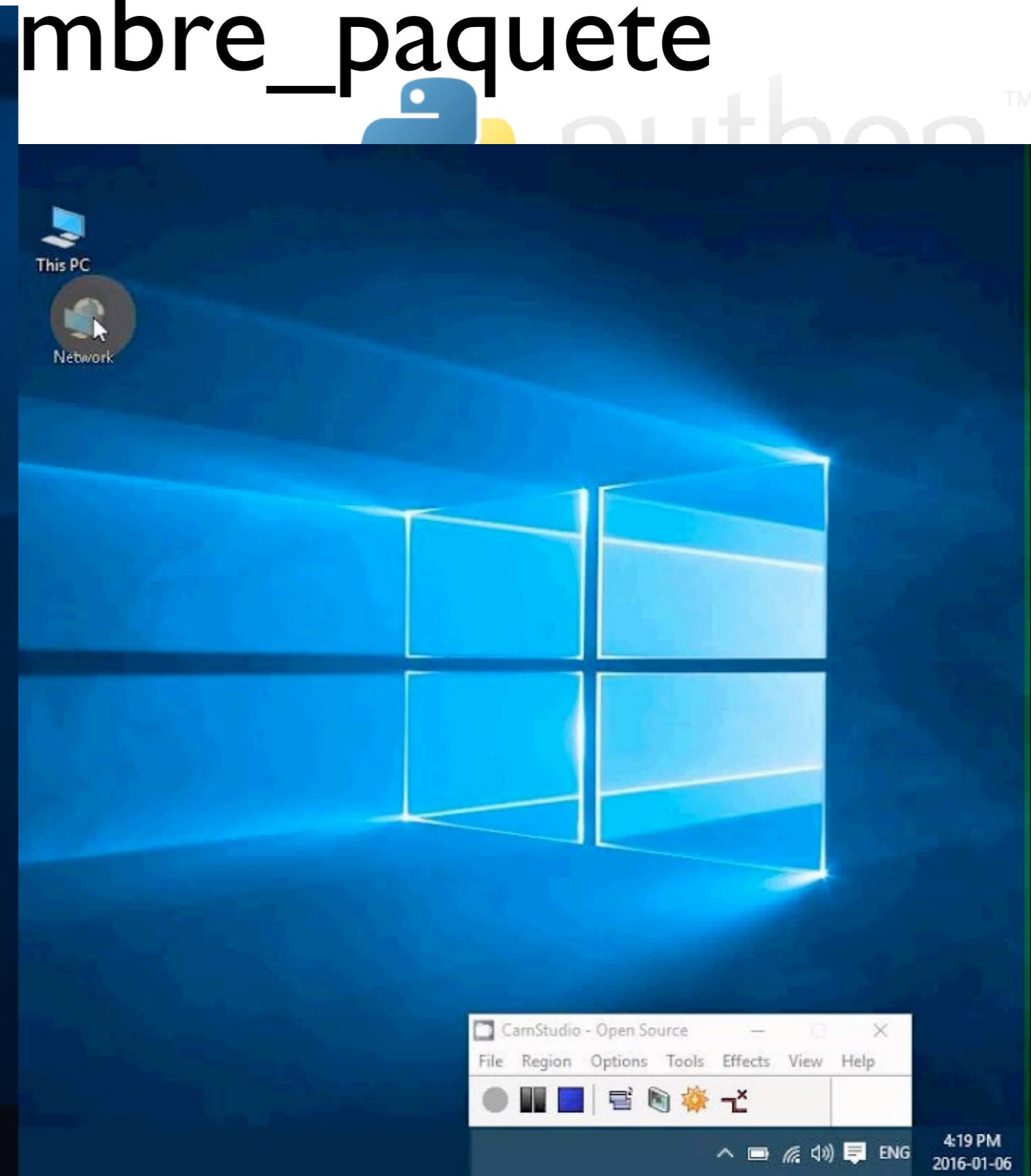
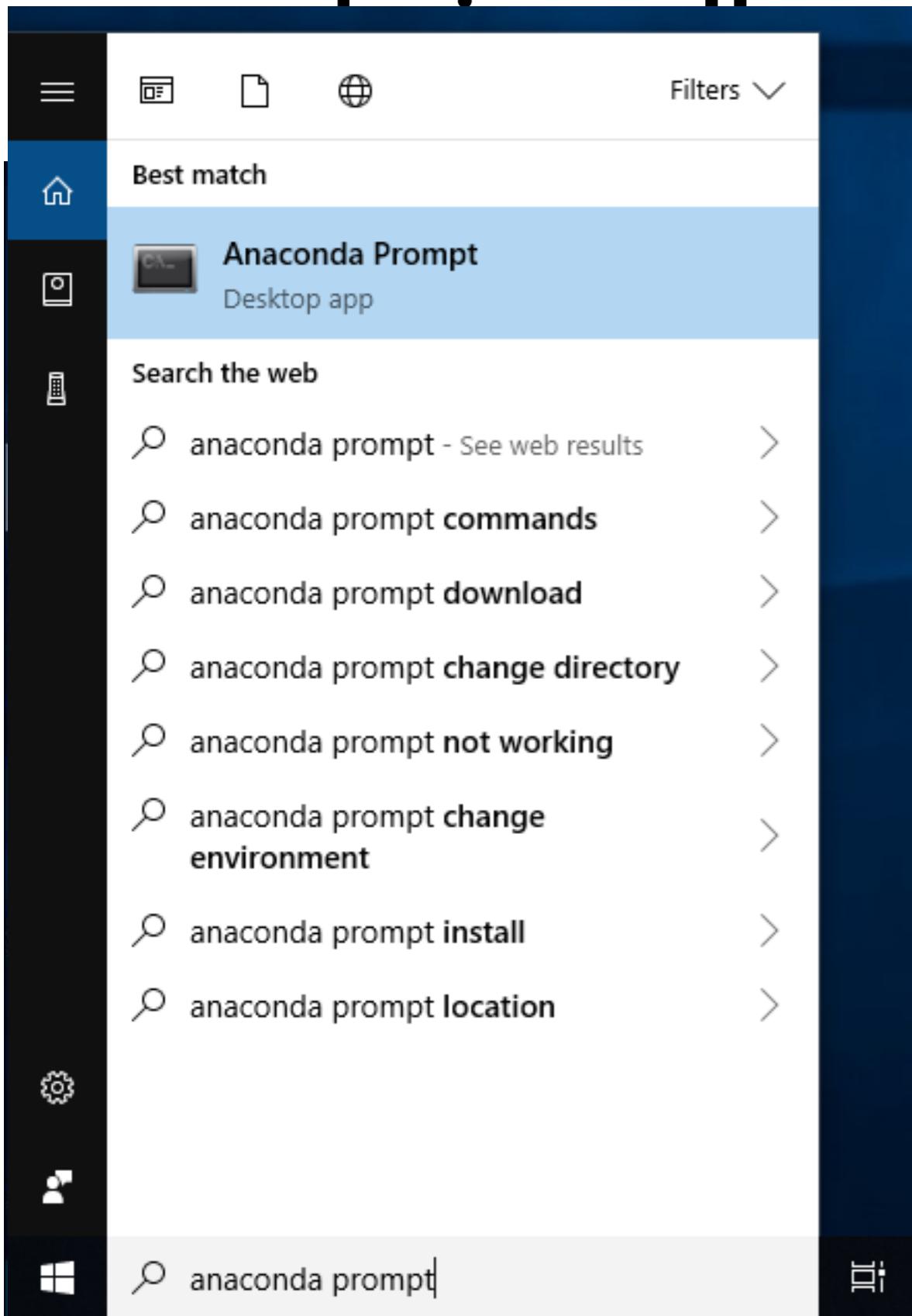
# Instalación de paquetes: conda install nombre\_paquete



TM  
python

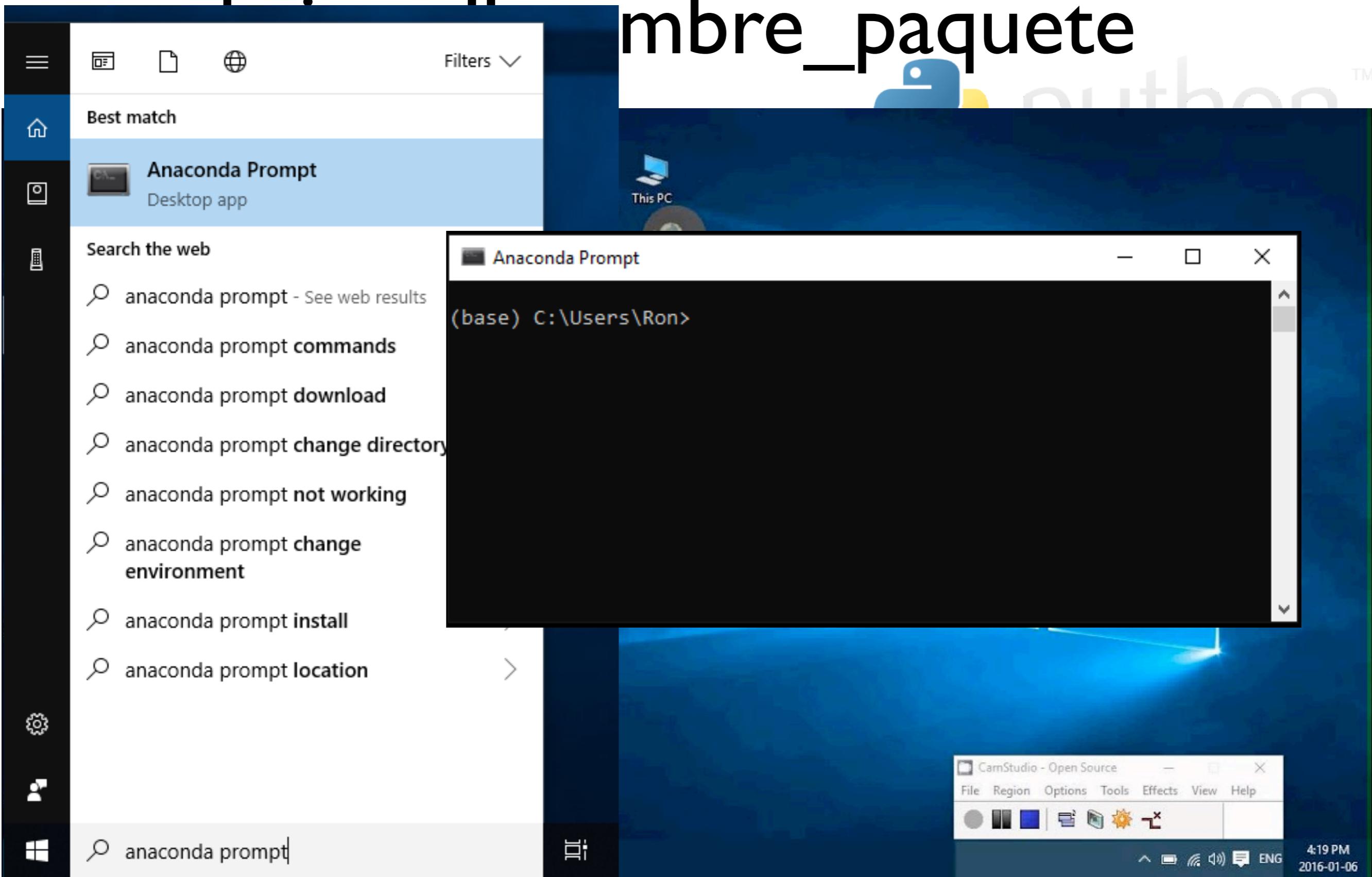


# Instalación de paquetes: nombre\_paquete



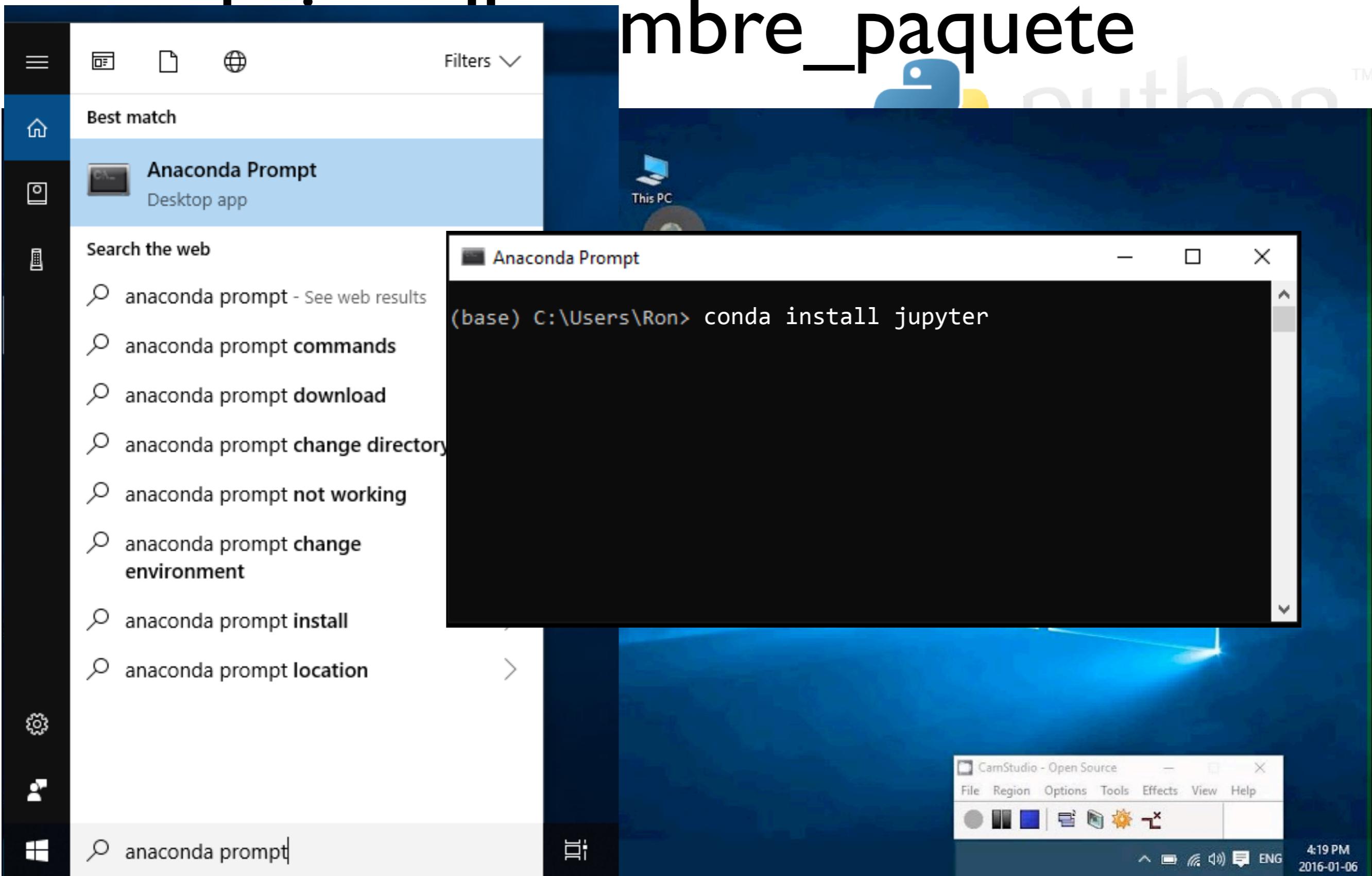
# Instalación de paquetes:

## Nombre\_paquete



# Instalación de paquetes:

## Nombre\_paquete



# Programación en Python

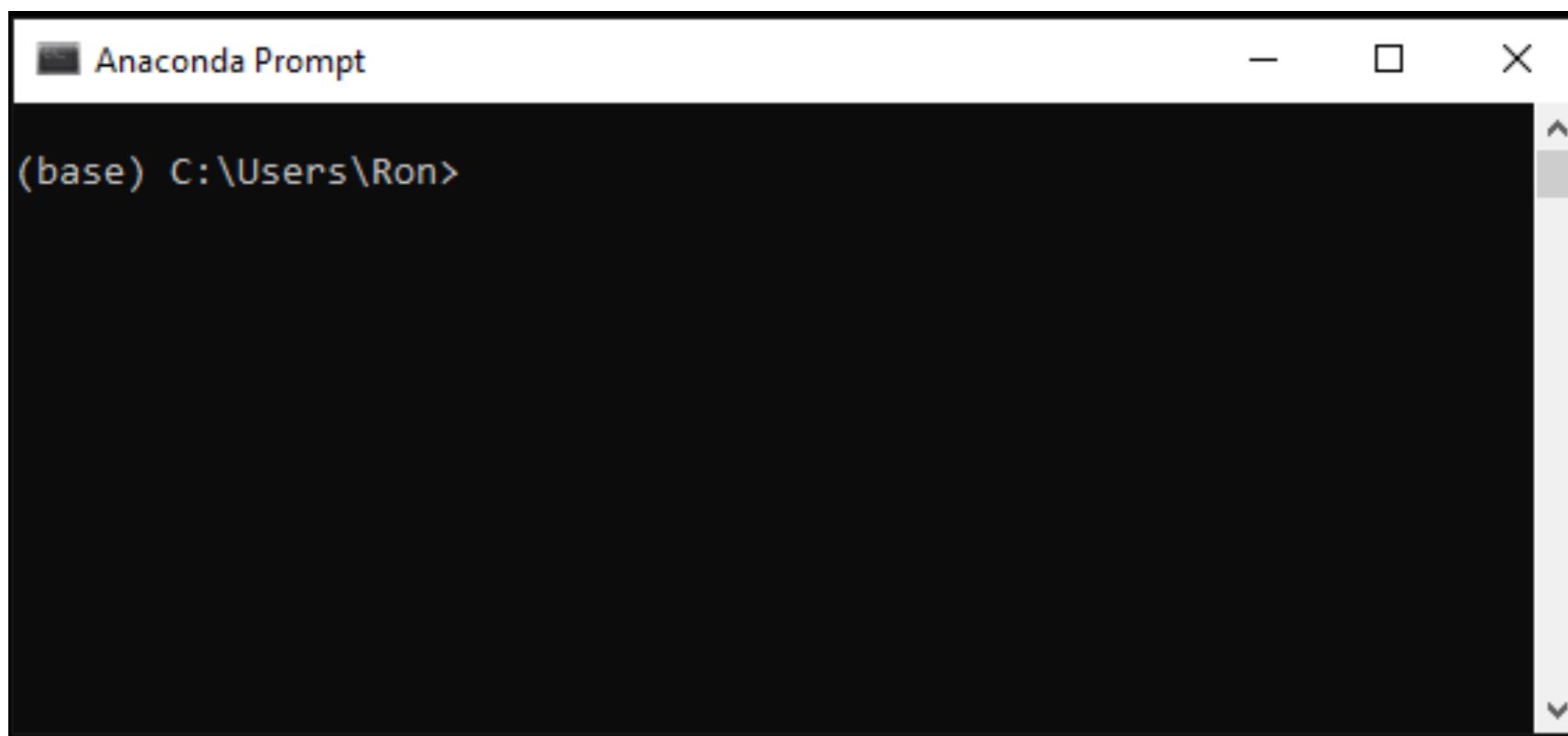


Abrimos la notebook (libreta)  
de Jupyter

# Programación en Python



Abrimos la notebook (libreta)  
de Jupyter



A screenshot of the Anaconda Prompt window. The title bar says "Anaconda Prompt". The command line shows "(base) C:\Users\Ron>". The window has a dark theme with a vertical scroll bar on the right side.

# Programación en Python



Abrimos la notebook (libreta)  
de Jupyter

A screenshot of the Anaconda Prompt window. The title bar says "Anaconda Prompt". The main area shows the command "(base) C:\Users\Ron> jupyter notebook" entered and ready to be run. The window has standard minimize, maximize, and close buttons at the top right.

# Programación en Python



<https://github.com/biuse/TallerEmpleo>