

# **Software libre para investigación: Una (muy breve) introducción**

Mario Gavidia-Calderón

# Hola!

IAGUSP

Instituto de Astronomia, Geofísica  
e Ciências Atmosféricas



IAG-USP

# **Modelos de calidad del aire**

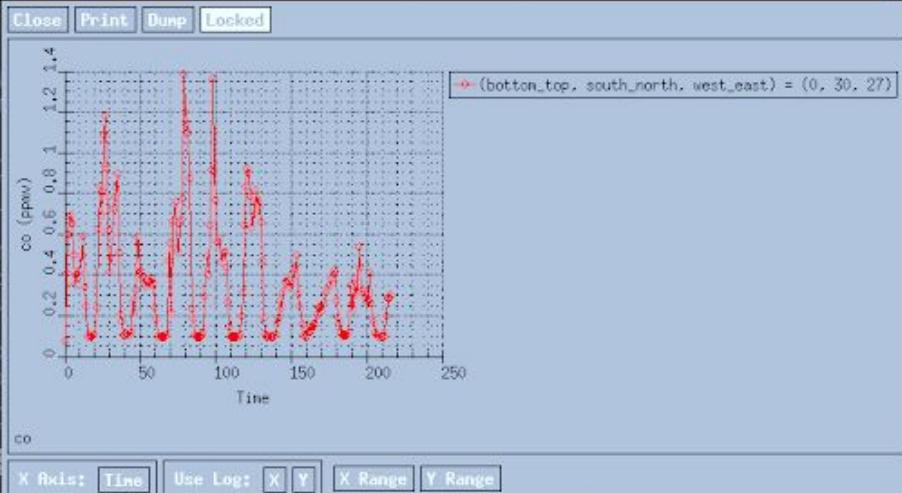
# Análisis de datos atmosféricos y dataviz

**Todo lo hago con  
Software libre**

mgavida@svante2: /scr2/mgavida/WRF4/WRF/test/em\_real

moavida@svante2: ~ 74x19

co (on svante2)



mgavida@svante2: /scr2/mgavida/WRF4/WRF/test/em\_real 74x19  
windturbines.txt  
wrfbdy\_d01  
wrfchemi\_12z\_d01  
wrfchemi\_d01\_2018-06-21\_00:00:00  
wrfchemi\_d02\_2018-06-21\_00:00:00  
wrf\_chem\_val  
wrf.exe  
wrfinput\_d01  
wrfinput\_d02  
wrfout\_d01\_2018-06-21\_00:00:00  
wrfout\_d01\_2018-06-21\_00:00:00\_2way\_test  
wrfout\_d01\_2018-06-21\_00:00:00\_aas4wrf\_test  
wrfout\_d01\_2018-06-21\_00:00:00\_first\_a4w  
wrfout\_d02\_2018-06-21\_00:00:00





“...es todo **software** cuyo código fuente puede ser estudiado, modificado, y **utilizado libremente** con cualquier fin y **redistribuido** con cambios y/o **mejoras** o sin ellas.”

**Por qué estamos aquí?**

# **1. Cuáles son las softwares libres más utilizados?**

- 1. Cuáles son las softwares libres más utilizados?**
- 2. Cuáles son los pro y contras?**

- 1. Cuáles son las softwares libres más utilizados?**
- 2. Cuáles son los pro y contras?**
- 3. Por qué son importantes para la investigación?**

**Software libres + utilizados en  
investigación**

# Hacer el experimento

**Hacer el experimento**

**Calcular/Analizar resultados (datos)**

**Hacer el experimento**

**Calcular/Analizar resultados (datos)**

**Visualizar resultados**

**Hacer el experimento**

**Calcular/Analizar resultados (datos)**

**Visualizar resultados**

**Reportar**

**Hacer el experimento**

**Calcular/Analizar resultados (datos)**

**Visualizar resultados**

**Reportar**

# Analizar Resultados

# Microsoft Excel

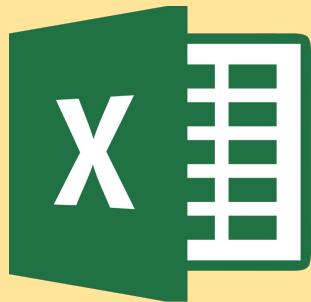


# Microsoft Excel



 **LibreOffice**  
The Document Foundation

# Microsoft Excel



 **LibreOffice**  
The Document Foundation



Google Sheets

# (Un paréntesis personal)



Google Drive



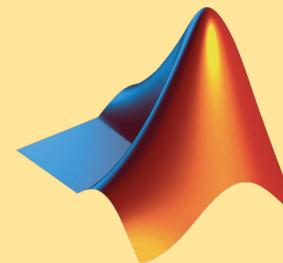
**LibreOffice**  
The Document Foundation





*Booo!!*

# A veces necesitamos armas más pesadas



# A veces necesitamos armas más pesadas



ArcGIS



AUTOCAD



# Pero hay software libre para todos los gustos



ArcGIS



AUTOCAD



# Pero hay software libre para todos los gustos



ArcGIS



AUTOCAD



**Y si no es necesario  
tanto software???**

**Lenguajes de  
programación al  
rescate!**

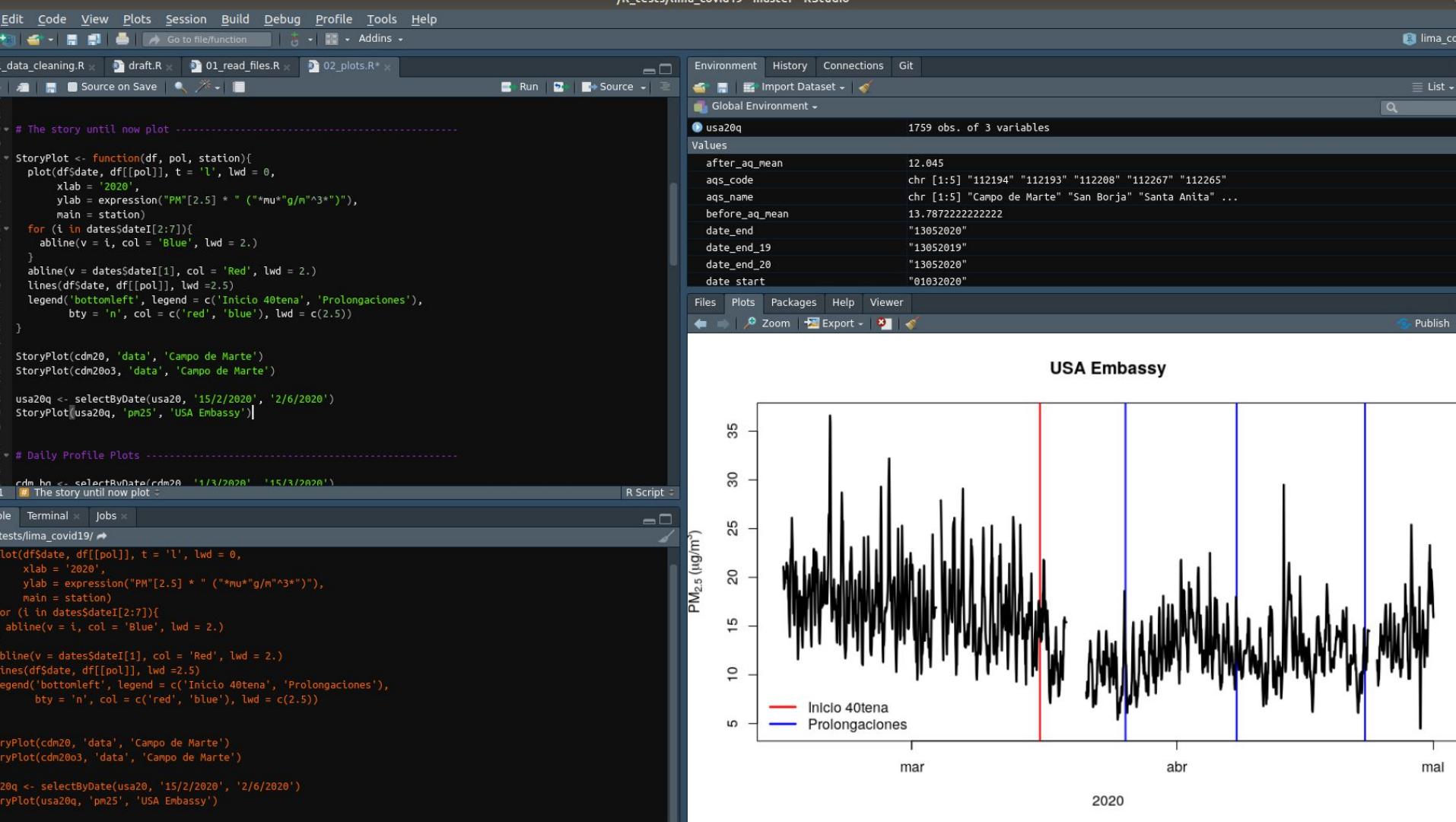


# (Otro paréntesis personal)



# R & Python

- Lenguajes de programación
- Opensource
- Multiplataforma
- Los más usados en ciencia de datos
- Muchos paquetes para diferentes necesidades



```

rf_eval
    _pycache_
    old_plots
    wrf_sp_eval
        cletes2017_lation.dat
        co_Cid.Universitária-USF
        co_Ibirapuera.png
        co_Parque.D.Pedro II.pn
        co_Pinheiros.png
        co_São Caetano do Sul.p
        data_preparation.py
        fig_test_3km.tar.gz
        fig_test.tar.gz
        met_20_06_2018-29_06
        model_eval_sp.py
        model_stats.py
        no_Cid.Universitária-USF
        no_Ibirapuera.png
        no_Parque.D.Pedro II.pn
        no_Pinheiros.png
        no_São Caetano do Sul.p
        no2_Cid.Universitária-US
        no2_Ibirapuera.png
        no2_Parque.D.Pedro II.p
        no2_Pinheiros.png
        no2_São Caetano do Sul
        o3_Cid.Universitária-USF
        o3_Ibirapuera.png
        o3_no_no2_co
        o3_no_no2_co_stats.csv
        o3_Parque.D.Pedro II.pn
        o3_Pinheiros.png
        o3_Pinheiros.svg
        o3_São Caetano do Sul.p
        photo_comp_Pinheiros.p
        pol_20_06_2018-29_06_
        qualar_py.py
        rh2_Cid.Universitária-US
        rh2_Ibirapuera.png
        rh2_Parque.D.Pedro II.pr
        rh2_Pinheiros.png
        rh2_São Caetano do Sul
        t2_Cid.Universitária-USP
        t2_Ibirapuera.png
        t2_Parque.D.Pedro II.pn
        t2_Pinheiros.png
        t2_rh2_ws_wd
        t2_rh2_ws_wd_stats.csv
        t2_São Caetano do Sul.p

```

/scr2/mgavida/python\_stunts/test/read\_emissions.py

Editor Spyder

```

Este é um arquivo de script temporário.

import pandas as pd
import xarray as xr

file_name1km = '/scr2/mgavida/WRF4/in_out/emissions_1km.txt'
file_name3km = '/scr2/mgavida/WRF4/in_out/emissions_3km.txt'
file_name9km = '/scr2/mgavida/WRF4/in_out/emissions_9km.txt'

df = pd.read_csv(file_name1km,
                  delimiter='t',
                  names=['i', 'lon', 'lat', 'E_CO', 'E_HCHO', 'E_C2H5OH',
                         'E_KET', 'E_NH3', 'E_XYL', 'E_TOL', 'E_ISO',
                         'E_OLI', 'E_OLT', 'E_OL2', 'E_HC8', 'E_HC5', 'E_ORA2',
                         'E_ETH', 'E_ALD', 'E_CSL', 'E_SO2', 'E_HC3', 'E_NO2',
                         'E_NO', 'E_CH3OH', 'E_PM25I', 'E_PM25J', 'E_SO4I',
                         'E_SO4J', 'E_NO3I', 'E_NO3J', 'E_ORGT', 'E_ORGJ',
                         'E_ECI', 'E_ECJ', 'E_SO4C', 'E_NO3C', 'E_ORGC',
                         'E_ECC'])

n_lon = len(df.lon.unique())
n_lat = len(df.lat.unique())

n_points = n_lon * n_lat
lonId = df['lon']: n_lon].values
latId = df['lat']: n_points: n_lon].values[:::-1]

date = pd.date_range('2018-06-21 00:00',
                     '2018-06-30 23:00',
                     freq='H')

emiss_names = df.columns[3:]

def create_dataarray_per_emi(emiss_df, pol, lat, lon, date):
    """
    Create a xarray dataarray for each emission species
    by reshaping the emission dataframe (emiss df) according the lat, lon
    and date, add emission_zdim dimension and add attribute names to lat and
    lon dimensions.

    Parameters
    -----
    emiss_df : pandas DataFrame
        create by read_csv emission_file
    pol : string
        name of emitted pollutant.
    lat : numpy ndarray
        latitudes of emission_file.
    lon : numpy ndarray
    """

    # Create a new DataFrame with the same columns as emiss_df, but with
    # the first three columns removed.
    # This will be used to store the reshaped data.
    new_df = pd.DataFrame(emiss_df)
    new_df = new_df.drop(emiss_names[:3], axis=1)

    # Reshape the DataFrame into a 3D array.
    # The first dimension is the date index.
    # The second dimension is the latitude.
    # The third dimension is the longitude.
    # The fourth dimension is the emission species.
    new_df = new_df.unstack([0, 1, 2]).T

    # Add the emission_zdim dimension.
    new_df = new_df.set_index(emiss_names[:3])

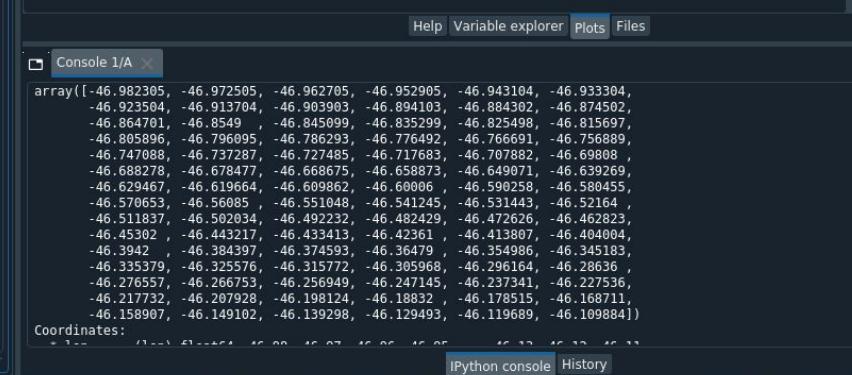
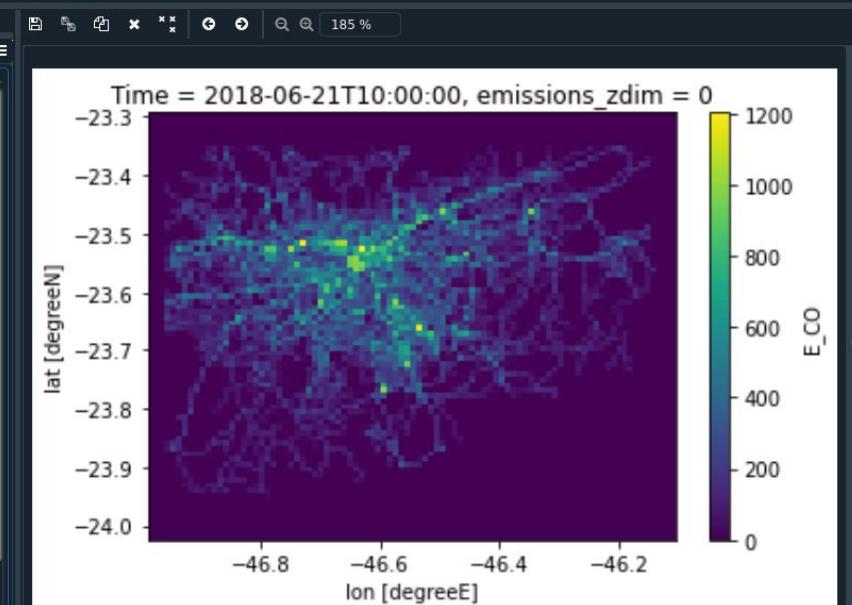
    # Add attributes to the latitude and longitude.
    new_df['lat'] = lat
    new_df['lon'] = lon
    new_df['date'] = date

    # Create a new xarray DataArray.
    da = xr.DataArray(new_df, dims=[emiss_names[:3], 'lat', 'lon', 'date'])

    # Set the coordinates.
    da['lat'].attrs['name'] = 'lat'
    da['lat'].attrs['units'] = 'degreeN'
    da['lon'].attrs['name'] = 'lon'
    da['lon'].attrs['units'] = 'degreeE'
    da['date'].attrs['name'] = 'date'
    da['date'].attrs['units'] = 'ISO8601'

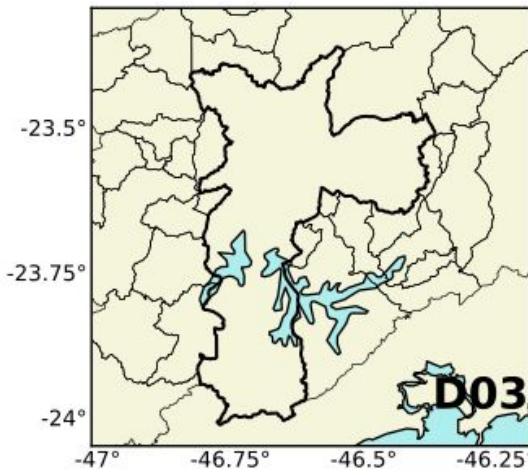
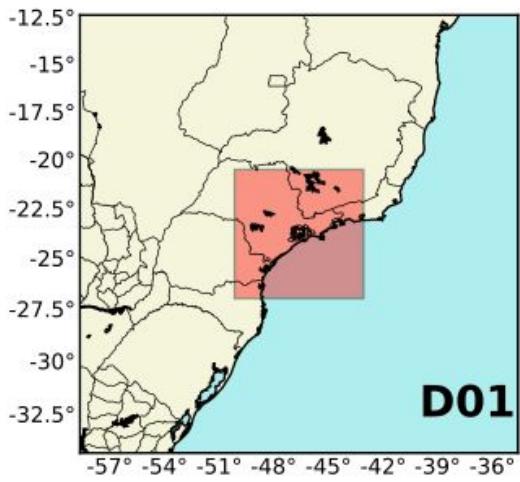
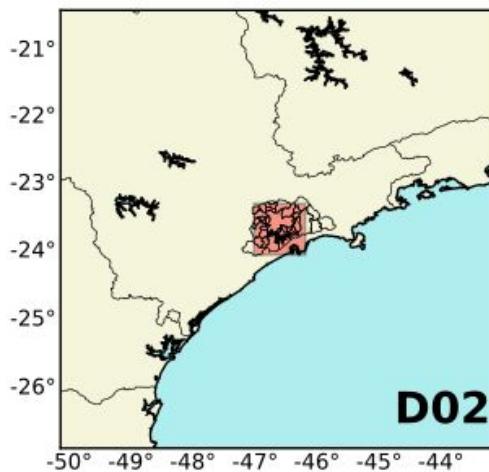
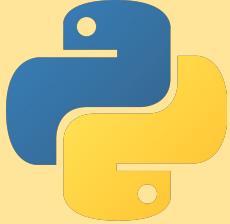
    return da

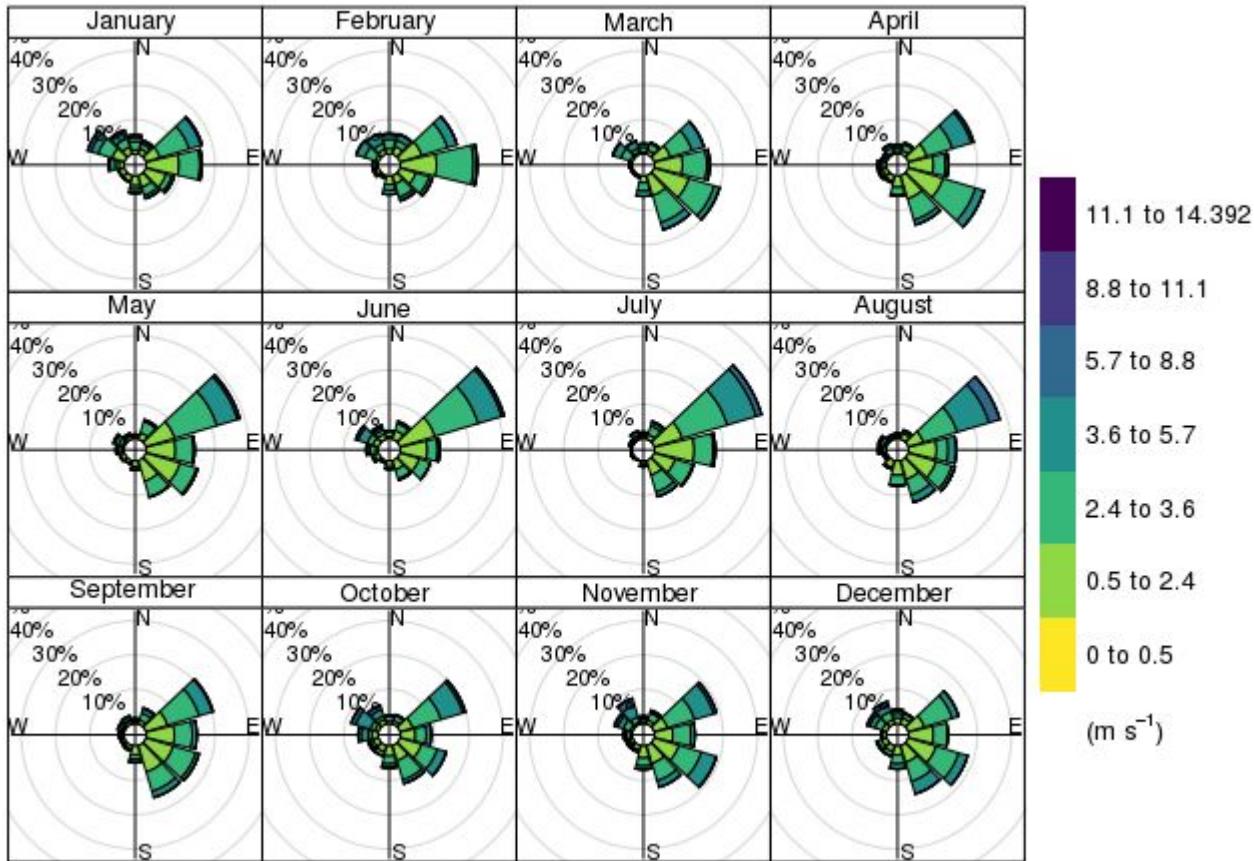
```

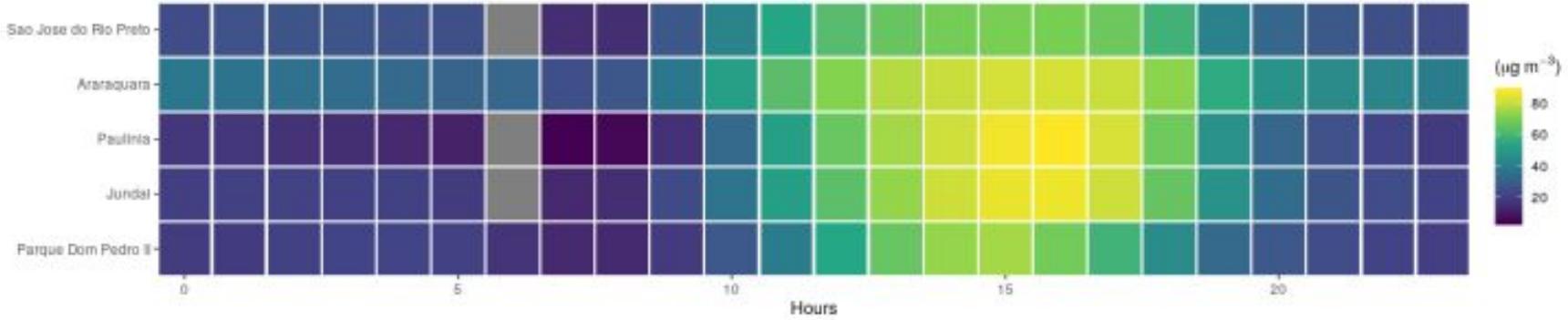


# Visualizar datos

**R y Python tienen  
excelentes paquetes  
para dataviz**







# Reportar

**Word puede hacerte  
enojar**

# Moving a picture in Microsoft Word



■ You  
the

mess up  
whole document

■ It  
actually  
does  
what  
you  
want



**LATEX** al rescate!

**What you see is what  
you get  
(WYSIWYG)**

**What you see is what  
you mean  
(WYSIWYM)**

```

8 + \begin{sidewaysfigure}[ht]
9   \includegraphics[width=\textwidth]{figs/interp_method_eval_pop_dens.pdf}
10  \caption{Test of different interpolation methods to calculate population density.}
11  \label{fig:pop_den}
12 \end{sidewaysfigure}

13
14 Another solution is to upscale PM2.5 simulation to the Gridded Population of the World
15 resolutions, because it has data at 30 sec, 2.5 min, 15 min, 30 min and 60 min resolution.

16
17
18 \subsection{Exposition curves}
19
20 The methodology was based on \cite{Gao2018}. First we have to calculate the Relative Risk ($RR$) by using
21 Integrated Exposure Response (IER) functions. With RR, we can calculate Population Attributable Fractions is
22 ($PAF$), which will be used to calculate Mortality ($\Delta M$).
23
24 \begin{equation*}
25   RR_{i,j,k}(C_l) = \begin{cases} 1 + \alpha_{i,j,k}(1 - e^{-\beta_{i,j,k}(C_l - C_0)})^{\gamma_{i,j,k}} & C_l > C_0 \\ 1 & C_l < C_0 \end{cases}
26 \end{equation*}
27 \label{eq:1}
28 \end{equation*}
29
30 \begin{equation*}
31   PAF_{i,j,k} = \frac{RR_{i,j,k}(C_l) - 1}{RR_{i,j,k}(C_l)}
32 \end{equation*}
33
34 \begin{equation*}
35   \Delta M = PAF_{i,j,k,l} * y_{0,i,j,k,l} * Pop_{i,j,l}
36 \end{equation*}
37
38 These equations $C_l$ is the annual PM2.5 concentration from WRF-Chem, $C_0$, the
39 counterfactual concentration (in \cite{Cohen2017} is the Theoretical minimum-risk exposure level);
40 $\alpha_{i,j,k}$, $\beta_{i,j,k}$ and $\gamma_{i,j,k}$ are the parameters that describe the IER curves in
41 $i^{th}$ age and $j^{th}$ sex group for the $k^{th}$ disease. In the
42 mortality equation $y_{0,i,j,k,l}$ is the current age-sex-specific mortality rate for the
43 $k^{th}$ and $l^{th}$ disease and $Pop_{i,j,l}$ is the exposed population in that grid cell ($\$).
44
45 In \cite{Gao2018}, they calculate the mortality for the following diseases: "ischemic heart disease
46 (IHD), stroke (STK, including both ischemic and hemorrhagic stroke), lung cancer (LC), and chronic obstructive
47 pulmonary disease (COPD), and for one disease among young children, acute lower respiratory infections
48 (LRI)".
49
50 Right now, we are still trying to get the parameters that shape the IER curves.
51
52
53 \subsection{Mortality Calculation}
54 On stand-by.

```

## 3.2 Spatial interpolation methods

Works that estimates mortality using model results used different spatial resolutions, going from 30 km [2] to 0.5 degrees [4]. So, to get comparable results with other studies, a process of **upsampling** has to be made upon WRF-Chem simulation and also population count data(Gridded Population of the World (GPW)) <sup>2</sup>.

This **upsampling** process rose problems. Because interpolation methods, such as bi-linear and nearest-neighbor, can create *new population*, we can end with more population than the original data. Other problem was that the spatial distribution of population doesn't preserves, getting less dense population hot spots.

One solution was to, instead of interpolation population count, it's best to interpolate population density, and calculate population count by multiplying it by the cell (grid) area (Figure 5, at the end of the document). Further, there is a **conservative method** for spatial interpolation that gives better results, but its implementation is not as easy as the other mentioned methods and it doesn't conserve 100%. This conservative method was used to upscale population count and PM<sub>2.5</sub> (Figure 4).

Another solution is to upscale PM<sub>2.5</sub> simulation to the Gridded Population of the World resolutions, because it has data at 30 sec, 2.5 min, 15 min, 30 min and 60 min resolution.

## 3.3 Exposition curves

The methodology was based on [2]. First we have to calculate the Relative Risk (\$RR\$) by using Integrated Exposure Response (IER) functions. With RR, we can calculate Population Attributable Fractions is (\$PAF\$), which will be used to calculate Mortality (\$\Delta M\$).

$$RR_{i,j,k}(C_l) = \begin{cases} 1 + \alpha_{i,j,k}(1 - e^{-\beta_{i,j,k}(C_l - C_0)})^{\gamma_{i,j,k}} & C_l > C_0 \\ 1 & C_l < C_0 \end{cases}$$

$$PAF_{i,j,k} = \frac{RR_{i,j,k}(C_l) - 1}{RR_{i,j,k}(C_l)}$$

$$\Delta M = PAF_{i,j,k,l} * y_{0,i,j,k,l} * Pop_{i,j,l}$$

In these equations \$C\_l\$ is the annual PM<sub>2.5</sub> concentration from WRF-Chem, \$C\_0\$, the counterfactual concentration (in [1] is the Theoretical minimum-risk exposure level); \$\alpha\_{i,j,k}\$, \$\beta\_{i,j,k}\$ and \$\gamma\_{i,j,k}\$ are the parameters that describe the IER curves in \$i^{th}\$ age and \$j^{th}\$ sex group for the \$k^{th}\$ disease. In the mortality equation \$y\_{0,i,j,k,l}\$ is the current age-sex-specific mortality rate for the \$k^{th}\$ and \$l^{th}\$ disease and \$Pop\_{i,j,l}\$ is the exposed population in that grid cell (\$l\$).

In [2], they calculate the mortality for the following diseases: "ischemic heart disease (IHD), stroke (STK, including both ischemic and hemorrhagic stroke), lung cancer (LC), and chronic obstructive pulmonary disease (COPD), and for one disease among young children, acute lower respiratory infections (LRI)".

Right now, we are still trying to get the parameters that shape the IER curves.

<sup>2</sup>Retrieved from: <https://sedac.ciesin.columbia.edu/datacollection/gpw-v4>

# Pros y Contras

# Contras

- No son muy populares en todos los ambientes (debería)
- A veces pueden ser difíciles de instalar
- Hay una curva de aprendizaje

# Pros

- Son Gratis!
- No hay límite de licencias
- Hay mucha pero mucha ayuda en línea
- Son versátiles
- Puedes contribuir
- Todos para uno y uno para todos!

# Importancia en investigación

**Sabían que no  
siempre es posible  
reproducir los  
resultados de una  
investigación?**

**Incluso las hechas por  
tu yo del pasado?**

**Ayudan a replicar la  
investigación  
(Reproducibility)**

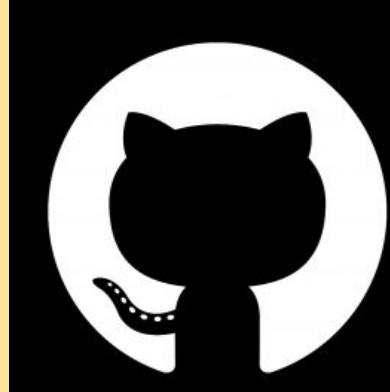
**Como no hay licencia,  
todos pueden acceder  
al lenguaje de  
programación**

**Se comparte el script,  
la data original y tu  
reporte para llegar al  
mismo resultado**

**Esto hace el proceso  
más transparente.**

**Aquí entra a la  
cancha...**

# Git & GitHub



**Software de control de cambios y plataforma de desarrollo colaborativo, respectivamente.**

**Esto suena  
complicado...**

# "FINAL".doc



FINAL.doc!



FINAL\_rev.2.doc



FINAL\_rev.6.COMMENTS.doc



FINAL\_rev.8.comments5.  
CORRECTIONS.doc



FINAL\_rev.8.comments5.  
CORRECTIONS.doc



FINAL\_rev.18.comments7.  
corrections9.MORE.30.doc



FINAL\_rev.22.comments49.  
corrections.10.#@\$%WHYDID  
ICOMETOGRAD SCHOOL????.doc



FINAL\_rev.22.comments49.  
corrections.10.#@\$%WHYDID  
ICOMETOGRAD SCHOOL????.doc

# (Otro paréntesis personal)

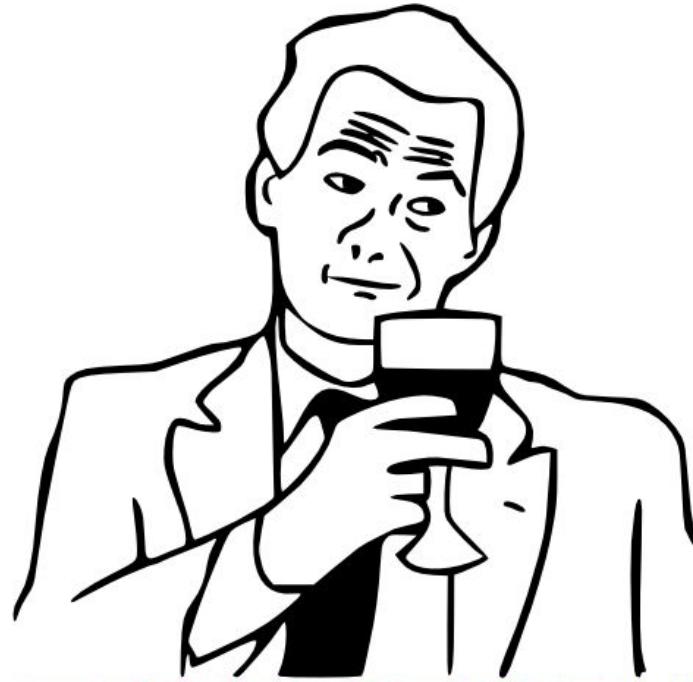


[tesis\\_borrador\\_final\\_final\\_porfin\\_delosporfines2.0\\_09052012.pdf](#)

# (Otro paréntesis personal)



tesis\_borra



**TRUE STORY**

WeKnowMemes

0\_09052012.pdf

# **Git evita esta situación**

**Mejor un ejemplo  
práctico**



github.com/quishqa/AAS4WRF.py

Unwatch 1 Star 0 Fork 0

Code Issues 0 Pull requests 0 Actions Projects 0 Wiki Security 0 Insights Settings

Another Assimilation System for WRF-Chem, python flavored. Edit

Manage topics

-o 20 commits 1 branch 0 packages 0 releases 1 contributor

Branch: master New pull request Create new file Upload files Find file Clone or download

quishqa	Now, conservative method is well implemented, cell area is not requir... [...]	Latest commit 1d512ff 2 days ago
	.gitignore	Correction in conservative method, multiply by cell area, regrid, and...
	README.md	Now, conservative method is well implemented, cell area is not requir...
	aas4wrf.py	Now, conservative method is well implemented, cell area is not requir...
	aas4wrf.yml	Change to YAML config file to better emission file set up, Print diff...
	aas4wrf_example.svg	Update README.md to use aas4wrf.yml, add img example
	emissions.txt	README.md updated and emissions.txt uploadad
	requirements.txt	Change to YAML config file to better emission file set up, Print diff...
	wrfchemi_zeros.py	script to create a wrfchemi with zero emissions

## README.md



# aas4wrf.py

`aas4wrf.py` is a Python script to create the `wrfchemi` file from local emissions needed to run WRF-Chem model. It's based on his older brother [AAS4WRF.ncl](#). Currently, It works with CBMZ/MOSAIC chemical mechanism and for surface emissions.

## Installation

You need to install the packages that `aas4wrf.py` needs. We recommend to use [miniconda](#) or [anaconda](#)

You can download this repo or clone it by:

```
git clone https://github.com/quishqa/AAS4WRF.py.git
```

Then add `conda-forge` channel by:

```
conda config --add channels conda-forge
```

To avoid conflicts during the Installation, we also recommend create a new environment to run `aas4wrf.py`:

```
conda create --name aas4wrf
conda activate aas4wrf
```

Commits · quishqa/AAS4WRF.py × +

github.com/quishtqa/AAS4WRF.py/commits/master

Branch: master

- Commits on Jun 9, 2020
  - Now, conservative method is well implemented, cell area is not requir...  1d512ff 
- Commits on Jun 8, 2020
  - Correction in conservative method, multiply by cell area, regrid, and...  9fa1c90 
- Commits on May 29, 2020
  - Fix usage msg and to accept yml with different names  7a4849d 
- Commits on May 26, 2020
  - Update README.md to use aas4wrf.yml, add img example  15c5186 
  - Change to YAML config file to better emission file set up, Print diff...  6ea08d5 
- Commits on May 25, 2020
  - fix a typo in Usage  4af8f2e 

## Fix usage msg and to accept yml with different names

Browse files

master

quishqa committed 13 days ago

1 parent 15c5186 commit 7a4849dbf75b679cb73f777b9a0f91b919f339ca

Showing 1 changed file with 4 additions and 2 deletions.

Unified Split

6 aas4wrf.py

...

```
@@ -264,11 +264,13 @@ def print_conservation(wrfchemi, emiss_input, pol):
264 264     import sys
265 265     import yaml
266 266     if len(sys.argv) < 2:
267 -     print('usage: python {} aasf4wrf.cfg'.format(sys.argv[0]))
267 +     print('usage: python {} aasf4wrf.yml'.format(sys.argv[0]))
268 268     sys.exit()
269 +
270 +     config_file = sys.argv[1]
269 271
270 272     # Retrieving parameters from config file
271 -     with open('aasf4wrf.yml') as file:
273 +     with open(config_file) as file:
272 274         config = yaml.load(file, Loader=yaml.FullLoader)
273 275
274 276     wrfout_file = config['Input']['wrfinput_file']
275
```

0 comments on commit 7a4849d

Lock conversation

- Un programa en python con git

- Un programa en python con git
- Se “subió” a github

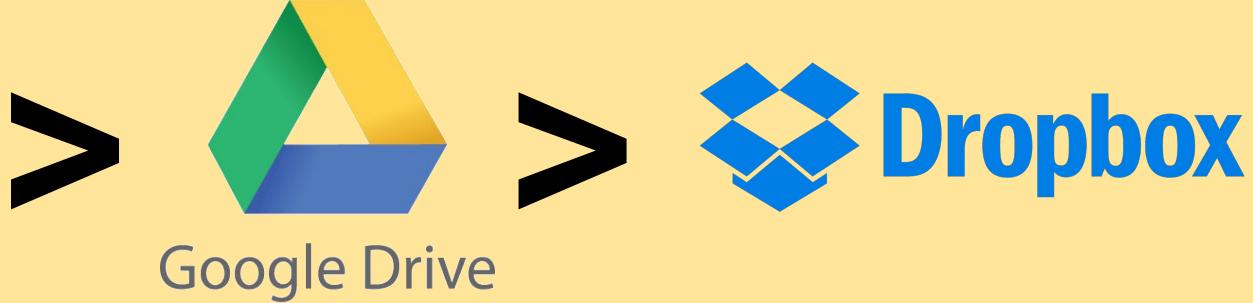
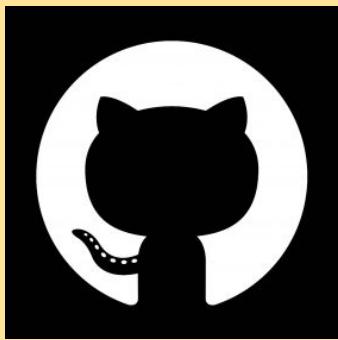
- Un programa en python con git
- Se “subió” a github
- Colegas pueden evaluar el  
método, siguiendo el README

- Un programa en python con git
- Se “subió” a github
- Colegas pueden replicar el método, siguiendo el README
- Pueden ver las correcciones

- Un programa en python con git
- Se “subió” a github
- Colegas pueden replicar el método, siguiendo el README
- Pueden ver las correcciones
- **Pueden crear copias, examinar, corregir y mejorar el programa.**

**Se replica, colabora y  
mejora :)**

(Último paréntesis personal)



# Resumiendo...

**Preferimos lenguajes  
científicos de  
programación.**

**Son gratis y  
versátiles. Hay una  
curva de aprendizaje,  
pero hay mucha ayuda  
online!**

**Importante en la  
reproducibilidad de la  
ciencia.**

# **Palabras finales**

**Por qué estoy aquí?**

**Invitarles a probar R o  
python, latex y git.  
Con paciencia y  
optimismo.**



Dude, sucking at something is the first step  
towards being sort of good at something

# Está presentación se ayudó de:

- Lowndes, J.S.S. et al. Our path to better science in less time using open data science tools. *Nat. Ecol. Evol.* 1, 0160 (2017).
- Introduction to Open Reproducible Science Workflows, Earth Lab Course, in: <https://www.earthdatascience.org/courses/intro-to-earth-data-science/open-reproducible-science/>
- Scopatz, A., & Huff, K. D. (2015). Effective computation in physics.
- Bartlett, Alice. 2016. “Git for Humans.” Financial Times, London; Talk at UX Brighton. (<https://speakerdeck.com/alicebartlett/git-for-humans>)
- Y varios memes del internet

# Lleve lleve casero!

- R (Junto con Rstudio)

<https://cran.r-project.org/bin/windows/base/>

- Python (Junto con Spyder o Jupyter)

<https://docs.conda.io/en/latest/miniconda.html>

- Latex

<https://www.overleaf.com>

# **gracias!**

[mario.calderon@iag.usp.br](mailto:mario.calderon@iag.usp.br)