# Rendez-vous de l'info scientifique

# Traitement de données avec Pandas & Jupyter notebooks

Pablo Iriarte – [pablo.iriarte@unige.ch](mailto:pablo.iriarte@unige.ch) / DIS

25 septembre 2024

**BIBLIOTHÈQUE**

UNIVERSITÉ DE GENÈVE

# Programme

## Introduction

- Historique
- Excel et les erreurs scientifiques
- Reproducibility Crisis & Data deluge

## Jupyter Notebooks

- Famille d'outils
- Accès au JupyterHub du cours ou installation via la distribution Anaconda
- Créer, organiser et partager des notebooks

## Pandas

- Importer et exporter des données
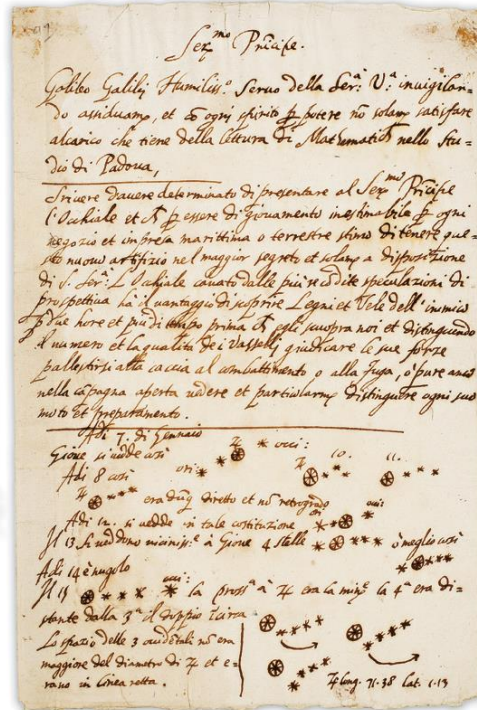- Manipuler et analyser les données
- Générer des graphiques

**BIBLIOTHÈQUE**

**UNIVERSITÉ DE GENÈVE**

# Introduction

Historique

- **iPython** (2001->) https://ipython.org/

- **Jupyter** (2014 ->) https://jupyter.org/

- **Famille d'outils**

  – **Jupyter Hub** : https://jupyterhub.readthedocs.io/

  – **Jupyter Lab** : https://jupyterlab.readthedocs.io/

  – **NB viewer** : https://nbviewer.jupyter.org/

  – **Binder** : https://mybinder.org/

UNIVERSITÉ DE GENÈVE

# Introduction

Historique : Jupiter



https://commons.wikimedia.org/wiki/File:Galileo_manuscript.png

UNIVERSITÉ DE GENÈVE

# Introduction

## Excel et les erreurs scientifiques

L'exemple du «Reinhart-Rogoff error»



*Reinhart, Carmen M., and Kenneth S. Rogoff. 2010. DOI:10.1257/aer.100.2.573*



https://www.nytimes.com/2013/04/19/opinion/krugman-the-excel-depression.html

UNIVERSITÉ DE GENÈVE

# Introduction

## Excel et les erreurs scientifiques

L'exemple de la conversion des données



Ziemann et al. 2016. DOI:10.1186/s13059-016-1044-7

UNIVERSITÉ DE GENÈVE

# Introduction

## **Reproductibilité et Open Science**

La science en crise?

1,500 scientists lift the lid on reproducibility

Baker 2016, Nature 533

https://doi.org/10.1038/533452a



IS THERE A REPRODUCIBILITY CRISIS?

7% Don't know

52% Yes, a significant crisis

3% No, there is no crisis

1,576 researchers surveyed

38% Yes, a slight crisis

©nature

UNIVERSITÉ DE GENÈVE

# Introduction

**Reproductibilité et
Open Science**

La science en crise?

1,500 scientists lift the lid on reproducibility

Baker 2016, Nature 533

https://doi.org/10.1038/533452a



WHAT FACTORS CONTRIBUTE TO IRREPRODUCIBLE RESEARCH?
Many top-rated factors relate to intense competition and time pressure.

- Always/often contribute · Sometimes contribute

Selective reporting
Pressure to publish
Low statistical power or poor analysis
Not replicated enough in original lab
Insufficient oversight/mentoring
Methods, code unavailable
Poor experimental design
Raw data not available from original lab
Fraud
Insufficient peer review
Problems with reproduction efforts
Technical expertise required for reproduction
Variability of standard reagents
Bad luck

©nature

**BIBLIOTHÈQUE**

**UNIVERSITÉ DE GENÈVE**

# Introduction

**Reproductibilité et Open Science**

Wired

https://www.wired.com/2017/04/want-fix-sciences-replication-crisis-replicate/



MEGAN MEYER OPINION 04.19.17 07:30 AM

WANT TO FIX SCIENCE'S REPLICATION CRISIS? THEN REPLICATE

UNIVERSITÉ DE GENÈVE

# Introduction

## Big Data et Open Data

Quantifying the Data Deluge and the Data Drought

https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2984851

Nombreux réservoirs ouverts

Kaggle : https://www.kaggle.com

Github: http://github.com

WikiData : https://www.wikidata.org

Statistique historique : https://hsso.ch

UNIVERSITÉ DE GENÈVE

# Introduction

**Excel : limitations**

Liste complète :

https://support.office.com/en-us/article/excel-specifications-and-limits-1672b34d-7043-467e-8e27-269d656771c3
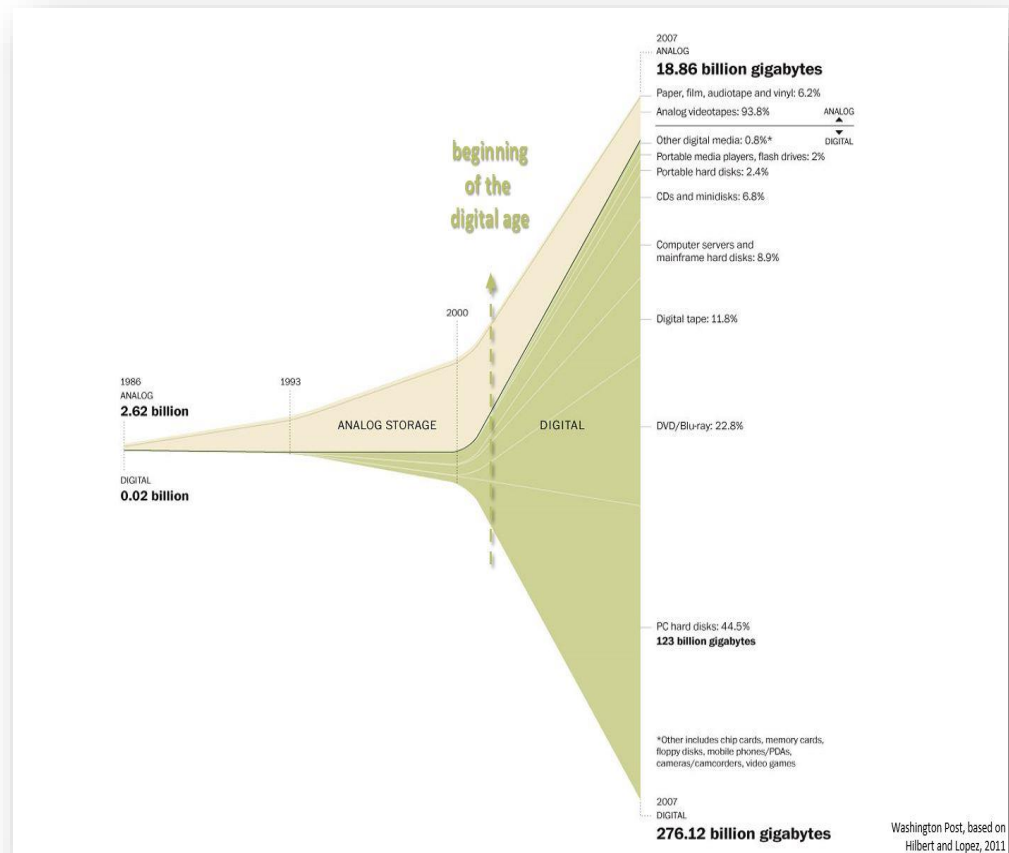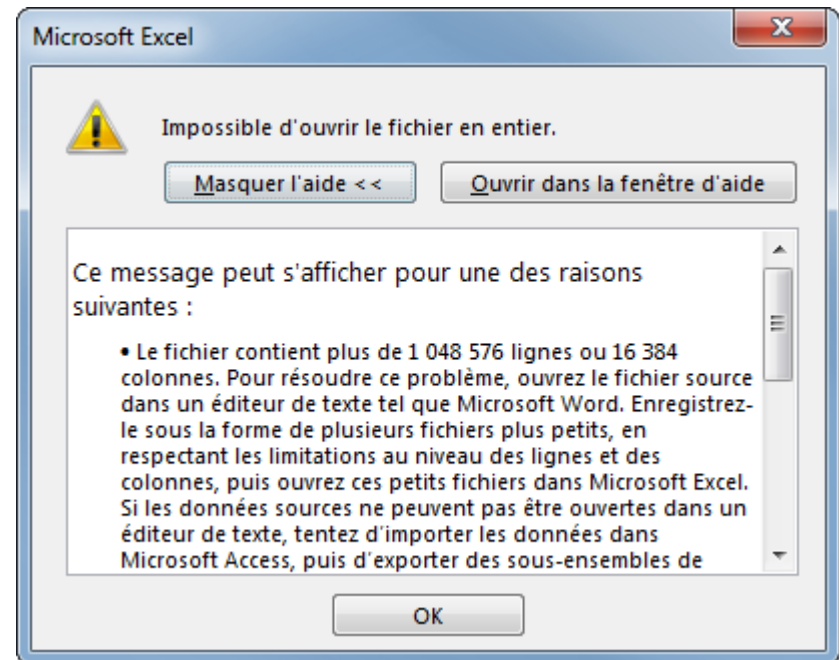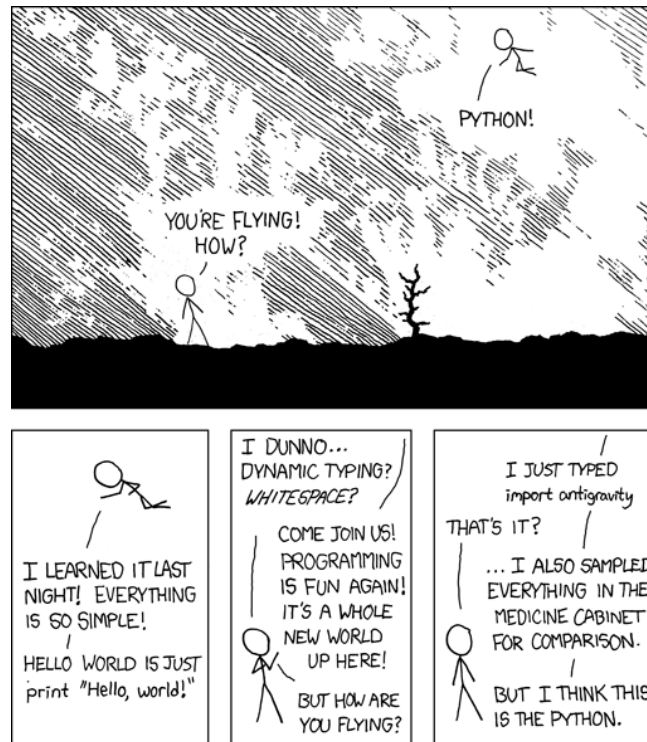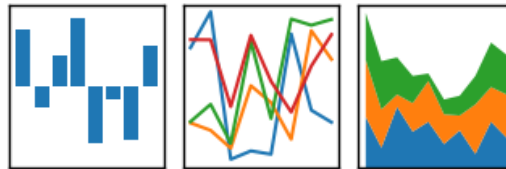
**Microsoft Excel**

⚠ Impossible d'ouvrir le fichier en entier.

[ Masquer l'aide << ]   [ Ouvrir dans la fenêtre d'aide ]

Ce message peut s'afficher pour une des raisons suivantes :

• Le fichier contient plus de 1 048 576 lignes ou 16 384 colonnes. Pour résoudre ce problème, ouvrez le fichier source dans un éditeur de texte tel que Microsoft Word. Enregistrez-le sous la forme de plusieurs fichiers plus petits, en respectant les limitations au niveau des lignes et des colonnes, puis ouvrez ces petits fichiers dans Microsoft Excel. Si les données sources ne peuvent pas être ouvertes dans un éditeur de texte, tentez d'importer les données dans Microsoft Access, puis d'exporter des sous-ensembles de

[ OK ]

**UNIVERSITÉ DE GENÈVE**

# Introduction



https://www.xkcd.com/353/

# Introduction

pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$

https://pandas.pydata.org

import pandas as pd

**BIBLIOTHÈQUE**

**UNIVERSITÉ DE GENÈVE**

# Introduction

**Reproductibilité et Open Science**

Nature

https://www.nature.com/articles/d41586-018-07196-1
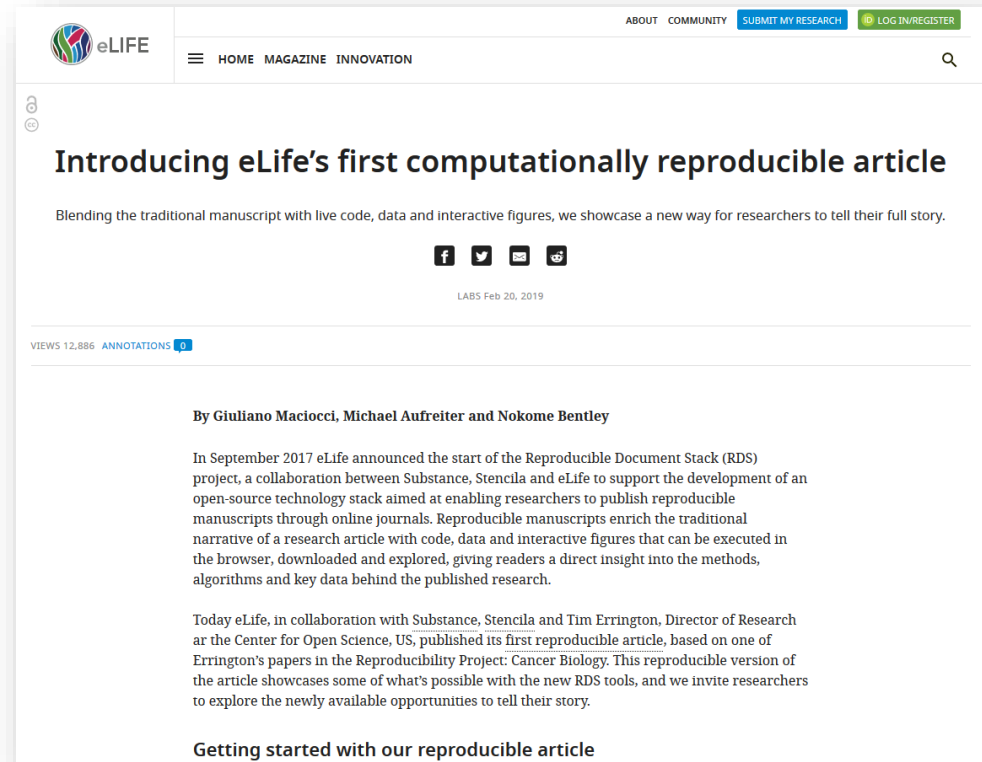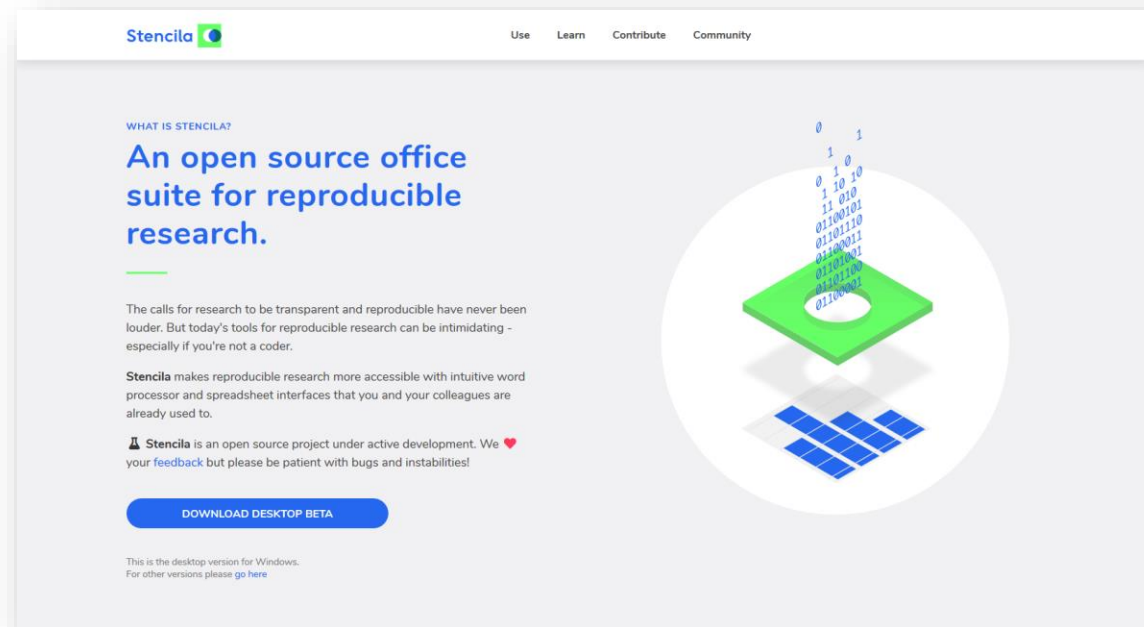
UNIVERSITÉ DE GENÈVE

# Introduction

## Reproductibilité et Open Science

eLife

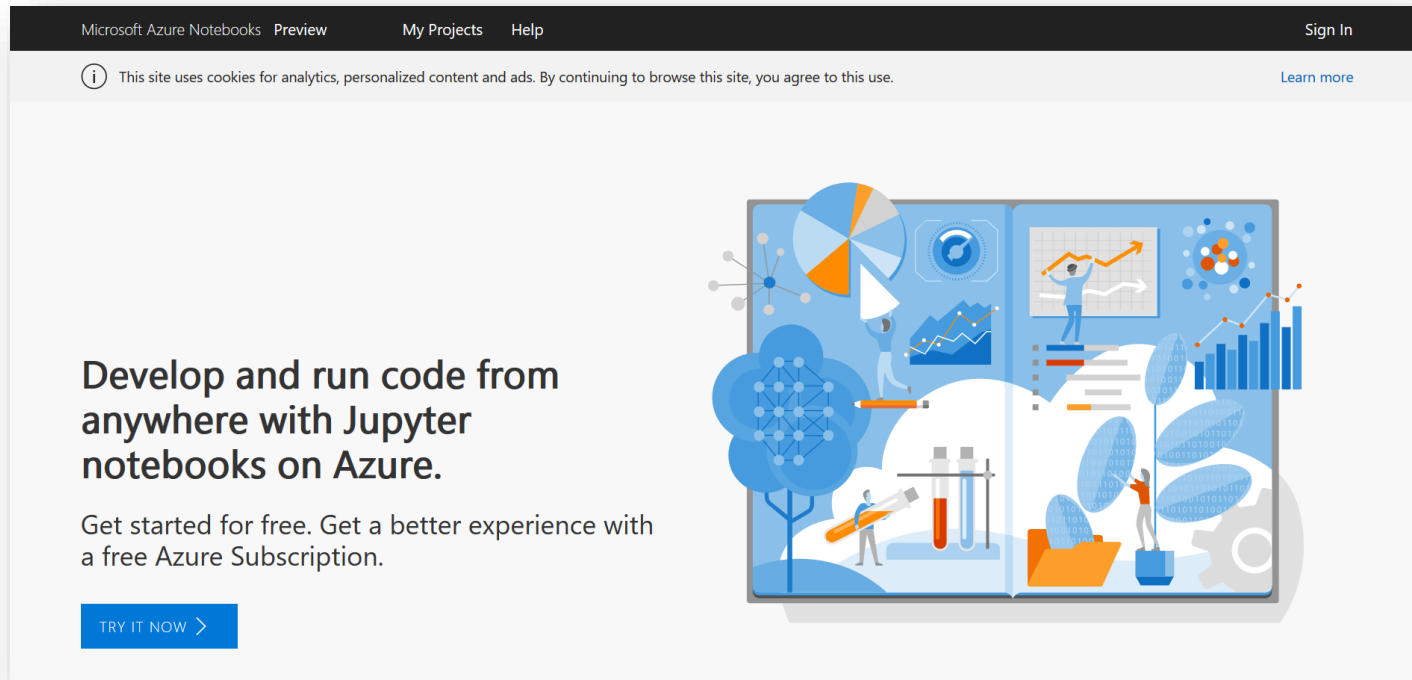https://elifesciences.org/labs/ad58f08d/introducing-elife-s-first-computationally-reproducible-article

UNIVERSITÉ DE GENÈVE

# Introduction

**Reproductibilité et Open Science**



https://stenci.la

UNIVERSITÉ
DE GENÈVE

# Introduction



https://notebooks.azure.com/

UNIVERSITÉ DE GENÈVE

# Introduction



https://colab.research.google.com

UNIVERSITÉ DE GENÈVE

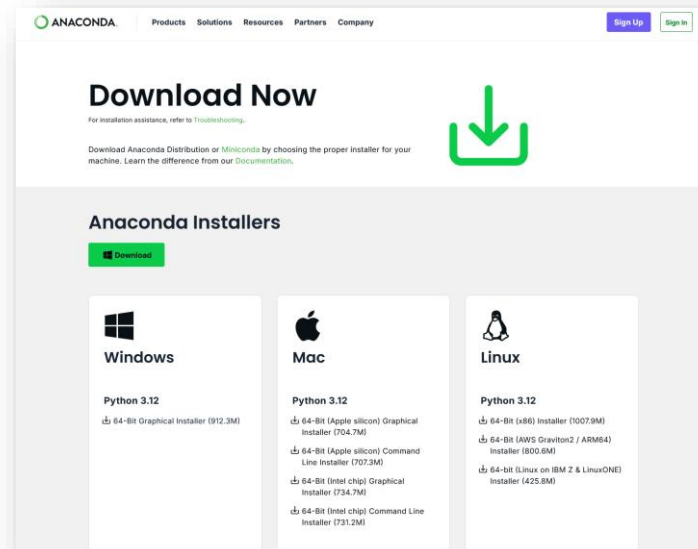# Jupyter Notebooks

**Installer Jupyter Notebooks et Pandas sur son poste** personnel avec la distribution « Anaconda » :
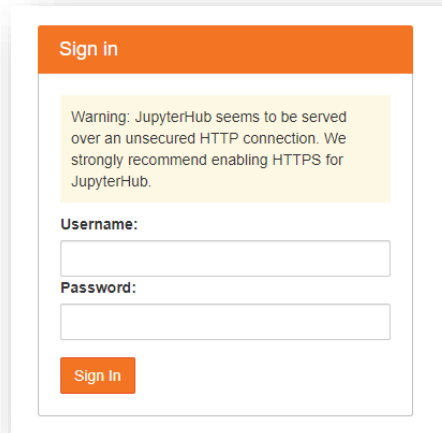
https://www.anaconda.com/download/

# Jupyter Notebooks

**Travail sur le JupyterHub du cours**

Se connecter sur cette adresse avec le login/pwd fourni pendant le cours :

## https://learning.jupyterhub.unige.ch

UNIVERSITÉ DE GENÈVE

# Jupyter Notebooks

Packages compris dans l'installation :

- Notebook (jupyter)

- Pandas

- NumPy

- Matplotlib

- NLTK

- …

Liste complète : https://docs.anaconda.com/anaconda/pkg-docs/

**UNIVERSITÉ DE GENÈVE**

# Jupyter Notebooks

**Créer, organiser et partager des notebooks**

Lancer Anaconda -> Jupyter Notebook

UNIVERSITÉ DE GENÈVE

# Jupyter Notebooks

**Si besoin :** créer un lien symbolique entre le « home » et le dossier avec les notebooks

1. Avec le shell se positionner sur le «home»
2. Créer le lien avec la commande :

    mklink /D Nom-du-lien Dossier-de-destination

Aide : https://www.howtogeek.com/howto/16226/complete-guide-to-symbolic-links-symlinks-on-windows-or-linux/

**BIBLIOTHÈQUE**

**UNIVERSITÉ DE GENÈVE**

# Jupyter Notebooks

**Se familiariser avec les notebooks**

Exercices

1. Ouvrir un notebook d'exemple (sur le dossier du cours)
2. Créer un nouveau notebook et le renommer
3. Ajouter une cellule de texte (markdown)
4. Ajouter une cellule de code python (calcul simple)
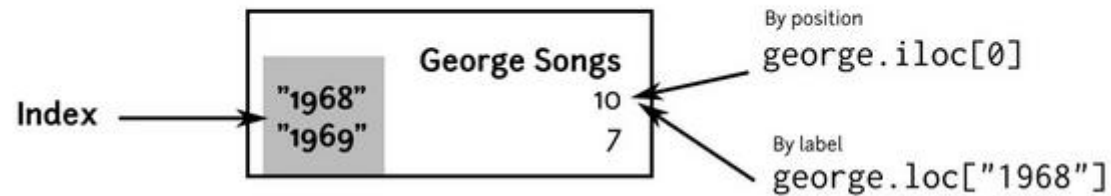5. L'exporter en format HTML

Aide markdown : https://guides.github.com/features/mastering-markdown/

Aide python : https://www.stavros.io/tutorials/python/

UNIVERSITÉ DE GENÈVE

# Pandas

Series : 1 dimension

# Pandas

Index : afficher des données par la position ou le nom de l'index

# Pandas

DataFrame : 2 dimensions

# Pandas

DataFrame : axes

# Pandas

DataFrame : slices

## Row & Column Slicing Examples

```
df.iloc[2:4, 0:1]     With a : return data frames
                      Position - Half-open interval
df.loc['d':, 'Units']  Without a : return series
                      Label - Closed interval
   Rows   Columns
```

**UNIVERSITÉ DE GENÈVE**

# Pandas

Opérations facilitées par les index : jointures automatiques

# Pandas

Opérations : GroupBy

UNIVERSITÉ
DE GENÈVE

# Pandas

Opérations : GroupBy

| Method | Result |
|---|---|
| .all | Boolean if all cells in group are True |
| .any | Boolean if any cells in group are True |
| .count | Count of non null values |
| .size | Size of group (includes null) |
| .idxmax | Index of maximum values |
| .idxmin | Index of minimum values |
| .quantile | Quantile (default of .5) of group |
| .agg(func) | Apply func to each group. If func returns scalar, then reducing |
| .apply(func) | Use split-apply-combine rules |
| .last | Last value |
| .nth | Nth row from group |
| .max | Maximum value |
| .min | Minimum value |
| .mean | Mean value |
| .median | Median value |
| .sem | Standard error of mean of group |
| .std | Standard deviation |
| .var | Variation of group |
| .prod | Product of group |
| .sum | Sum of group |

UNIVERSITÉ DE GENÈVE

# Pandas

Index multidimensionnels

UNIVERSITÉ
DE GENÈVE

# Pandas

Pivoter les tables

# Pandas

Jointures
(merge)



Visualizing Joins

UNIVERSITÉ
DE GENÈVE

# Pour avancer « pas à pas »



https://datacarpentry.org/python-socialsci/

# Pour aller plus loin



https://github.com/jupyter/jupyter/wiki#a-gallery-of-interesting-jupyter-notebooks

UNIVERSITÉ DE GENÈVE

# Pour aller plus loin

# Pandas

Exercices pratiques disponibles ici :
https://github.com/dis-unige/formations

1.  Importer des données
2.  Analyser des données
3.  Travailler avec différents types de données et des données manquantes
4.  Exporter des données
5.  Créer des graphiques simples

Aide Pandas : https://pandas.pydata.org/pandas-docs/stable/user_guide/10min.html

UNIVERSITÉ
DE GENÈVE

# Sources

Cheat Sheets distribués dans le cours :

- Jupyter notebook :

  https://www.datacamp.com/cheat-sheet/jupyter-notebook-cheat-sheet

- Markdown :

  http://geog.uoregon.edu/bartlein/courses/geog607/Rmd/MDquick-refcard.pdf

- Pandas :

  https://pandas.pydata.org/Pandas_Cheat_Sheet.pdf

**UNIVERSITÉ DE GENÈVE**