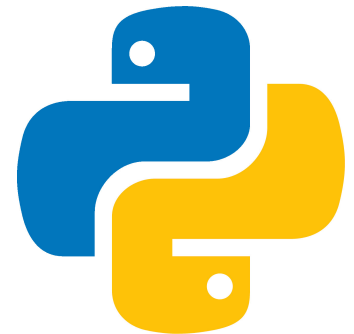


Outils génériques de visualisation de données pour Python

Par Vincent Roger

23/11/2023





À propos de moi

PhD en informatique.

Data scientist de Kiviak instrument et spécialiste de l'apprentissage automatique.

LinkedIn

<https://www.linkedin.com/in/vroger11/>

Site web

<https://website.vincent-roger.fr/>

Pourquoi la visualisation de données avec Python ?

- L'extraction et la préparation des données sont simplifiées par des outils tels que Pandas, Dask et d'autres.
- Il s'agit d'un langage incontournable pour les scientifiques des données, l'apprentissage automatique et les personnes travaillant sur les MLOPS.
- Alors pourquoi apprendre un nouveau langage juste pour la visualisation de données ?

Plan

- Non orienté web
- Orienté web

Données que nous allons utiliser

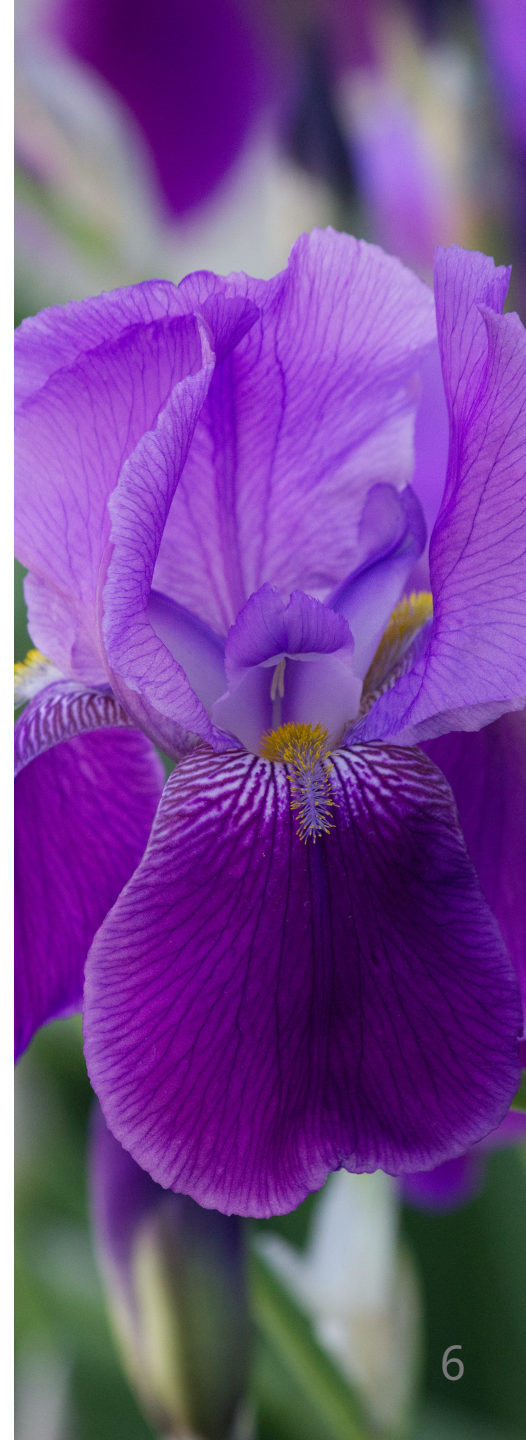
```
import pandas as pd

df_iris = pd.read_csv('https://raw.githubusercontent.com/mwaskom/seaborn-data/master/iris.csv')

df_iris
```

```
   sepal_length  sepal_width  petal_length  petal_width  species
0             5.1           3.5           1.4           0.2    setosa
1             4.9           3.0           1.4           0.2    setosa
2             4.7           3.2           1.3           0.2    setosa
3             4.6           3.1           1.5           0.2    setosa
4             5.0           3.6           1.4           0.2    setosa
..           ...           ...           ...           ...     ...
145            6.7           3.0           5.2           2.3  virginica
146            6.3           2.5           5.0           1.9  virginica
147            6.5           3.0           5.2           2.0  virginica
148            6.2           3.4           5.4           2.3  virginica
149            5.9           3.0           5.1           1.8  virginica
```

```
[150 rows x 5 columns]
```

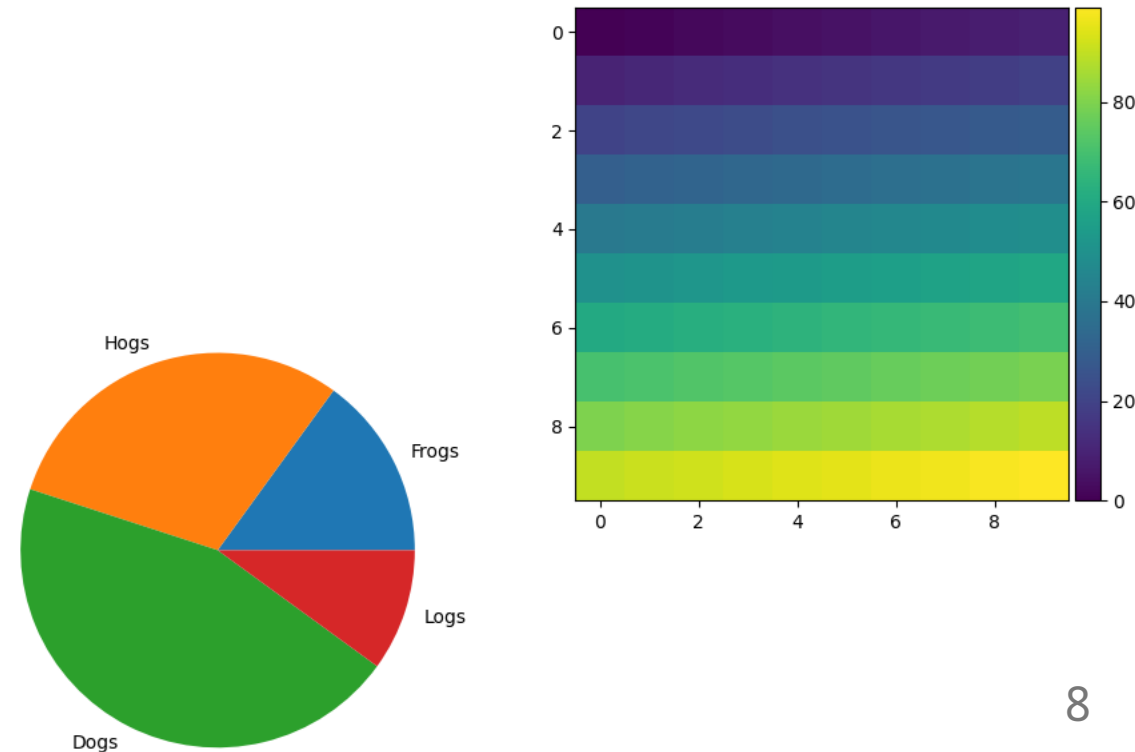
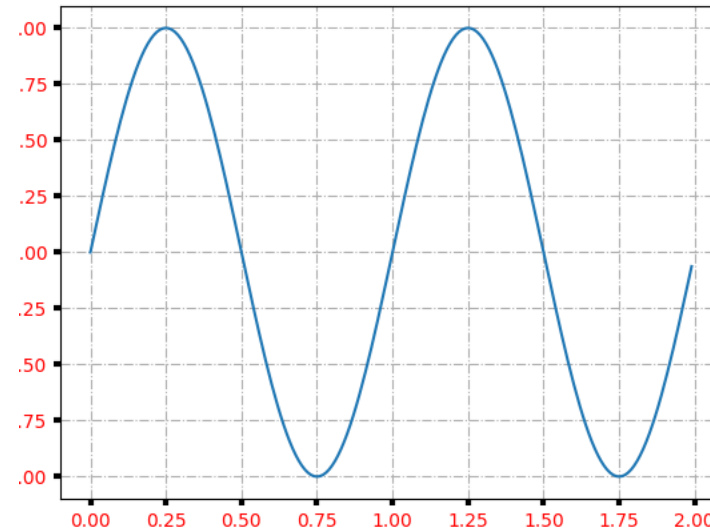


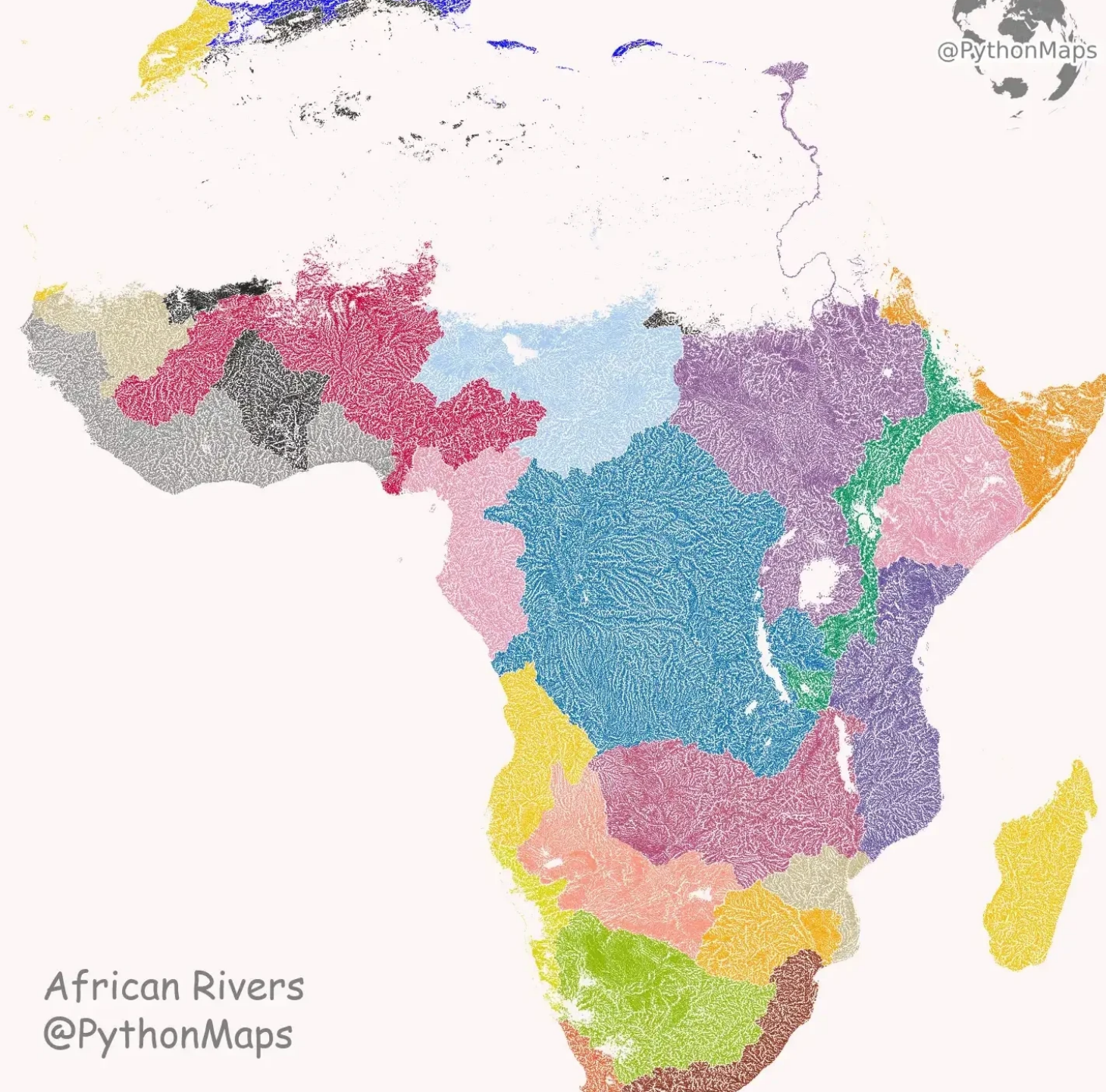


Non orienté web

Matplotlib

- Basé sur Qt.
- Bibliothèque standard.
- Par défaut pour de nombreux outils comme Pandas.





African Rivers
@PythonMaps

Matplotlib

- Puissant si maîtrisé !

Matplotlib - diagramme de dispersion des sépales de l'iris

```
import matplotlib.pyplot as plt

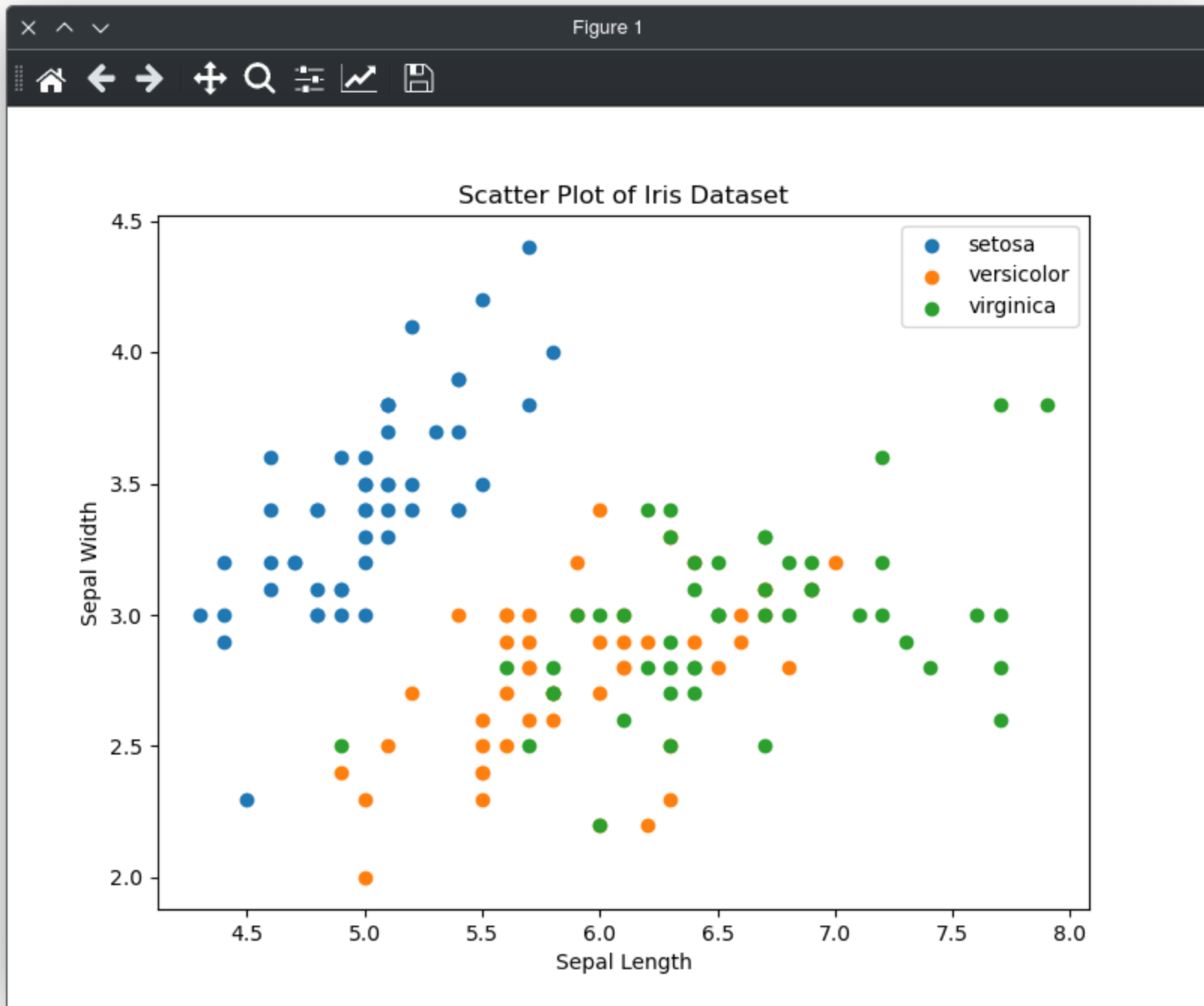
# Create a scatter plot with legends based on species
plt.figure(figsize=(8, 6))

# Iterate over each species and plot scatter points
for species in df_iris["species"].unique():
    species_data = df_iris[df_iris['species'] == species]
    plt.scatter(species_data['sepal_length'], species_data['sepal_width'], label=species)

# Add labels and title
plt.xlabel('Sepal Length')
plt.ylabel('Sepal Width')
plt.title('Scatter Plot of Iris Dataset')

# Add legend
plt.legend()

# Show the plot
plt.show()
```



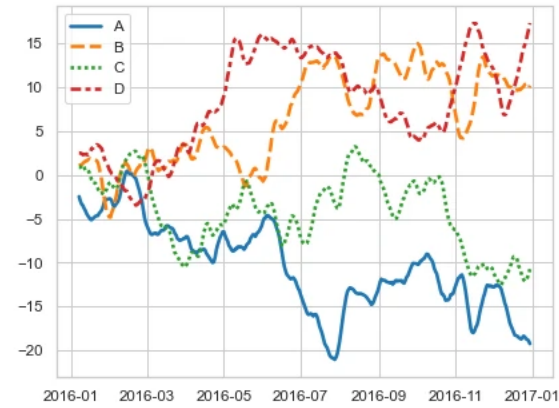
Matplotlib

Avantages :

- La librairie la plus utilisée.
- Travailler localement.
- Grande communauté.

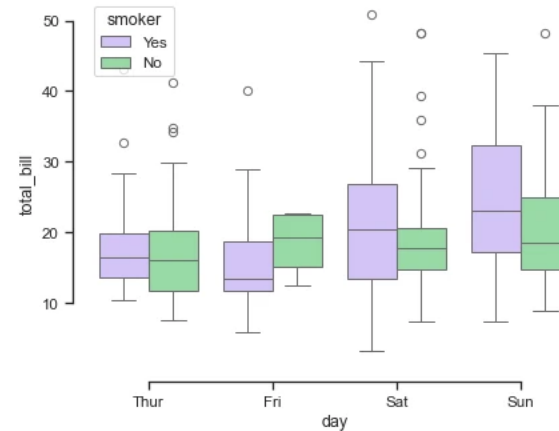
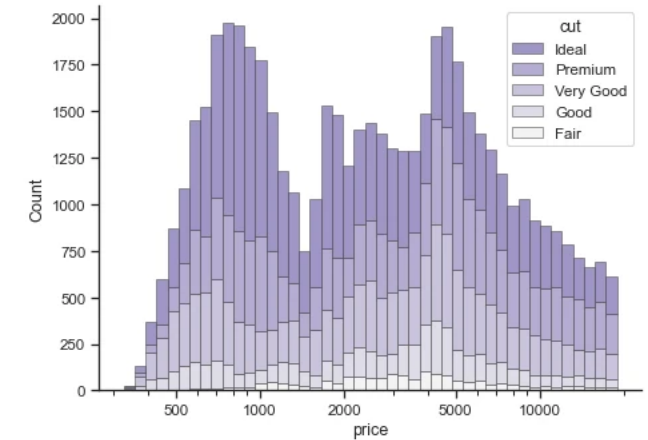
Inconvénients :

- Verbeux.
- N'est pas bien adapté aux tracés dynamiques (en particulier avec de nombreux points de données).
- Il peut être difficile de réaliser certains tracés spécifiques.



Seaborn

- Basé sur Matplotlib.
- Conçu pour les graphiques statistiques.



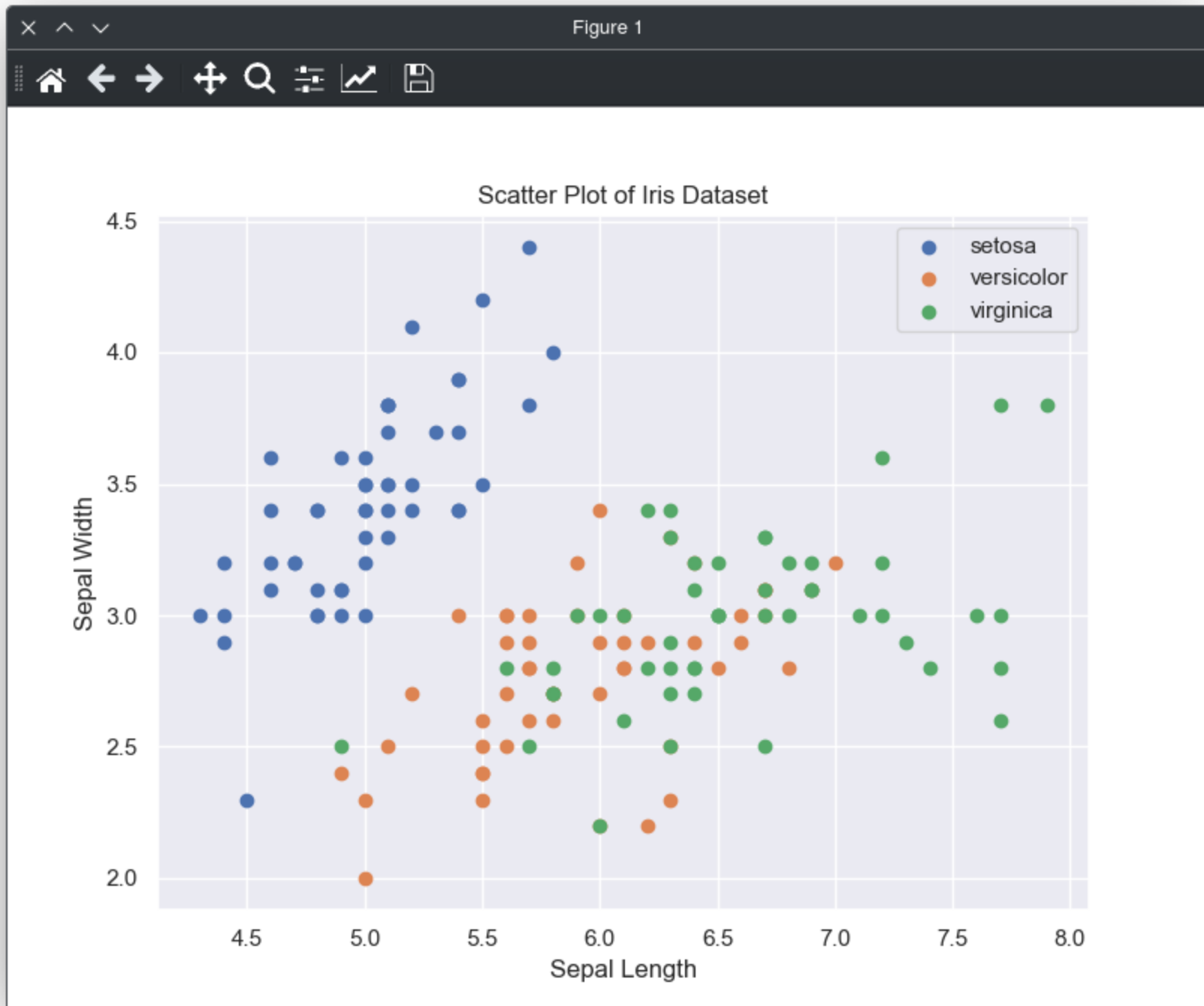
Seaborn - diagramme de dispersion des sépales de l'iris

```
import matplotlib.pyplot as plt
import seaborn as sns

# Create a scatter plot
sns.scatterplot(data=df_iris, x="sepal_length", y="sepal_width", hue="species")

# Add labels and title
plt.xlabel('Sepal length')
plt.ylabel('Sepal width')
plt.title('Scatter Plot of Sepal')

# Display the plot
plt.show()
```

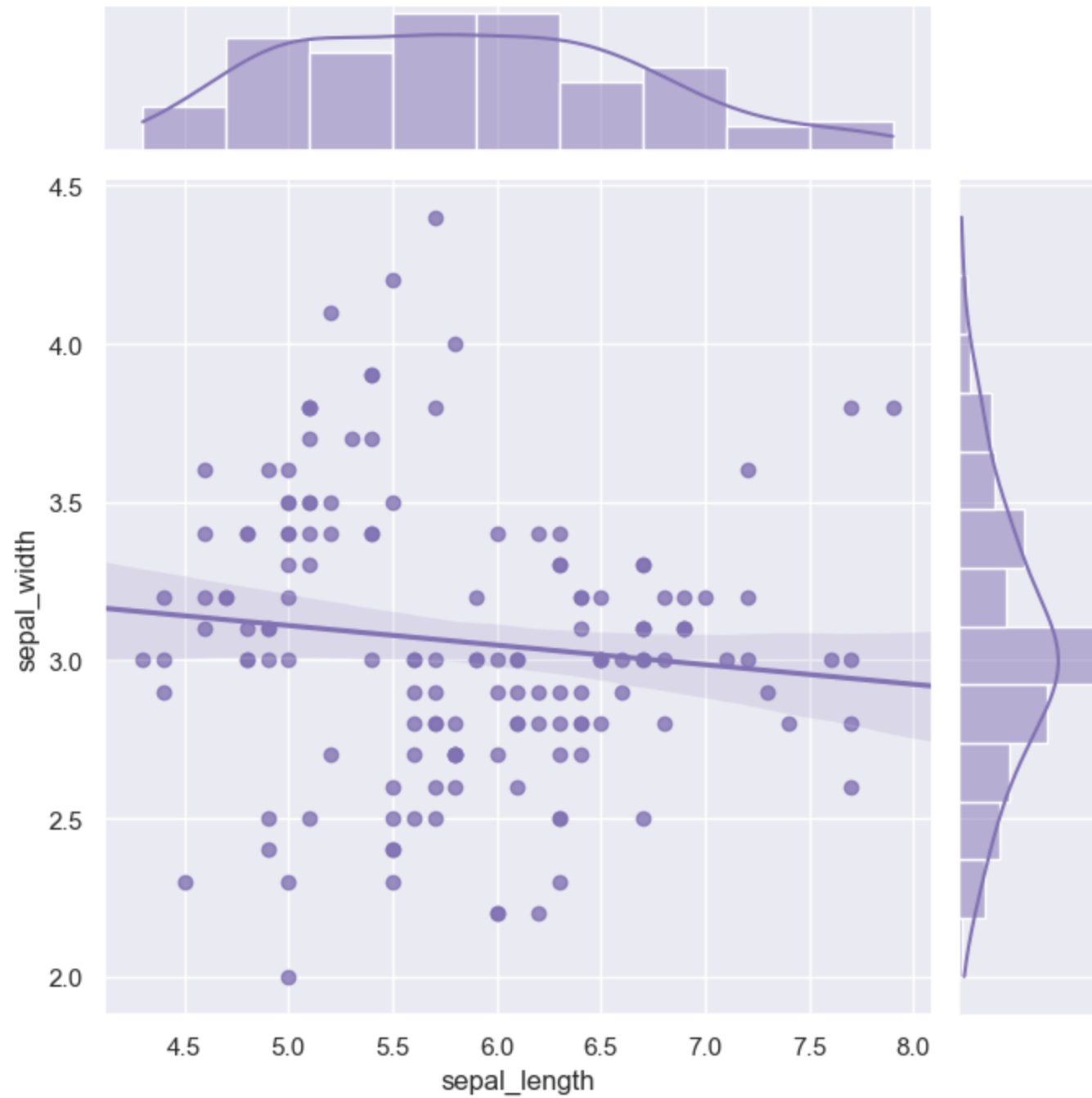



Seaborn - faire rapidement des statistiques

```
import matplotlib.pyplot as plt
import seaborn as sns

# Compute jointplot
sns.set_theme(style="darkgrid")
g = sns.jointplot(data=df_iris, x="sepal_length", y="sepal_width",
                  kind="reg", truncate=False,
                  color="m", height=7)

# Display the plot
plt.show()
```



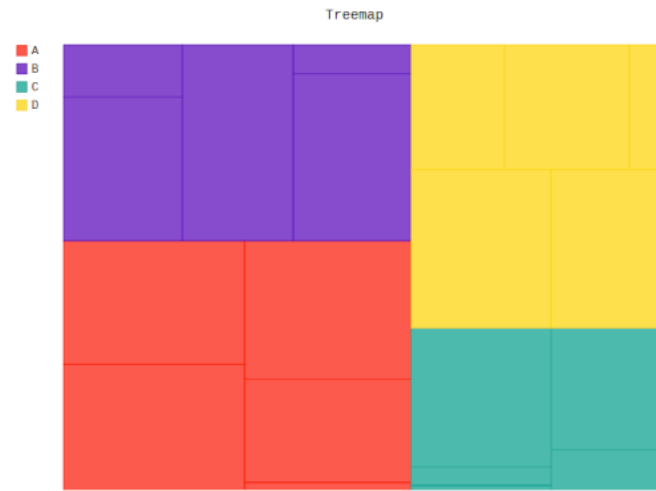
Seaborn

Avantages :

- Fournit des thèmes.
- Fournit des fonctions de haut niveau pour l'analyse statistique afin de faciliter certains tracés.

Inconvénients :

- Toujours un peu verbeux lors de la modification de certains tracés.
- N'est pas conçu pour les tracés dynamiques.
- Ne permet pas de générer tous les types de tracés.

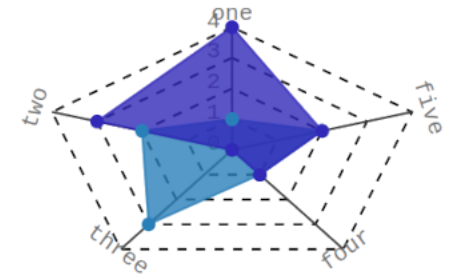


Pygal

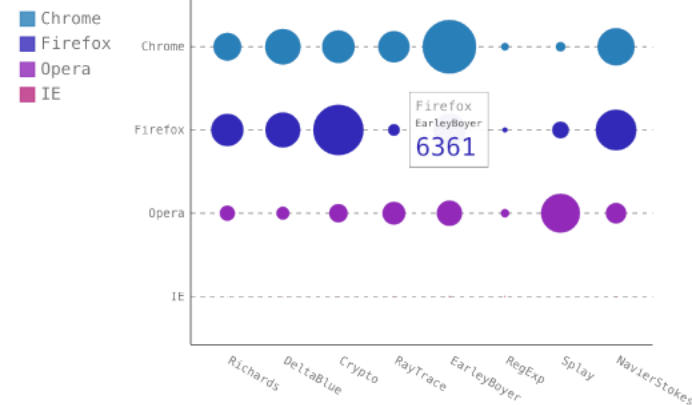
- Produire des SVG dynamiques.
- Écrit en Python.

Radar Chart with Fill

alpha
beta



V8 benchmark results



Pygal - diagramme de dispersion des sépales de l'iris

```
import pygal

scatter_chart = pygal.XY(stroke=False)
for specie in df_iris["species"].unique():
    df_s = df_iris[df_iris["species"] == specie]
    scatter_chart.add(f'{specie}', list(zip(df_s["sepal_length"], df_s["sepal_width"])))

scatter_chart.title = 'Sepal Scatter Plot Example'
scatter_chart.x_title = 'Sepal length'
scatter_chart.y_title = 'Sepal width'
scatter_chart.render_to_file('scatterPlot.svg')
```


Pygal

Avantages

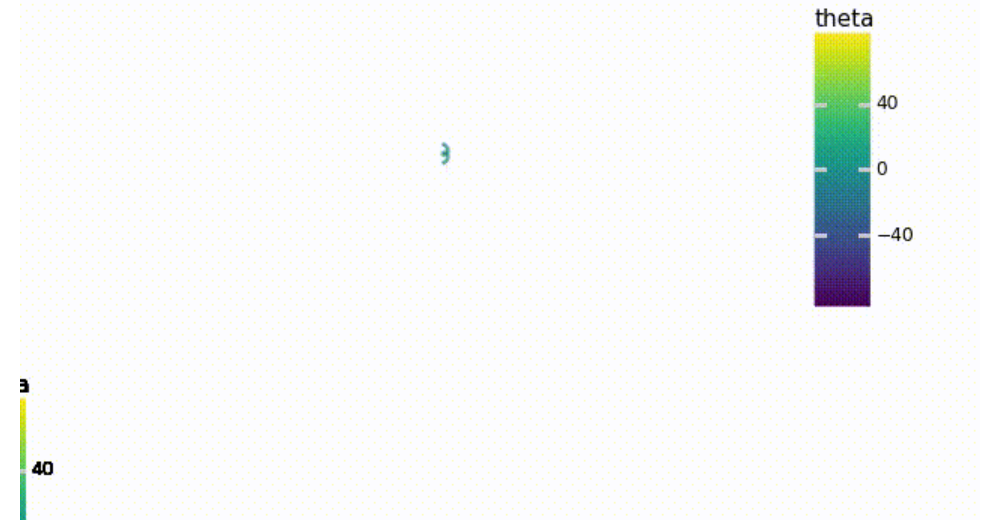
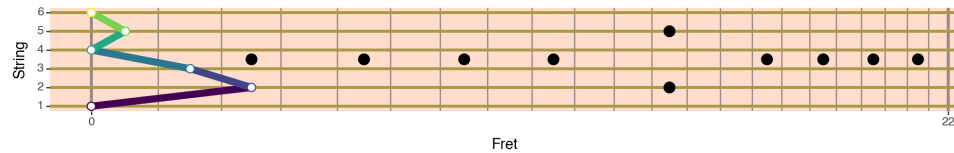
- Même si le résultat est un fichier SVG, il est interactif par défaut !
- Grande variété de graphiques.

Inconvénients

- SVG compatible avec les navigateurs web, mais pas avec des logiciels comme Inkscape ou notebook (même avec rendu HTML).
- Les graphiques ne sont pas clairs avec les valeurs par défaut (en particulier la grille).

Plotnine

- Implémentation de la grammaire de ggplot2 en Python.



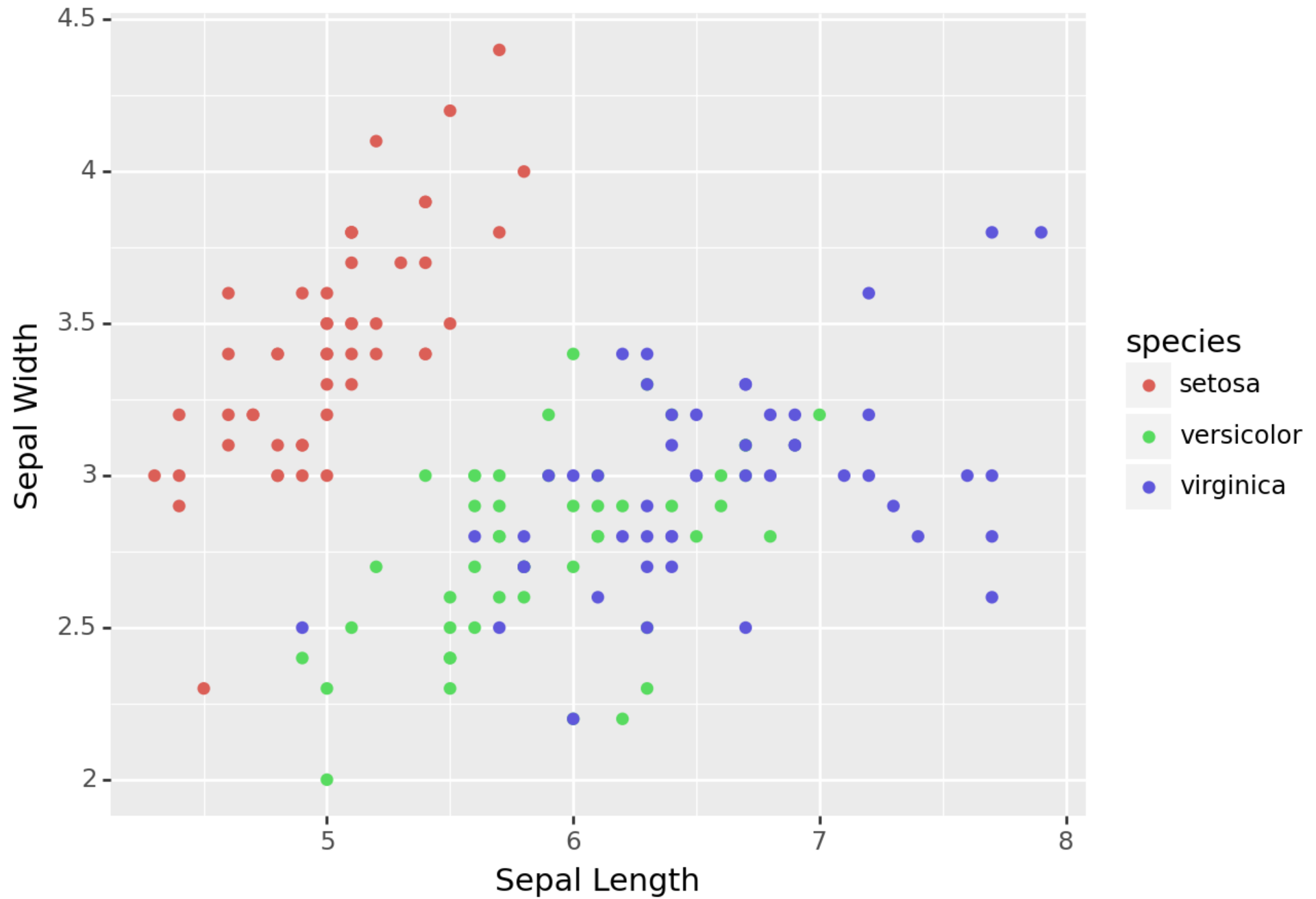
Plotnine - diagramme de dispersion des sépales de l'iris

```
from plotnine import ggplot, aes, geom_point, labs

# Create a scatter plot of sepal size
scatter_plot = (
    ggplot(df_iris, aes(x='sepal_length', y='sepal_width', color='species')) +
    geom_point() +
    labs(title='Iris Dataset Scatter Plot', x='Sepal Length', y='Sepal Width')
)

# Display the plot
print(scatter_plot)
```


Scatter Plot of Sepal Size



Plotnine

Avantages

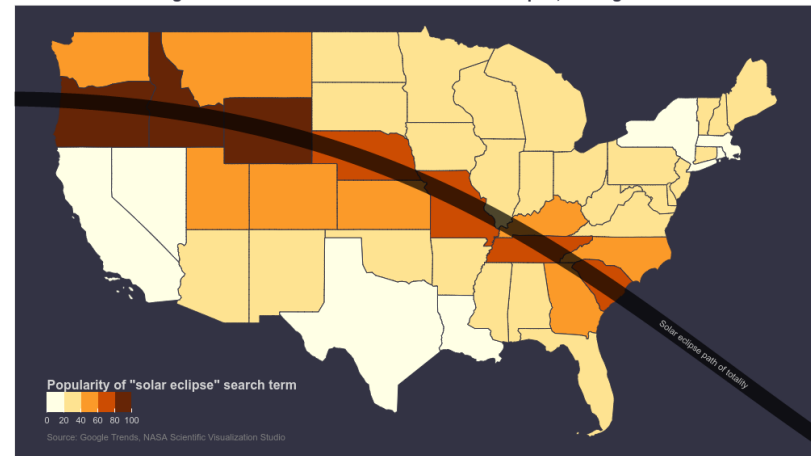
- Aussi puissant que ggplot2.
- Outil adapté aux utilisateurs de R utilisant Python.
- Grande variété de tracés possibles.

Inconvénients

- La syntaxe n'est pas Python.
- Ne permet pas de tracer de grands ensembles de données.
- L'interactivité des tracés est limitée.



Orienté Web

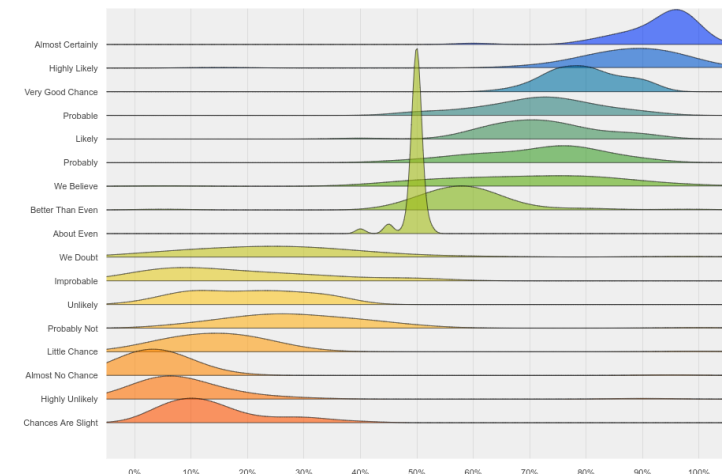


Bokeh

- Utilisable dans Jupyter et IPython (sortie HTML par défaut).
- Capacité de créer des tableaux de bord.
- L'utilisation de WebGL est possible.

Periodic Table (omitting LA and AC Series)

 nonmetal noble gas alkali metal alkaline earth metal metalloid halogen metal transition metal																																
I																	2															
1																	2															
H 1.00794 Hydrogen																	He 4.002602 Helium															
II																	10															
3	4																	5	6	7	8	9	10									
Li 6.941 Lithium	Be 9.012182 Beryllium																	B 10.811 Boron	C 12.0107 Carbon	N 14.0067 Nitrogen	O 15.9994 Oxygen	F 18.9984032 Fluorine	Ne 20.1797 Neon									
III																	17	18														
11	12																	13	14	15	16	17	18									
Na 22.98976928 Sodium	Mg 24.305 Magnesium																	Al 26.9815386 Aluminum	Si 28.086 Silicon	P 30.973762 Phosphorus	S 32.065 Sulfur	Cl 35.453 Chlorine	Ar 39.948 Argon									
IV																	29	30	31	32	33	34	35	36								
19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36															
K 39.0983 Potassium	Ca 40.078 Calcium	Sc 44.955912 Scandium	Ti 47.867 Titanium	V 50.9415 Vanadium	Cr 51.9961 Chromium	Mn 54.938043 Manganese	Fe 55.845 Iron	Co 58.933195 Cobalt	Ni 58.6934 Nickel	Cu 63.546 Copper	Zn 65.38 Zinc	Ga 69.723 Gallium	Ge 72.64 Germanium	As 74.92160 Arsenic	Se 78.96 Selenium	Br 79.904 Bromine	Kr 83.795 Krypton															
V																	47	48	49	50	51	52	53	54								
37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54															
Rb 85.4678 Rubidium	Sr 87.62 Strontium	Y 88.90585 Yttrium	Zr 91.224 Zirconium	Nb 92.90638 Niobium	Mo 95.96 Molybdenum	Tc [98] Technetium	Ru 101.07 Ruthenium	Rh 102.90550 Rhodium	Pd 106.42 Palladium	Ag 107.8682 Silver	Cd 112.411 Cadmium	In 114.818 Indium	Sn 117.710 Tin	Sb 121.760 Antimony	Te 127.60 Tellurium	I 126.90447 Iodine	Xe 131.293 Xenon															
VI																	79	80	81	82	83	84	85	86								
55	56																	79	80	81	82	83	84	85	86							
Cs 132.9054519 Cesium	Ba 137.327 Barium																	Hf 178.49 Hafnium	Ta 180.94736 Tantalum	W 183.84 Tungsten	Re 186.207 Rhenium	Os 190.23 Osmium	Ir 192.222 Iridium	Pt 195.084 Platinum	Au 196.966569 Gold	Hg 200.59 Mercury	Tl 204.3833 Thallium	Pb 207.2 Lead	Bi 208.98040 Bismuth	Po [209] Polonium	At [210] Astatine	Rn [222] Radon
VII																	104	105	106	107	108	109	110	111	112	113	114	115	116	117	118	
87	88																	104	105	106	107	108	109	110	111	112	113	114	115	116	117	118
Fr [223] Francium	Ra [226] Radium																	Rf [261] Rutherfordium	Db [262] Dubnium	Sg [263] Seaborgium	Bh [264] Bohrium	Hs [265] Hassium	Mt [266] Meitnerium	Ds [267] Darmstadtium	Rg [268] Roentgenium	Cn [269] Copernicium	Nh [270] Nihonium	Fl [271] Flerovium	Mc [272] Moscovium	Lv [273] Livermorium	Ts [274] Tennessine	Og [276] Oganesson



Bokeh - diagramme de dispersion des sépales de l'iris

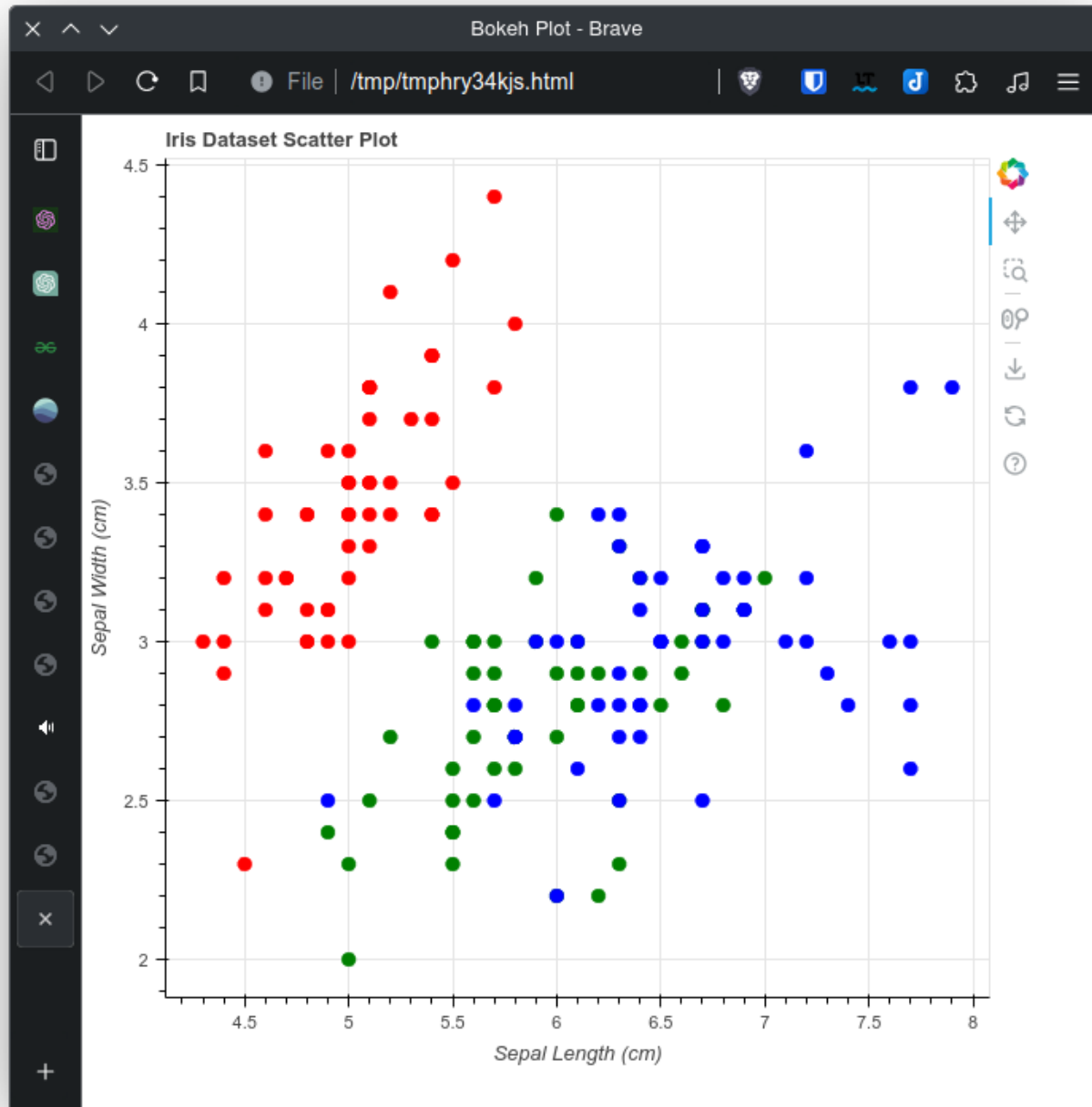
```
from bokeh.plotting import figure, show
from bokeh.transform import factor_cmap
from bokeh.models import ColumnDataSource

# Create a ColumnDataSource
source = ColumnDataSource(df_iris)

# Create a scatter plot with titles
p = figure(title='Iris Dataset Scatter Plot', x_axis_label='Sepal Length (cm)', y_axis_label='Sepal Width (cm)')

p.circle(x='sepal_length', y='sepal_width', source=source,
         size=8,
         color=factor_cmap(
             'species',
             palette=['red', 'green', 'blue'],
             factors=df_iris["species"].unique()
         )
        )

# Show the plot
show(p)
```

Bokeh

Avantages

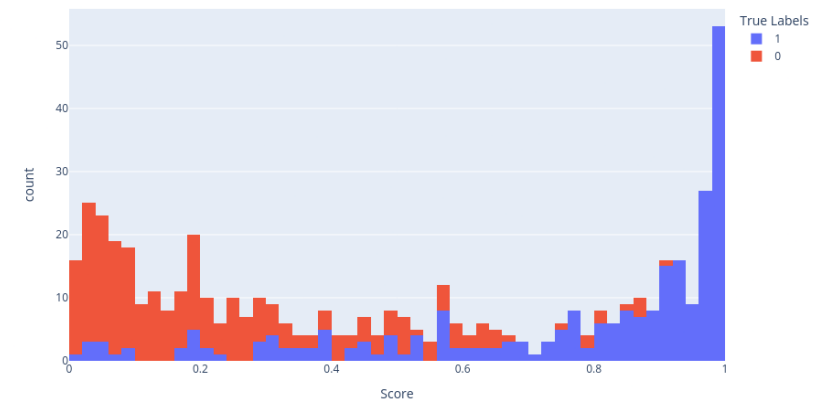
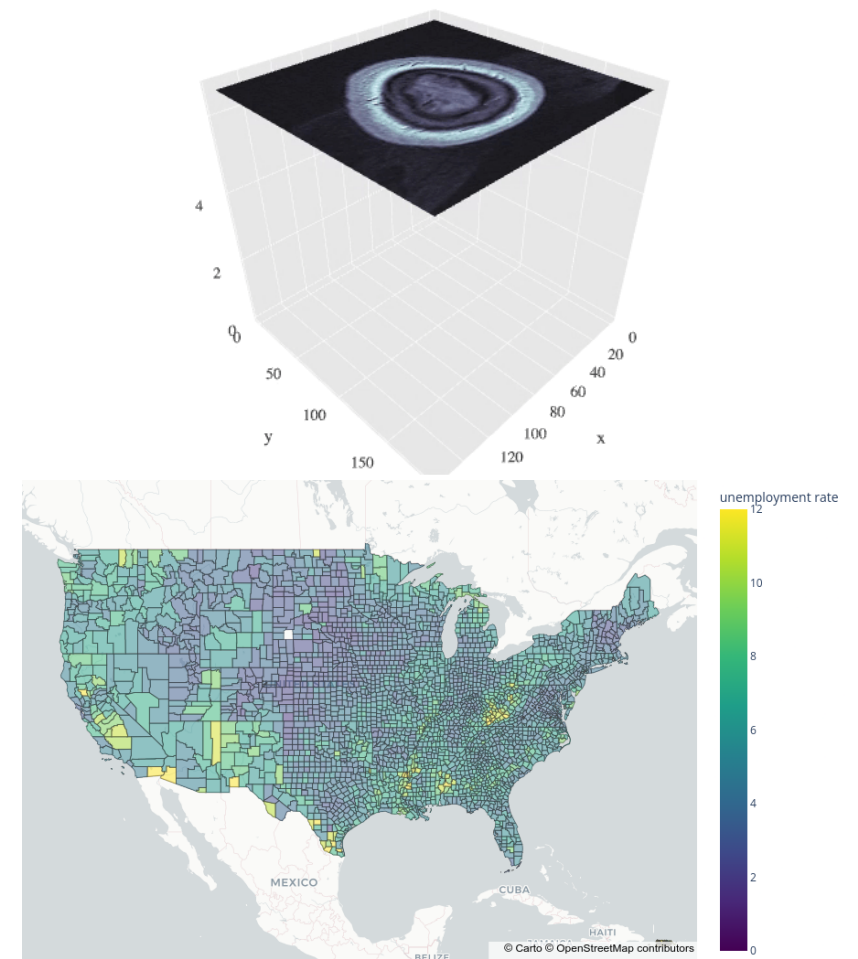
- API de bas niveau et de haut niveau.
- Peut facilement gérer un million de points.
- Interactivité (widgets et tableaux de bord).
- Capacité à réaliser de nombreux tracés différents.

Inconvénients

- Plus difficile à apprendre que d'autres bibliothèques.
- Moins complet pour les tracés en 3D.
- Si les serveurs de Bokeh sont hors service, vos tracés peuvent ne pas s'afficher correctement (ou pas du tout).

Plotly

- Conçu pour les visualisations interactives.
- L'utilisation de WebGL est possible.
- Utilisable avec Python, R et JavaScript.



Plotly - diagramme de dispersion des sépales de l'iris

```
import plotly.express as px

fig = px.scatter(df_iris, x='sepal_length', y='sepal_width', color='species',
                 title='Iris Dataset Scatter Plot', labels={'sepal_length': 'Sepal Length', 'sepal_width': 'Sepal Width'})

# Show the plot
fig.show()
```



Plotly

Avantages

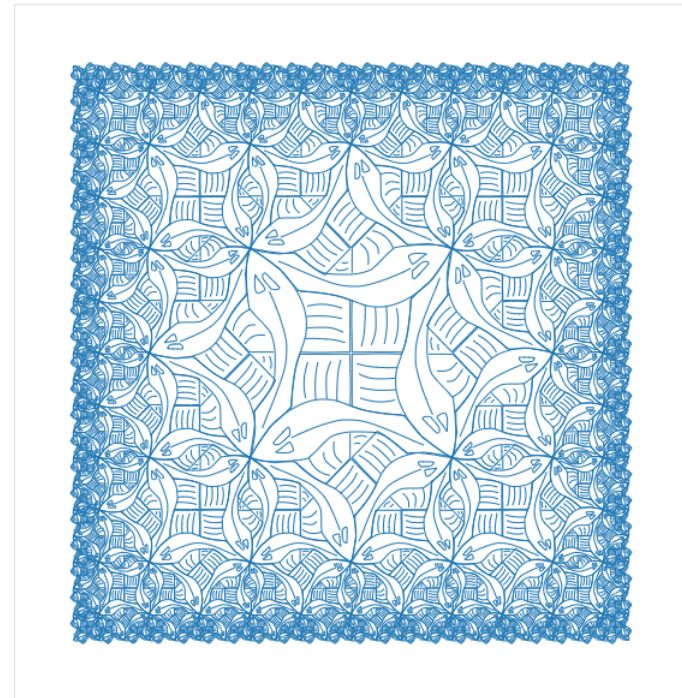
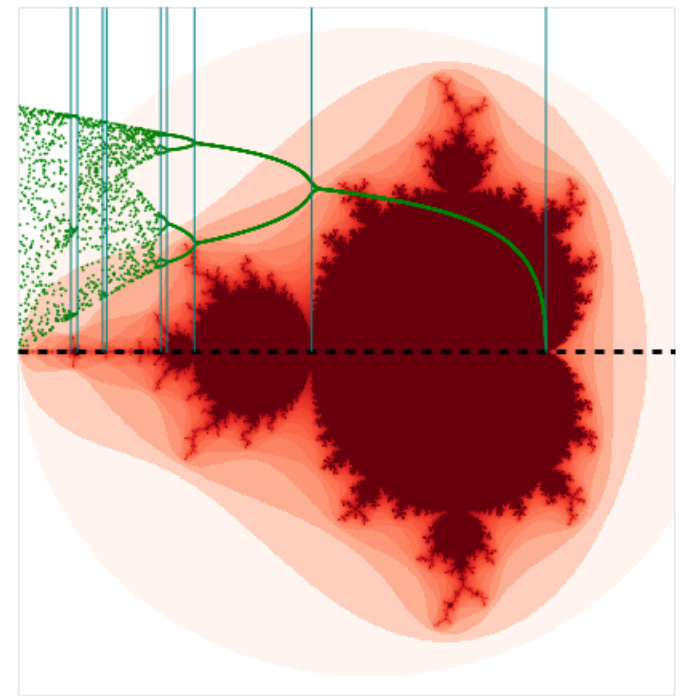
- La documentation est claire comme de l'eau de roche pour les tracés simples.
- Tableau de bord possible avec Dash (bien intégré).
- Redimensionnement automatique du graphique par défaut.

Inconvénients

- Si les serveurs de Plotly sont en panne, vos tracés peuvent ne pas s'afficher correctement (ou pas du tout).
- Difficile d'utiliser des images comme tiques.

Holoviews

Une bibliothèque de haut niveau pour Bokeh, Matplotlib et Plotly.



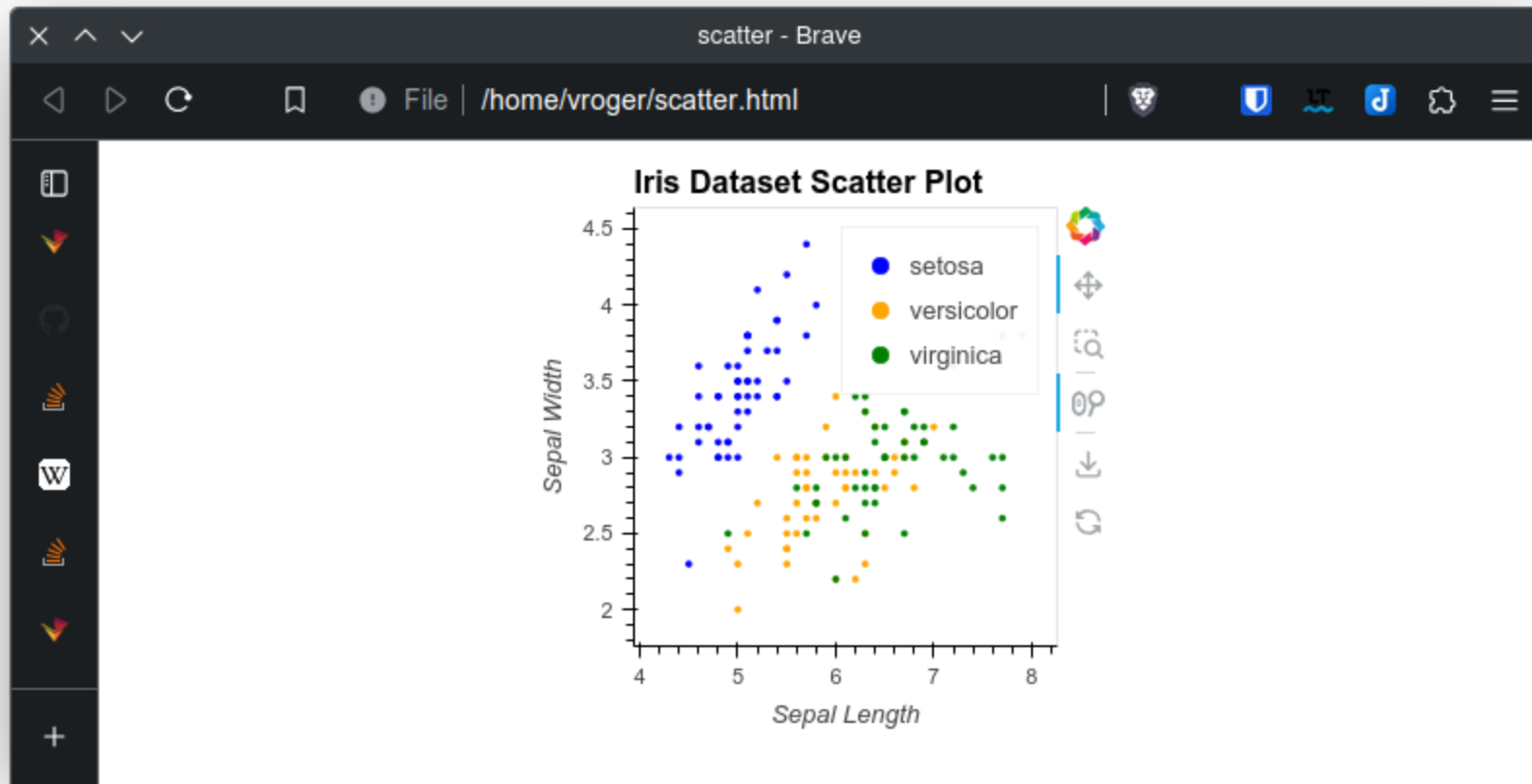
Holoviews - diagramme de dispersion des sépales de l'iris

```
import holoviews as hv

scatter = hv.Scatter(
    data=df_iris, kdims=['sepal_length'], vdims=['sepal_width', 'species'],
).opts(color='species', cmap=['blue', 'orange', 'green'])

scatter = scatter.opts(
    title="Iris Dataset Scatter Plot",
    xlabel='Sepal Length',
    ylabel='Sepal Width',
)

hv.save(scatter, 'scatter.html', backend='bokeh')
```



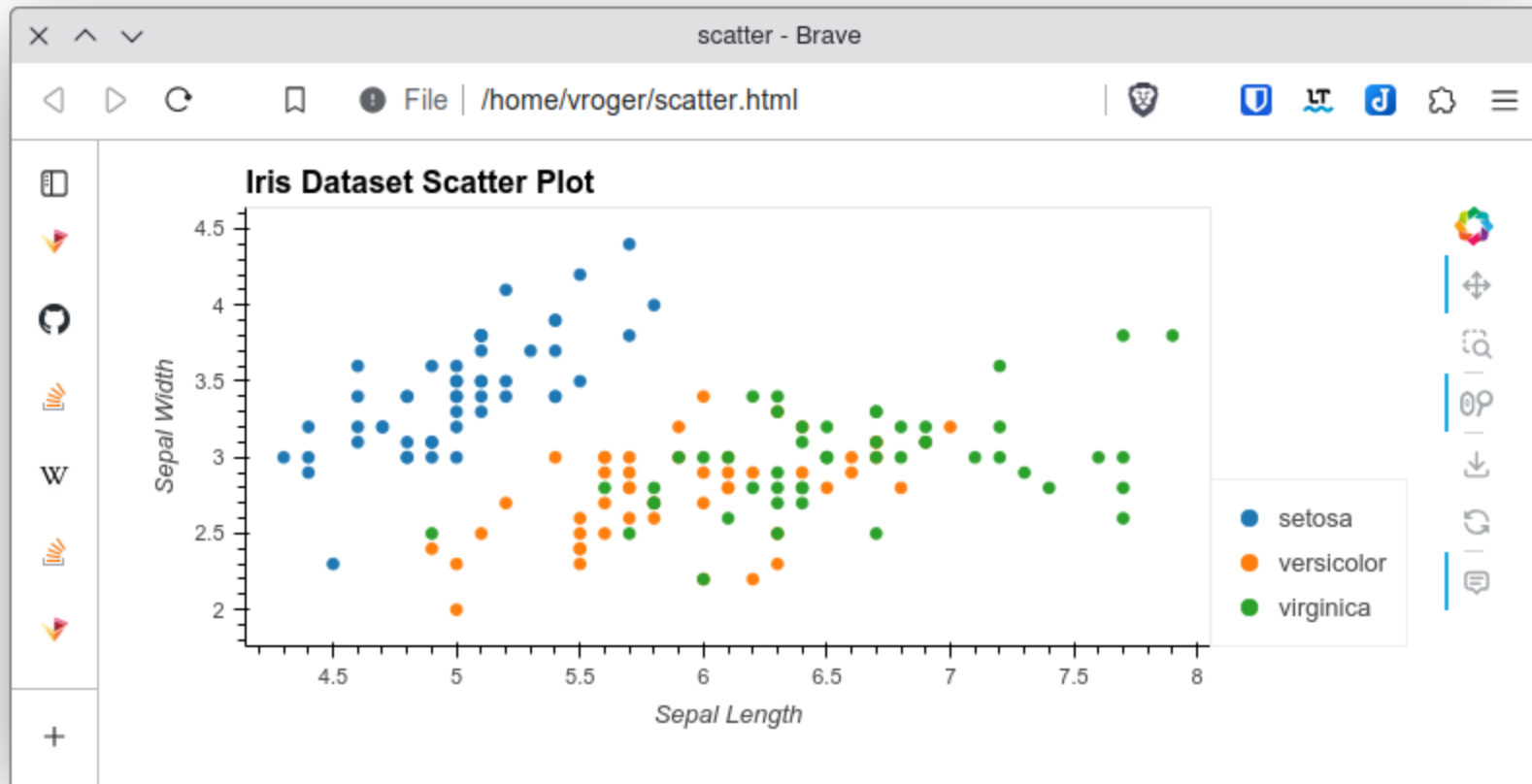
Holoviews - diagramme de dispersion des sépales de l'iris avec l'API Pandas

```
import hvplot
import hvplot.pandas
import holoviews as hv

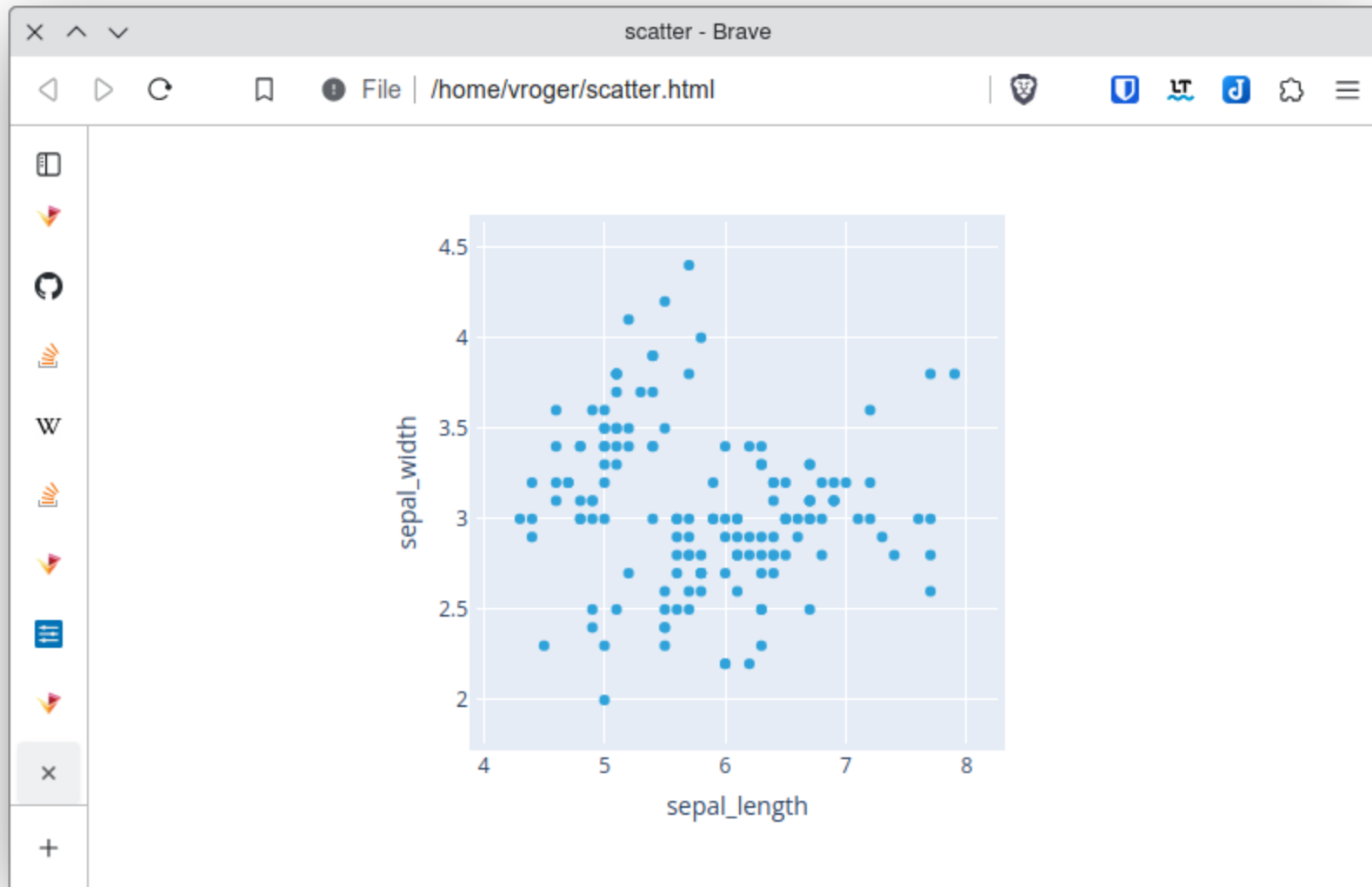
scatter = df_iris.hvplot(kind='scatter', x='sepal_length', y='sepal_width', color='species')

scatter = scatter.opts(
    title="Iris Dataset Scatter Plot",
    xlabel='Sepal Length',
    ylabel='Sepal Width',
)

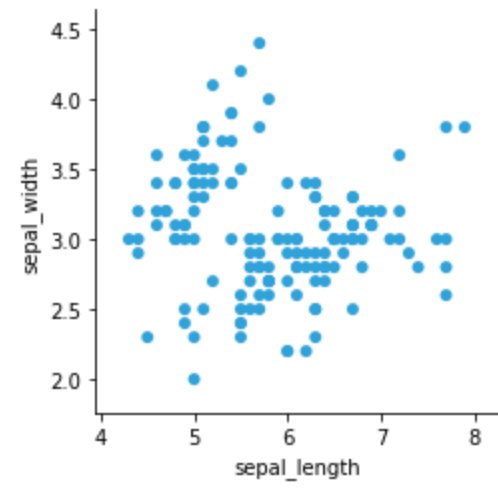
hv.save(scatter, 'scatter.html', backend='bokeh')
```

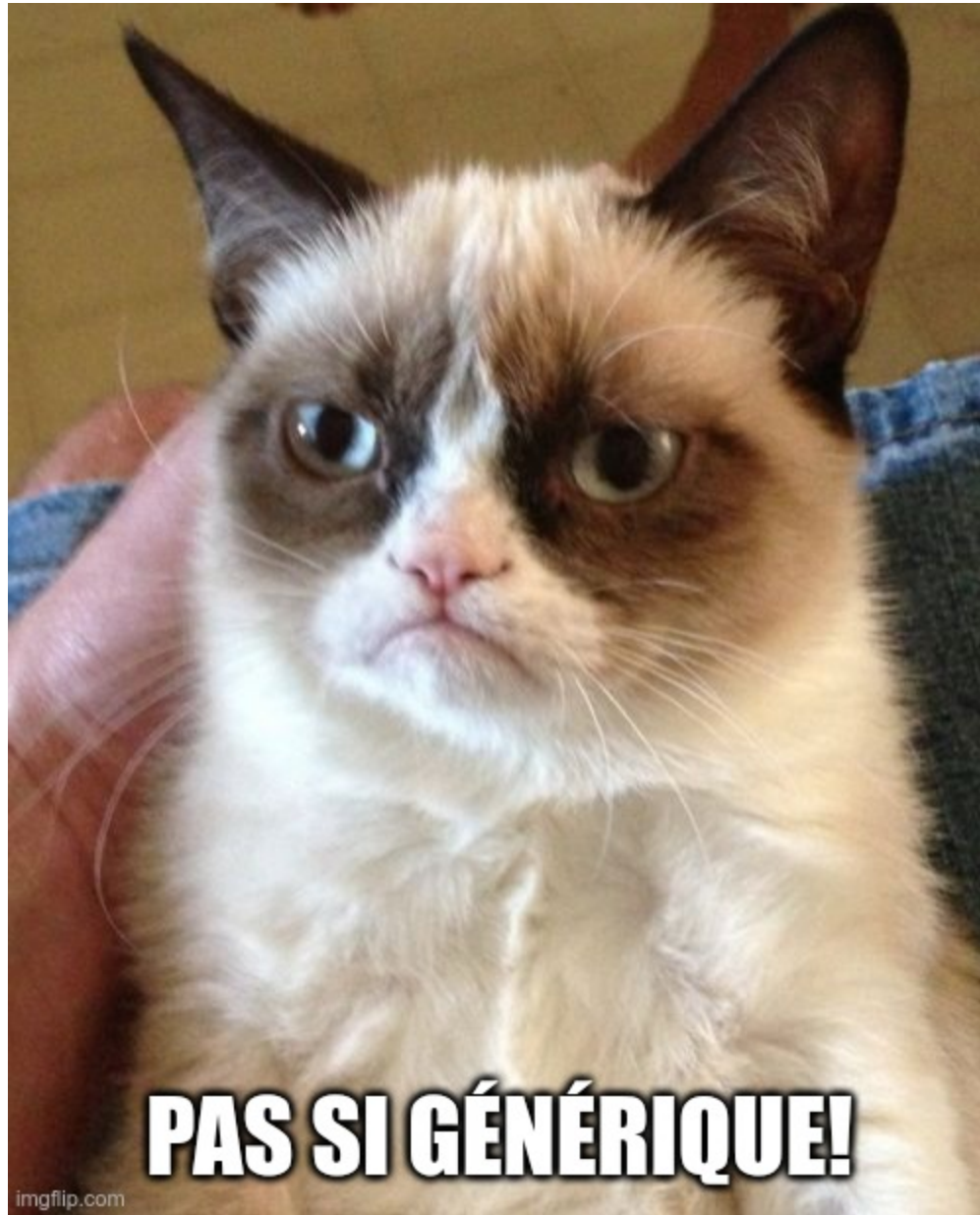


```
hv.save(scatter, 'scatter.html', backend='plotly')
```




```
hv.save(scatter, 'scatter.png', backend='matplotlib')
```





Holoviews

Avantages

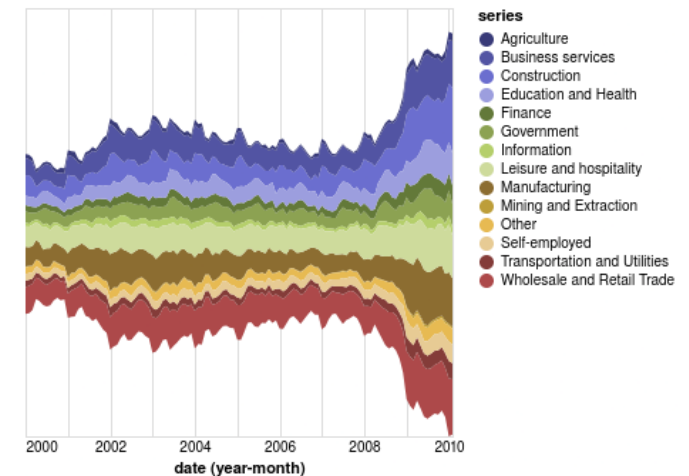
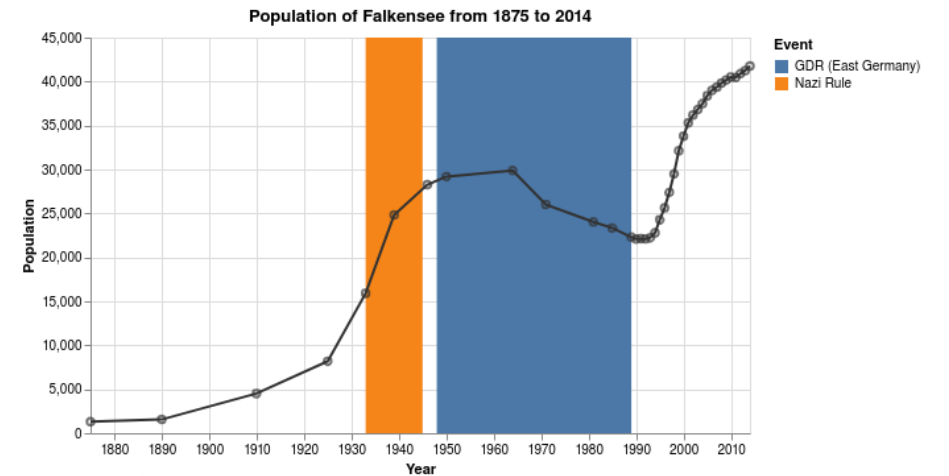
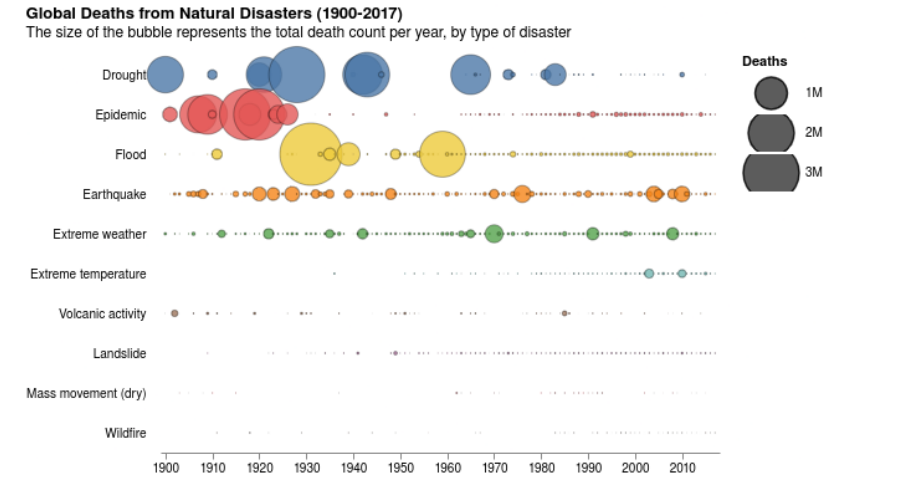
- Fonctionne bien avec Bokeh.
- API de haut niveau.

Désavantages

- Pas générique (du moins pour l'instant).

Altair

- Conçu pour les ordinateurs portables.
- Utilise le moteur Vega (qui produit du code JavaScript).

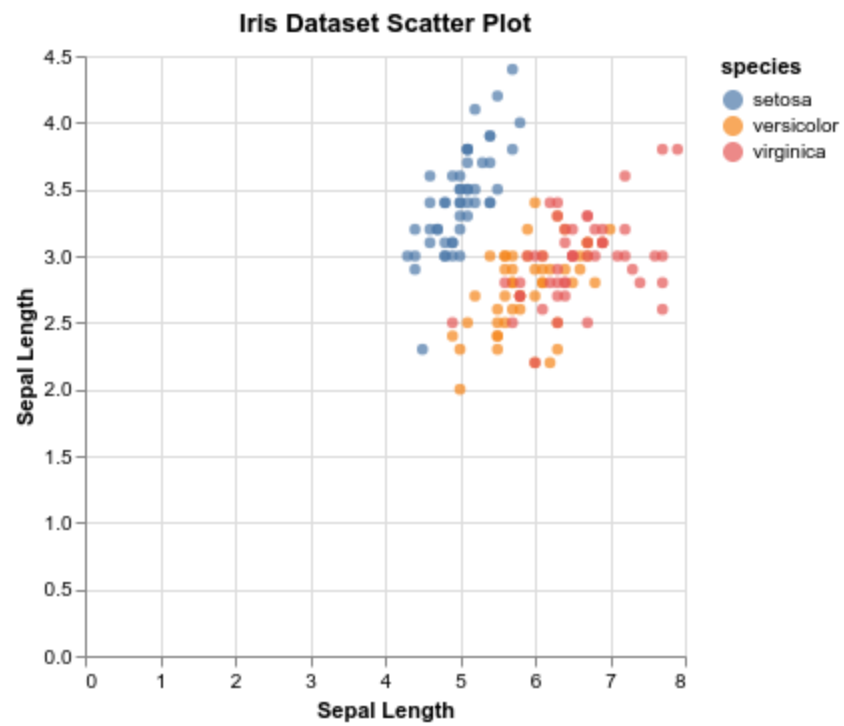


Altair - diagramme de dispersion des sépales de l'iris

```
import altair as alt

# Create the scatter plot
scatter_plot = alt.Chart(df_iris).mark_circle().encode(
    x=alt.X('sepal_length', title='Sepal Length'),
    y=alt.Y('sepal_width', title='Sepal Length'),
    color='species'
).properties(
    title='Iris Dataset Scatter Plot'
)

# Show the plot
scatter_plot
```



Altair

Avantages

- Utilisation simple des dataframes.
- Performances plus agréables qu'un Plotly dans un notebook VS Code.

Inconvénients

- Les valeurs par défaut ne sont pas très bonnes pour les graphiques statistiques.
- Pas de visualiseur pour les dernières versions (5.x).
- Personnalisation limitée des graphiques.

Conclusion



Mon choix

Plotly

Avantages

- Excellentes valeurs par défaut qui réduisent mes codes.
- Excellente documentation.
- Supporté par MLFlow.

Inconvénients

- Beaucoup de tracés dans les notebooks VSCode implique des ralentissements.
- Avoir des images en tant que ticks est difficile à mettre en œuvre (si vous connaissez un moyen simple, n'hésitez pas).

Bonus dashboard/webapp avec Streamlit

Configuration


Choose your preferred API:


echarts

Choose an example

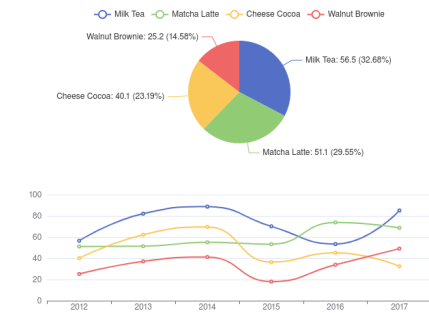
Dataset: Share Dataset

ECharts demos are extracted from <https://echarts.apache.org/examples/en/index.html> by copying/formatting the 'option' json object into st_echarts. Definitely check the echarts example page, convert the JSDN specs to Python Dicts and you should get a nice viz.

Made in  by [@andfania](#)

 Buy me a coffee

Streamlit ECharts Demo



Source Code

Credit: <https://echarts.apache.org/examples/en/editor.html?c=dataset-link>

Made with Streamlit



Merci d'avoir écouté !