
Phân tích dữ liệu

Data Analytics and Visualization with Python

TS. Đỗ Như Tài
dntai@sgu.edu.vn

Outline

- NumPy
 - Numerical Python N-dimensional array
- Pandas
 - Data Analytics
- Matplotlib
 - Basic Data Visualization
- Seaborn
 - Advanced Visualization



W3Schools Python Numpy

HTML CSS JAVASCRIPT SQL **PYTHON** JAVA PHP HOW TO W3.CSS C C++ C# BOOTSTRAP REACT

Python Modules

- NumPy Tutorial
- Pandas Tutorial
- SciPy Tutorial
- Django Tutorial

NumPy Tutorial

NumPy HOME

- NumPy Intro
- NumPy Getting Started
- NumPy Creating Arrays
- NumPy Array Indexing
- NumPy Array Slicing
- NumPy Data Types
- NumPy Copy vs View
- NumPy Array Shape
- NumPy Array Reshape
- NumPy Array Iterating
- NumPy Array Join
- NumPy Array Split
- NumPy Array Search
- NumPy Array Sort
- NumPy Array Filter

NumPy Random

- Random Intro
- Data Distribution

< Home

Next >

NumPy is a Python library.

NumPy is used for working with arrays.

NumPy is short for "Numerical Python".

[+ :]

Learning by Reading

We have created 43 tutorial pages for you to learn more about NumPy.

Starting with a basic introduction and ends up with creating and plotting random data sets, and working with NumPy functions:

Basic

Random

ufunc

<https://www.w3schools.com/python/numpy/default.asp>



Python Modules

NumPy Tutorial
Pandas Tutorial
SciPy Tutorial
Django Tutorial

Pandas Tutorial

Pandas HOME
Pandas Intro
Pandas Getting Started
Pandas Series
Pandas DataFrames
Pandas Read CSV
Pandas Read JSON
Pandas Analyzing Data

Cleaning Data

Cleaning Data
Cleaning Empty Cells
Cleaning Wrong Format
Cleaning Wrong Data
Removing Duplicates

Correlations

Pandas Correlations

Plotting

Pandas Plotting

W3Schools Python Pandas

Learning by Reading

Pandas Tutorial

We have created 14 tutorial pages for you to learn more about Pandas.

Starting with a basic introduction and ends up with cleaning and plotting data:

Basic

Introduction

Getting Started

Pandas Series

DataFrames

Read CSV

Read JSON

Analyze Data

Cleaning Data

Clean Data

Clean Empty Cells

Clean Wrong Format

Clean Wrong Data

Remove Duplicates

Advanced

Correlations

Plotting

<https://www.w3schools.com/python/pandas/default.asp>



W3Schools Python

HTML CSS JAVASCRIPT SQL **PYTHON** JAVA PHP HOW TO W3.CSS C C++ C# BOOTSTRAP REACT

Python Modules

NumPy Tutorial
Pandas Tutorial
SciPy Tutorial
Django Tutorial

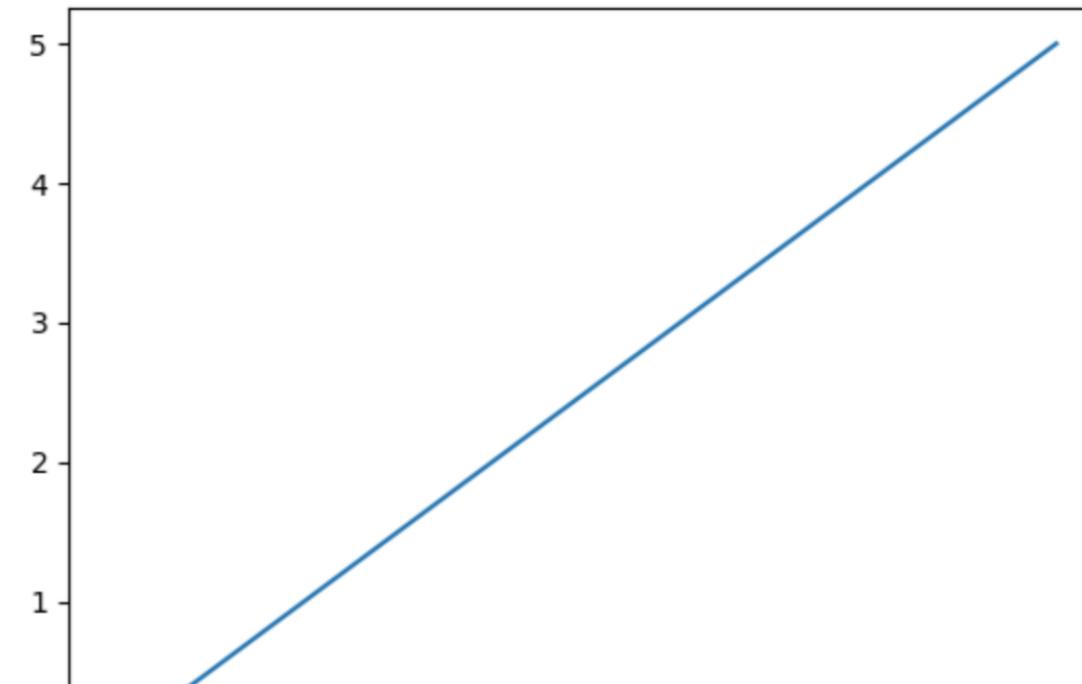
Python Matplotlib

Matplotlib Intro
Matplotlib Get Started
Matplotlib Pyplot
Matplotlib Plotting
Matplotlib Markers
Matplotlib Line
Matplotlib Labels
Matplotlib Grid
Matplotlib Subplot
Matplotlib Scatter
Matplotlib Bars
Matplotlib Histograms
Matplotlib Pie Charts

Matplotlib Tutorial

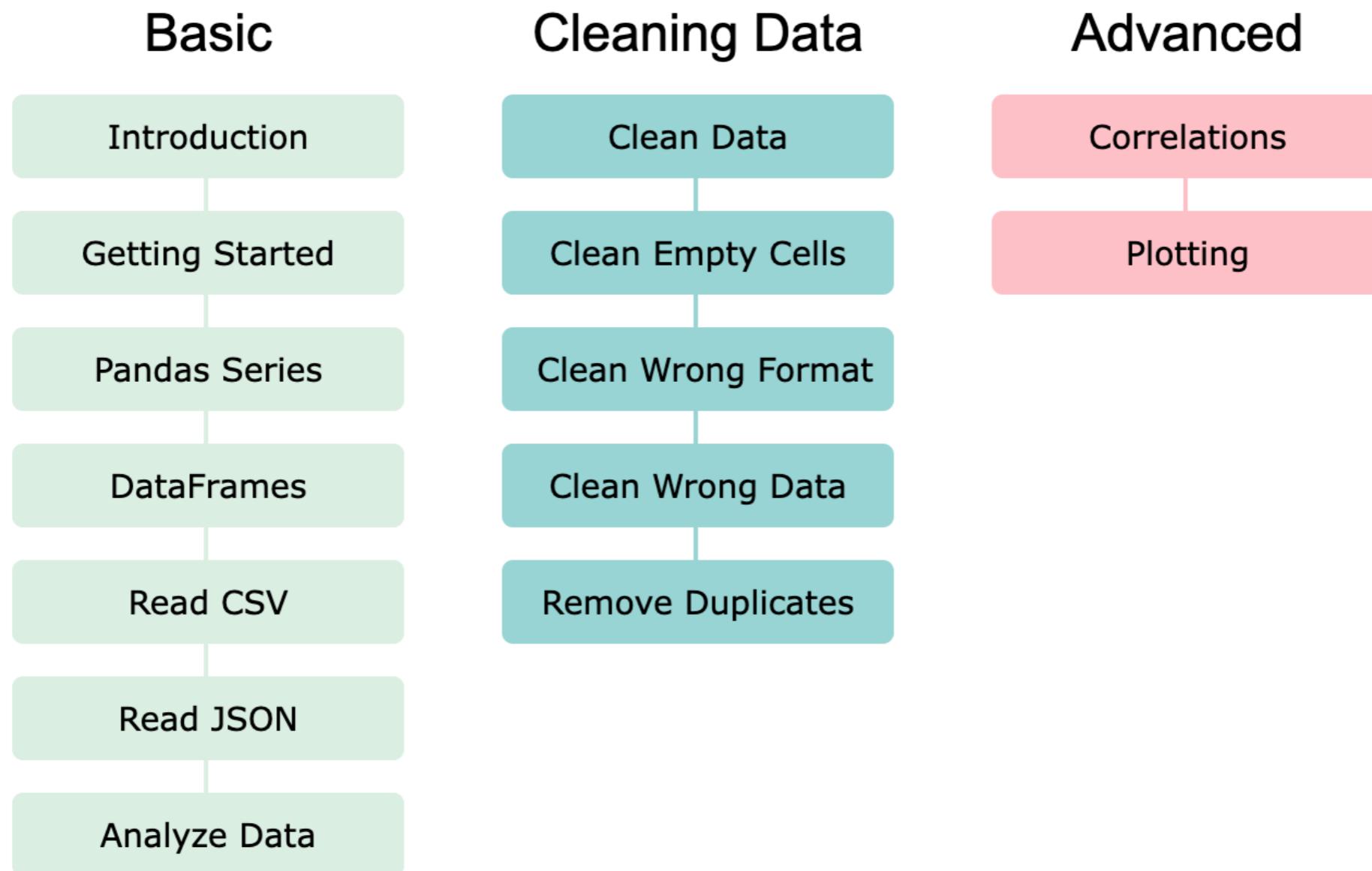
◀ Previous

Next ▶



<https://www.w3schools.com/python/>

Pandas: Data Analytics and Visualization



<https://www.w3schools.com/python/pandas/default.asp>

Wes McKinney (2022), "Python for Data Analysis: Data Wrangling with pandas, NumPy, and Jupyter", 3rd Edition, O'Reilly Media.

[wesm/pydata-book](https://github.com/wesm/pydata-book) Public

Code Issues Pull requests Actions Projects Wiki Security Insights

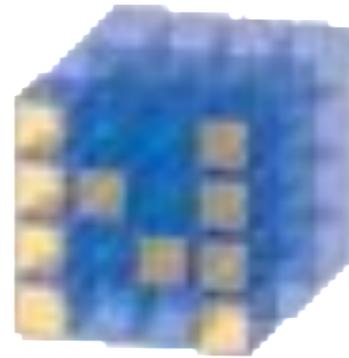
3rd-edition 3 branches 0 tags Go to file Code About

File	Description	Time
wesm	Upload cleaner notebook files without internal build toolchai...	3 days ago
datasets	Add fec.parquet	10 months ago
examples	Simplifying datasets	10 months ago
.gitignore	Add gitignore	8 years ago
COPYING	Update COPYING	4 months ago
README.md	Update notebooks in advance of publication	7 months ago
appa.ipynb	Upload cleaner notebook files without internal build toolchai...	3 days ago
appb.ipynb	Upload cleaner notebook files without internal build toolchai...	3 days ago
ch02.ipynb	Upload cleaner notebook files without internal build toolchai...	3 days ago
ch03.ipynb	Upload cleaner notebook files without internal build toolchai...	3 days ago
ch04.ipynb	Upload cleaner notebook files without internal build toolchai...	3 days ago
ch05.ipynb	Upload cleaner notebook files without internal build toolchai...	3 days ago
ch06.ipynb	Upload cleaner notebook files without internal build toolchai...	3 days ago

O'REILLY
Python for Data Analysis
Data Wrangling with pandas, NumPy & Jupyter
Third Edition
Wes McKinney

<https://github.com/wesm/pydata-book>

Numpy



NumPy
Base
N-dimensional array
package

NumPy
is the
fundamental package
for
scientific computing
with Python.

Source: <http://www.numpy.org/>

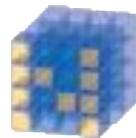


NumPy

NumPy

- NumPy provides a **multidimensional array** object to store homogenous or heterogeneous data; it also provides **optimized functions/methods** to operate on this array object.

Source: Yves Hilpisch (2014), Python for Finance: Analyze Big Financial Data, O'Reilly



NumPy

NumPy ndarray

One-dimensional Array (1-D Array)

0	1	n-1
1	2	3

Two-dimensional Array (2-D Array)

0	1	n-1			
0	1	2	3	4	5
1	6	7	8	9	10
2	11	12	13	14	15
m-1	16	17	18	19	20



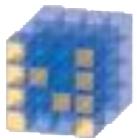
NumPy

NumPy

```
v = list(range(1, 6))
v
2 * v

import numpy as np
v = np.arange(1, 6)
v
2 * v
```

Source: Yves Hilpisch (2014), Python for Finance: Analyze Big Financial Data, O'Reilly



NumPy
Base
N-dimensional
array package

```
1 v = list(range(1, 6))  
2 v
```

```
[1, 2, 3, 4, 5]
```

```
1 2 * v
```

```
[1, 2, 3, 4, 5, 1, 2, 3, 4, 5]
```

```
1 import numpy as np  
2 v = np.arange(1, 6)  
3 v
```

```
array([1, 2, 3, 4, 5])
```

```
1 2 * v
```

```
array([ 2,  4,  6,  8, 10])
```

Python Data Structures

```
fruits = ["apple", "banana", "cherry"] #lists []
colors = ("red", "green", "blue") #tuples ()
animals = {'cat', 'dog'} #sets {}
person = {"name": "Tom", "age": 20} #dictionaries { }
```



Lists []

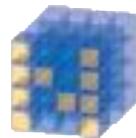
```
x = [60, 70, 80, 90]  
print(len(x))  
print(x[0])  
print(x[1])  
print(x[-1])
```

4

60

70

90



NumPy

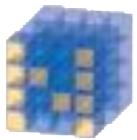
NumPy Create Array

```
import numpy as np
a = np.array([1, 2, 3])
b = np.array([4, 5, 6])
c = a * b
c
```

```
import numpy as np
a = np.array([1, 2, 3])
b = np.array([4, 5, 6])
c = a * b
c

array([ 4, 10, 18])
```

Source: Yves Hilpisch (2014), Python for Finance: Analyze Big Financial Data, O'Reilly



NumPy

NumPy

```
1 import numpy as np
2
3 a = np.zeros((2,2)) # Create an array of all zeros
4 print(a)           # Prints "[[ 0.  0.]
5                  #          [ 0.  0.]]"
6
7 b = np.ones((1,2)) # Create an array of all ones
8 print(b)           # Prints "[[ 1.  1.]]"
9
10 c = np.full((2,2), 7) # Create a constant array
11 print(c)            # Prints "[[ 7.  7.]
12                  #          [ 7.  7.]]"
13
14 d = np.eye(2)       # Create a 2x2 identity matrix
15 print(d)            # Prints "[[ 1.  0.]
16                  #          [ 0.  1.]]"
17
18 e = np.random.random((2,2)) # Create an array filled with random values
19 print(e)             # Might print "[[ 0.91940167  0.08143941]
20                  #          [ 0.68744134  0.87236687]]"
```

```
[[0.  0.]
 [0.  0.]]
[[1.  1.]]
[[7 7]
 [7 7]]
[[1. 0.]
 [0. 1.]]
[[0.66258211 0.65552598]
 [0.00429934 0.21695824]]
```

Source: <http://cs231n.github.io/python-numpy-tutorial/>

```
import numpy as np  
a = np.arange(15).reshape(3, 5)
```

a.shape
a.ndim
a.dtype.name

```
import numpy as np  
a = np.arange(15).reshape(3, 5)  
  
a  
  
array([[ 0,  1,  2,  3,  4],  
       [ 5,  6,  7,  8,  9],  
       [10, 11, 12, 13, 14]])
```

```
print(a.shape)
```

(3, 5)

```
a.ndim
```

2

```
a.dtype.name
```

'int64'

Source: <https://docs.scipy.org/doc/numpy-dev/user/quickstart.html>

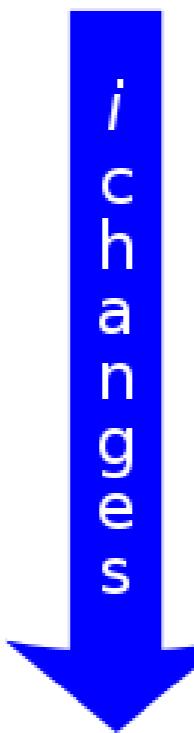
Matrix

m-by-*n* matrix

$a_{i,j}$

n columns j changes

m
rows

i
changes


$a_{1,1}$	$a_{1,2}$	$a_{1,3}$	\dots
$a_{2,1}$	$a_{2,2}$	$a_{2,3}$	\dots
$a_{3,1}$	$a_{3,2}$	$a_{3,3}$	\dots
\vdots	\vdots	\vdots	\ddots

Source: [https://simple.wikipedia.org/wiki/Matrix_\(mathematics\)](https://simple.wikipedia.org/wiki/Matrix_(mathematics))

NumPy ndarray: Multidimensional Array Object

NumPy ndarray

One-dimensional Array (1-D Array)

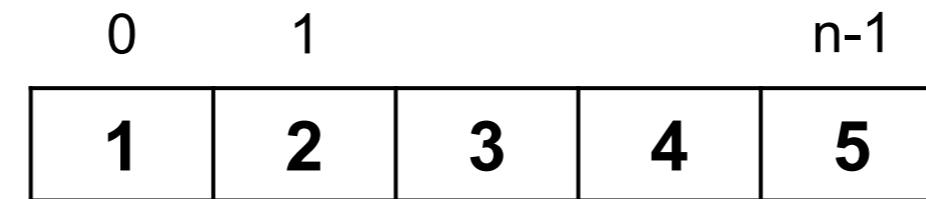
0	1	n-1
1	2	3

Two-dimensional Array (2-D Array)

0	1	n-1			
0	1	2	3	4	5
1	6	7	8	9	10
2	11	12	13	14	15
m-1	16	17	18	19	20

```
import numpy as np  
a = np.array([1,2,3,4,5])
```

One-dimensional Array (1-D Array)



```
a = np.array([1,2,3,4,5])  
a
```

```
array([1, 2, 3, 4, 5])
```

```
a = np.array([ [1,2,3,4,5],[6,7,8,9,10],[11,12,13,14,15],[16,17,18,19,20] ])
```

Two-dimensional Array (2-D Array)

	0	1		n-1	
0	1	2	3	4	5
1	6	7	8	9	10
	11	12	13	14	15
m-1	16	17	18	19	20

```
a = np.array([[1,2,3,4,5],[6,7,8,9,10],[11,12,13,14,15],[16,17,18,19,20]])  
a  
  
array([[ 1,  2,  3,  4,  5],  
       [ 6,  7,  8,  9, 10],  
       [11, 12, 13, 14, 15],  
       [16, 17, 18, 19, 20]])
```

```
import numpy as np  
a = np.array([[0, 1, 2, 3],  
[10, 11, 12, 13],  
[20, 21, 22, 23]])  
a
```

0	1	2	3
10	11	12	13
20	21	22	23

```
a = np.array([[0, 1, 2, 3], [10, 11, 12, 13], [20, 21, 22, 23]])
```

```
a = np.array([[0, 1, 2, 3], [10, 11, 12, 13], [20, 21, 22, 23]])  
a
```

```
array([[ 0,  1,  2,  3],  
       [10, 11, 12, 13],  
       [20, 21, 22, 23]])
```

```
print(a.ndim)
```

```
2
```

```
print(a.shape)
```

```
(3, 4)
```

0	1	2	3
10	11	12	13
20	21	22	23

NumPy Basics: Arrays and Vectorized Computation

Source: <https://www.safaribooksonline.com/library/view/python-for-data/9781449323592/ch04.html>

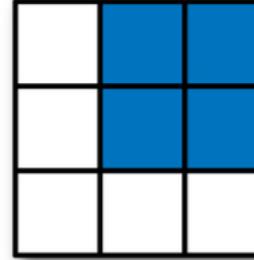
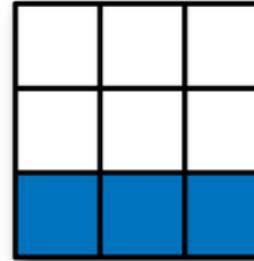
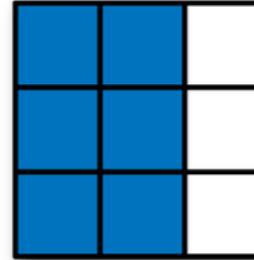
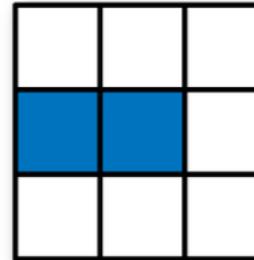
NumPy Array

axis 1

		0	1	2
		0, 0	0, 1	0, 2
		1, 0	1, 1	1, 2
axis 0	0	0, 0	0, 1	0, 2
	1	1, 0	1, 1	1, 2
	2	2, 0	2, 1	2, 2

Source: <https://www.safaribooksonline.com/library/view/python-for-data/9781449323592/ch04.html>

Numpy Array

Expression	Shape
	<code>arr[:2, 1:]</code> (2, 2)
	<code>arr[2]</code> (3,) <code>arr[2, :]</code> (3,) <code>arr[2:, :]</code> (1, 3)
	<code>arr[:, :2]</code> (3, 2)
	<code>arr[1, :2]</code> (2,) <code>arr[1:2, :2]</code> (1, 2)

Source: <https://www.safaribooksonline.com/library/view/python-for-data/9781449323592/ch04.html>



Tensor

- **3**
 - a rank 0 tensor; this is a **scalar** with shape []
- **[1., 2., 3.]**
 - a rank 1 tensor; this is a **vector** with shape [3]
- **[[1., 2., 3.], [4., 5., 6.]]**
 - a rank 2 tensor; a **matrix** with shape [2, 3]
- **[[[1., 2., 3.]], [[7., 8., 9.]]]**
 - a rank 3 **tensor** with shape [2, 1, 3]

<https://www.tensorflow.org/>

Scalar

80

Vector

[50 60 70]

Matrix

$$\begin{bmatrix} 50 & 60 & 70 \\ 55 & 65 & 75 \end{bmatrix}$$

Tensor

$$\begin{bmatrix} [50 & 60 & 70] & [70 & 80 & 90] \\ [55 & 65 & 75] & [75 & 85 & 95] \end{bmatrix}$$



pandas

Python Data Analysis Library

providing high-performance, easy-to-use
data structures and data analysis tools
for the Python programming language.

Source: <http://pandas.pydata.org/>



pandas: powerful Python data analysis toolkit

- Tabular data with heterogeneously-typed columns, as in an SQL table or Excel spreadsheet
- Ordered and unordered (not necessarily fixed-frequency) time series data.
- Arbitrary matrix data (homogeneously typed or heterogeneous) with row and column labels
- Any other form of observational / statistical data sets. The data actually need not be labeled at all to be placed into a pandas data structure

Source: <http://pandas.pydata.org/pandas-docs/stable/>



Series DataFrame

- Primary data structures of pandas
 - Series (1-dimensional)
 - DataFrame (2-dimensional)
- Handle the vast majority of typical use cases in **finance**, statistics, social science, and many areas of engineering.

Source: <http://pandas.pydata.org/pandas-docs/stable/>

pandas DataFrame



- **DataFrame** provides everything that R's `data.frame` provides and much more.
- pandas is built on top of **NumPy** and is intended to integrate well within a scientific computing environment with many other 3rd party libraries.



pandas

Comparison with SAS

pandas	SAS
DataFrame	data set
column	variable
row	observation
groupby	BY-group
NaN	.

Source: http://pandas.pydata.org/pandas-docs/stable/comparison_with_sas.html



Python Pandas Cheat Sheet

Data Wrangling
with pandas
Cheat Sheet
<http://pandas.pydata.org>

Syntax – Creating DataFrames

```
df = pd.DataFrame(
    {"a" : [4, 5, 6],
     "b" : [7, 8, 9],
     "c" : [10, 11, 12]},
    index = [1, 2, 3])
Specify values for each column.

df = pd.DataFrame(
    [[4, 7, 10],
     [5, 8, 11],
     [6, 9, 12]],
    index=[1, 2, 3],
    columns=['a', 'b', 'c'])
Specify values for each row.
```

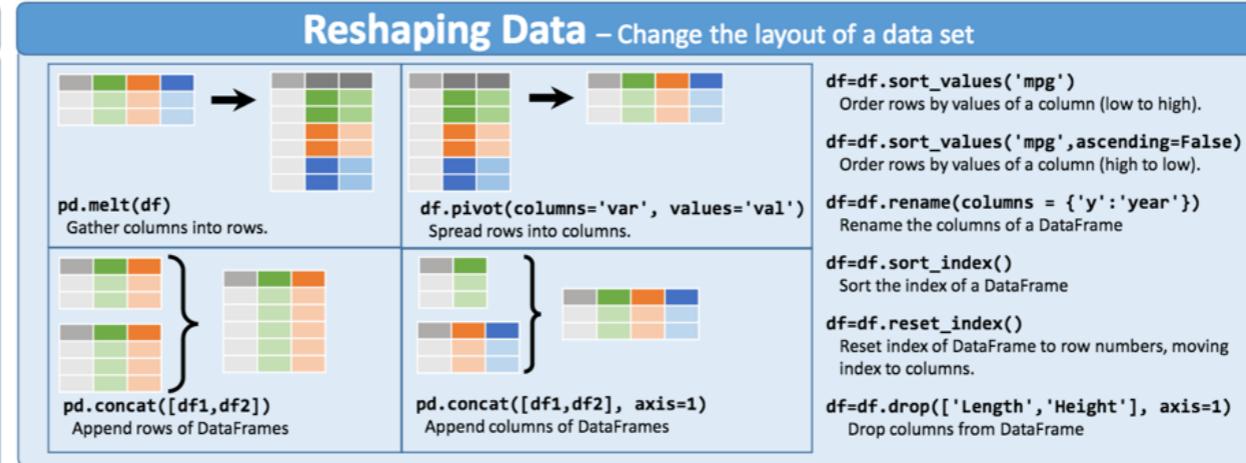
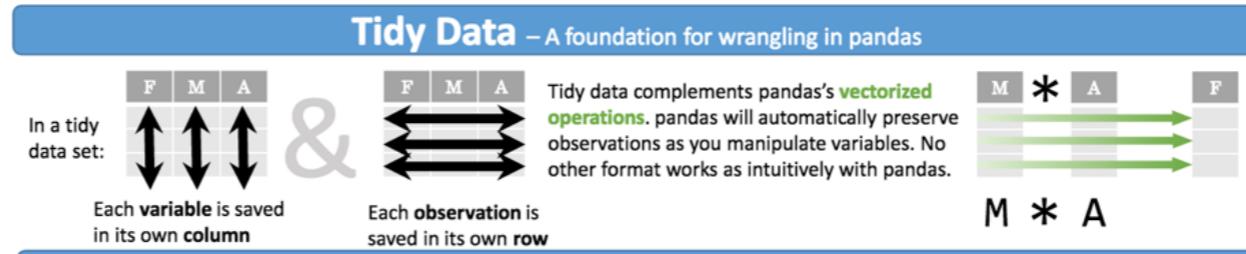
	a	b	c
1	4	7	10
2	5	8	11
3	6	9	12

```
df = pd.DataFrame(
    {"a" : [4, 5, 6],
     "b" : [7, 8, 9],
     "c" : [10, 11, 12]},
    index = pd.MultiIndex.from_tuples(
        [('d',1),('d',2),('e',2)],
        names=['n','v']))
Create DataFrame with a MultiIndex
```

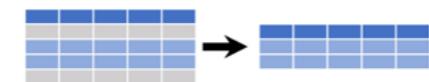
Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

```
df = (pd.melt(df)
      .rename(columns={
          'variable' : 'var',
          'value' : 'val'})
      .query('val >= 200')
     )
```



Subset Observations (Rows)



```
df[df.Length > 7]
Extract rows that meet logical criteria.

df.drop_duplicates()
Remove duplicate rows (only considers columns).

df.head(n)
Select first n rows.

df.tail(n)
Select last n rows.
```

`df.sample(frac=0.5)`
Randomly select fraction of rows.

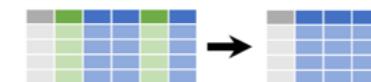
`df.sample(n=10)`
Randomly select n rows.

`df.iloc[10:20]`
Select rows by position.

`df.nlargest(n, 'value')`
Select and order top n entries.

`df.nsmallest(n, 'value')`
Select and order bottom n entries.

Subset Variables (Columns)



`df[['width', 'length', 'species']]`
Select multiple columns with specific names.

`df['width'] or df.width`
Select single column with specific name.

`df.filter(regex='regex')`
Select columns whose name matches regular expression regex.

regex (Regular Expressions) Examples

'.'	Matches strings containing a period '.'.
'Length\$'	Matches strings ending with word 'Length'
'^Sepal'	Matches strings beginning with the word 'Sepal'
'^x[1-5]\$'	Matches strings beginning with 'x' and ending with 1,2,3,4,5
'^(?!Species\$).*''	Matches strings except the string 'Species'

`df.loc[:, 'x2': 'x4']`
Select all columns between x2 and x4 (inclusive).

`df.iloc[:, [1,2,5]]`
Select columns in positions 1, 2 and 5 (first column is 0).

`df.loc[df['a'] > 10, ['a', 'c']]`
Select rows meeting logical condition, and only the specific columns.

Source: https://github.com/pandas-dev/pandas/blob/master/doc/cheatsheet/Pandas_Cheat_Sheet.pdf

<http://pandas.pydata.org/> This cheat sheet inspired by Rstudio Data Wrangling Cheatsheet (<https://www.rstudio.com/wp-content/uploads/2015/02/data-wrangling-cheatsheet.pdf>) Written by Irv Lustig, Princeton Consultants



Creating pd.DataFrame

	a	b	c
1	4	7	10
2	5	8	11
3	6	9	12

```
In [1]: import numpy as np
import pandas as pd
df = pd.DataFrame({"a": [4, 5, 6],
                   "b": [7, 8, 9],
                   "c": [10, 11, 12]},
                  index = [1, 2, 3])

df
```

Out[1]:

	a	b	c
1	4	7	10
2	5	8	11
3	6	9	12

```
import pandas as pd
df = pd.DataFrame({"a": [4, 5, 6],
                   "b": [7, 8, 9],
                   "c": [10, 11, 12]},
                  index = [1, 2, 3])
```



Pandas DataFrame

```
type(df)
```

```
type(df)
```

pandas.core.frame.DataFrame



```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
print('pandas imported')
```

```
s = pd.Series([1,3,5,np.nan,6,8])
s
```

```
dates = pd.date_range('20181001',
periods=6)
dates
```

Source: <http://pandas.pydata.org/pandas-docs/stable/10min.html>



```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 print('pandas imported')
```

```
pandas imported
```

```
1 s = pd.Series([1,3,5,np.nan,6,8])
2 s
```

```
0    1.0
1    3.0
2    5.0
3    NaN
4    6.0
5    8.0
dtype: float64
```

```
1 dates = pd.date_range('20181001', periods=6)
2 dates
```

```
DatetimeIndex(['2018-10-01', '2018-10-02', '2018-10-03', '2018-10-04',
                 '2018-10-05', '2018-10-06'],
                dtype='datetime64[ns]', freq='D')
```



```
df = pd.DataFrame(np.random.randn(6,4),  
index=dates, columns=list('ABCD'))  
df
```

```
1 df = pd.DataFrame(np.random.randn(6,4), index=dates, columns=list('ABCD'))  
2 df
```

	A	B	C	D
2018-10-01	-0.336188	0.584621	-1.061433	-0.036278
2018-10-02	0.903683	-0.839723	-0.270219	-1.099606
2018-10-03	0.920208	-0.240353	-0.818598	-1.105489
2018-10-04	0.221045	-0.314589	0.042071	-1.447280
2018-10-05	0.946862	-1.570305	-1.009180	-0.375659
2018-10-06	-0.225148	0.510691	2.002372	-0.335005



```
df = pd.DataFrame(np.random.randn(3,5),  
index=['student1','student2','student3'],  
columns=list('ABCDE'))  
df
```

```
1 df = pd.DataFrame(np.random.randn(3,5), index=['student1','student2','student3'], columns=list('ABCDE'))  
2 df
```

	A	B	C	D	E
student1	-0.346884	-1.232934	-0.302072	-1.345084	-0.723880
student2	1.090955	-0.010483	1.280072	-0.253958	-0.030604
student3	0.325660	0.808956	-0.395820	-1.498926	1.603471



```
df2 = pd.DataFrame({ 'A' : 1.,
'B' : pd.Timestamp('20181001'),
'C' : pd.Series(2.5,index=list(range(4)),dtype='float32'),
'D' : np.array([3] * 4,dtype='int32'),
'E' : pd.Categorical(["test","train","test","train"]),
'F' : 'foo' })
df2
```

```
1 df2 = pd.DataFrame({ 'A' : 1.,
2 'B' : pd.Timestamp('20181001'),
3 'C' : pd.Series(2.5,index=list(range(4)),dtype='float32'),
4 'D' : np.array([3] * 4,dtype='int32'),
5 'E' : pd.Categorical(["test","train","test","train"]),
6 'F' : 'foo' })
7 df2
```

	A	B	C	D	E	F
0	1.0	2018-10-01	2.5	3	test	foo
1	1.0	2018-10-01	2.5	3	train	foo
2	1.0	2018-10-01	2.5	3	test	foo
3	1.0	2018-10-01	2.5	3	train	foo



df2.dtypes

df2.dtypes

```
A          float64
B    datetime64[ns]
C          float32
D          int32
E      category
F          object
dtype: object
```

Python Finance Application with Pandas

```
import pandas as pd

# Create a DataFrame to store transactions
columns = ['Date', 'Description', 'Amount']
ledger = pd.DataFrame(columns=columns)

# Function to add a transaction
def add_transaction(date, description, amount):
    global ledger
    new_transaction = pd.DataFrame([[date, description, amount]], columns=columns)
    ledger = pd.concat([ledger, new_transaction], ignore_index=True)

# Function to view the ledger
def view_ledger():
    print(ledger)

# Function to get the current balance
def get_balance():
    return ledger['Amount'].sum()

# Adding sample transactions
add_transaction('2023-11-01', 'Income', 1000)
add_transaction('2023-11-02', 'Groceries', -200)
add_transaction('2023-11-03', 'Utilities', -100)

# Viewing the ledger
view_ledger()

# Checking the current balance
print("Current Balance:", get_balance())
```

	Date	Description	Amount
0	2023-11-01	Income	1000
1	2023-11-02	Groceries	-200
2	2023-11-03	Utilities	-100
		Current Balance:	700

Python Data Analysis and Visualization



Python Pandas



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

Python matplotlib

matplotlib



Source: <https://matplotlib.org/>

Python seaborn



seaborn

Source: <https://seaborn.pydata.org/>

Python plotly



plotly

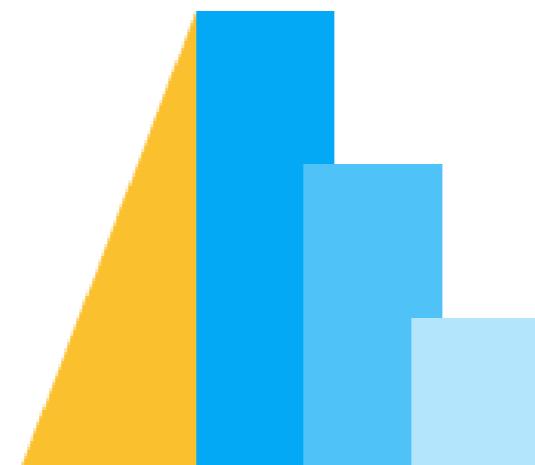
Source: <https://plotly.com/python/>

Python bokeh



Source: <https://bokeh.org/>

Python Altair



Altair

Source: <https://altair-viz.github.io/>

Python matplotlib

The screenshot shows the official Matplotlib website. At the top, there's a large logo for "matplotlib" with "Version 3.3.4" below it. To the right of the logo is an orange "Fork me on GitHub" button. Below the logo is a navigation bar with links for "Installation", "Documentation", "Examples", "Tutorials", and "Contributing". A search bar is also present. The main content area features the title "Matplotlib: Visualization with Python" and a brief description: "Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python." Below the description are four small images illustrating different types of plots: a line plot with oscillations, a histogram-like plot with a peak, a heatmap, and a 3D surface plot. A sidebar on the right contains links for the "Latest stable release" (version 3.3.4), "Last release for Python 2" (version 2.2.5), and the "Development version". There's also a "Support Matplotlib" button at the bottom of the sidebar.

matplotlib Version 3.3.4

Fork me on GitHub

Installation Documentation Examples Tutorials Contributing Search

home | contents » Matplotlib: Python plotting modules | index

Matplotlib: Visualization with Python

Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python.



Matplotlib makes easy things easy and hard things possible.

Create

- Develop publication quality plots with just a few lines of code
- Use interactive figures that can zoom, pan, update...

Customize

- Take full control of line styles, font properties, axes properties...
- Export and embed to a number of file formats and interactive environments

Extend

- Explore tailored functionality provided by third party packages
- Learn more about Matplotlib through the many external learning resources

Latest stable release
3.3.4: [docs](#) | [changelog](#)

Last release for Python 2
2.2.5: [docs](#) | [changelog](#)

Development version
[docs](#)

Support Matplotlib

<https://matplotlib.org/>



seaborn Python Seaborn



0.11.1

Gallery

Tutorial

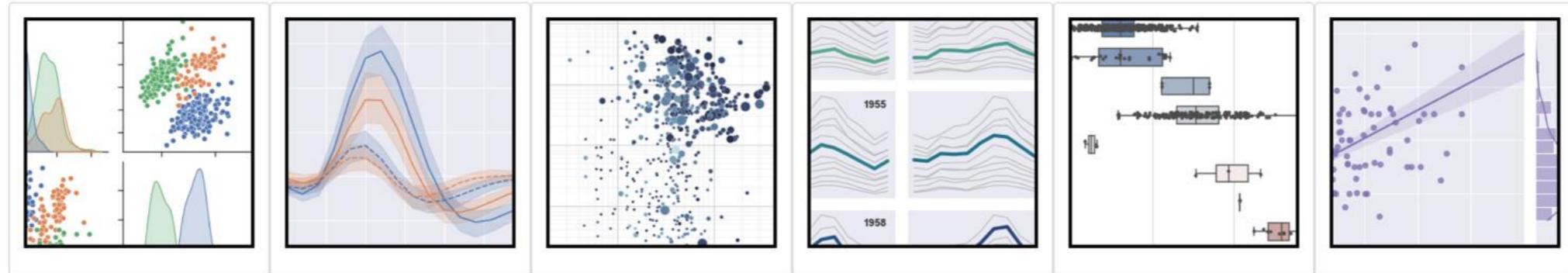
API

Site ▾

Page ▾

Search

seaborn: statistical data visualization



Seaborn is a Python data visualization library based on [matplotlib](#). It provides a high-level interface for drawing attractive and informative statistical graphics.

For a brief introduction to the ideas behind the library, you can read the [introductory notes](#). Visit the [installation page](#) to see how you can download the package and get started with it. You can browse the [example gallery](#) to see what you can do with seaborn, and then check out the [tutorial](#) and [API reference](#) to find out how.

To see the code or report a bug, please visit the [GitHub repository](#). General support questions are most at home on [stackoverflow](#) or [discourse](#), which have dedicated channels for seaborn.

Contents

- [Introduction](#)
- [Release notes](#)
- [Installing](#)
- [Example gallery](#)
- [Tutorial](#)
- [API reference](#)

Features

- Relational: [API](#) | [Tutorial](#)
- Distribution: [API](#) | [Tutorial](#)
- Categorical: [API](#) | [Tutorial](#)
- Regression: [API](#) | [Tutorial](#)
- Multiples: [API](#) | [Tutorial](#)
- Style: [API](#) | [Tutorial](#)
- Color: [API](#) | [Tutorial](#)



Python Plotly Graphing Library



Star 9,085

[DO MORE WITH DASH](#)

Search...

Quick Start

- [Getting Started](#)
- [Is Plotly Free?](#)
- [Figure Reference](#)
- [API Reference](#)
- [Dash](#)
- [GitHub](#)
- [community.plotly.com](#)

Examples

- [Fundamentals](#)
- [Basic Charts](#)
- [Statistical Charts](#)
- [Artificial Intelligence and Machine Learning](#)
- [Scientific Charts](#)
- [Financial Charts](#)
- [Maps](#)
- [3D Charts](#)

Plotly Python Open Source Graphing Library

Plotly's Python graphing library makes interactive, publication-quality graphs. Examples of how to make line plots, scatter plots, area charts, bar charts, error bars, box plots, histograms, heatmaps, subplots, multiple-axes, polar charts, and bubble charts.

Plotly.py is [free and open source](#) and you can [view the source](#), [report issues](#) or [contribute on GitHub](#).

Our recommended IDE for Plotly's Python graphing library is Dash Enterprise's [Data Science Workspaces](#), which has both Jupyter notebook and Python code file support.

[Find out if your company is using Dash Enterprise.](#)

[Install Dash Enterprise on Azure](#) | [Install Dash Enterprise on AWS](#)

Fundamentals

The Figure Data Structure

Creating and Updating Figures

Displaying Figures

Plotly Express

Analytical Apps with Dash

[More Fundamentals »](#)

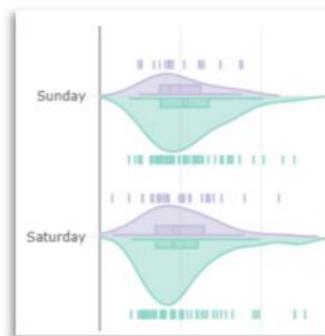
<https://plotly.com/python/>



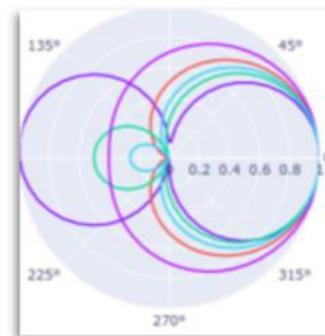
Python Plotly Graphing Library

Fundamentals

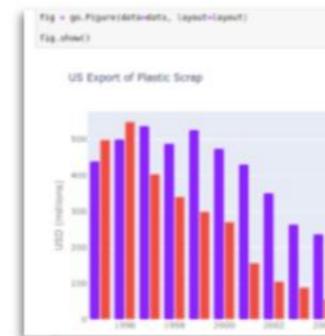
[More Fundamentals »](#)



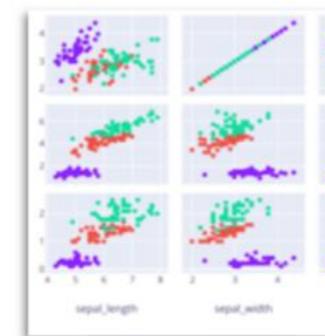
The Figure Data Structure



Creating and Updating Figures



Displaying Figures



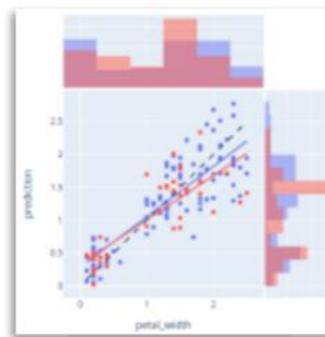
Plotly Express



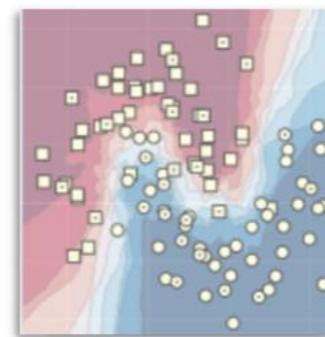
Analytical Apps with Dash

Artificial Intelligence and Machine Learning

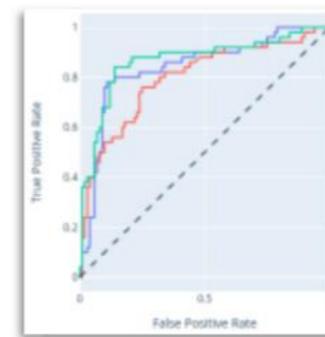
[More AI and ML »](#)



ML Regression



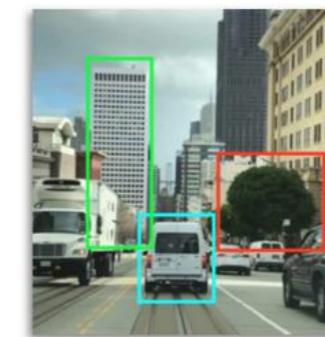
kNN Classification



ROC and PR Curves



PCA Visualization



AI/ML Apps with Dash

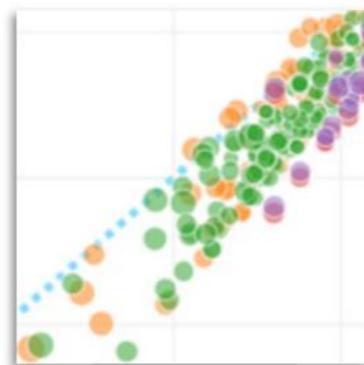
<https://plotly.com/python/>



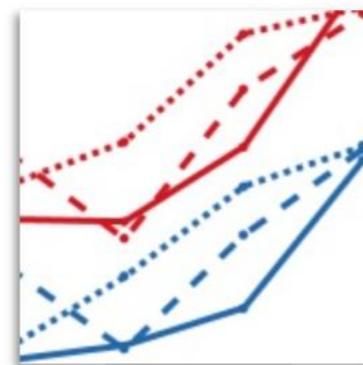
Python Plotly Graphing Library

Basic Charts

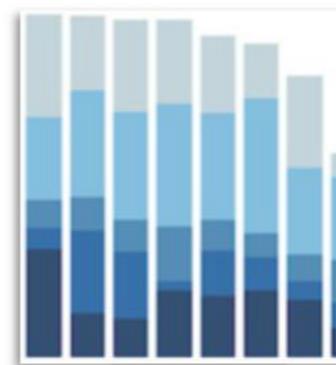
[More Basic Charts »](#)



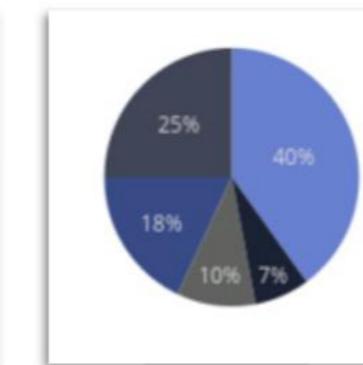
Scatter Plots



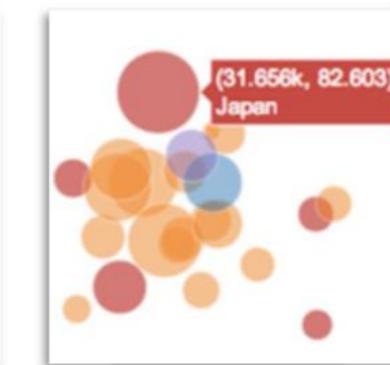
Line Charts



Bar Charts



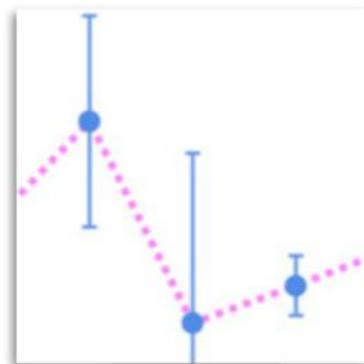
Pie Charts



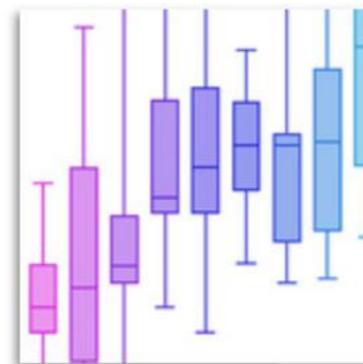
Bubble Charts

Statistical Charts

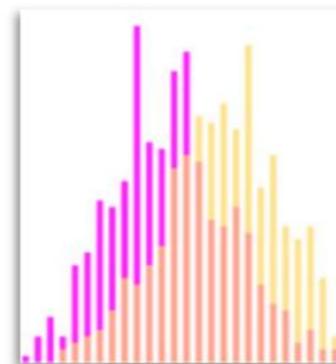
[More Statistical Charts »](#)



Error Bars



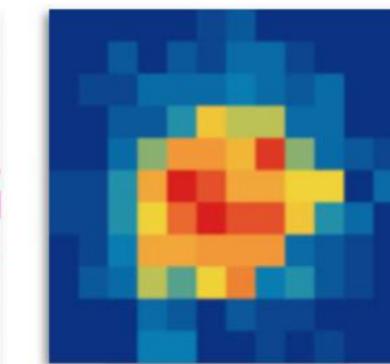
Box Plots



Histograms



Distplots



2D Histograms

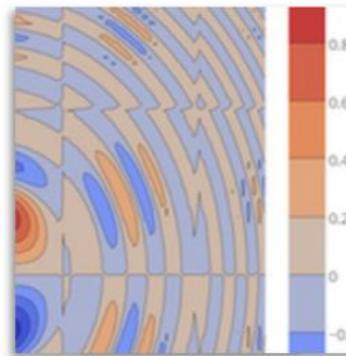
<https://plotly.com/python/>



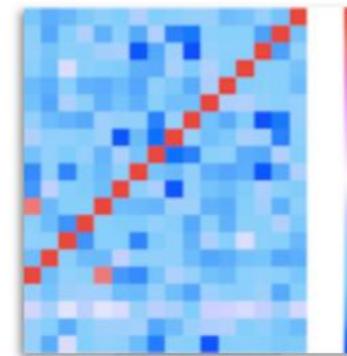
Python Plotly Graphing Library

Scientific Charts

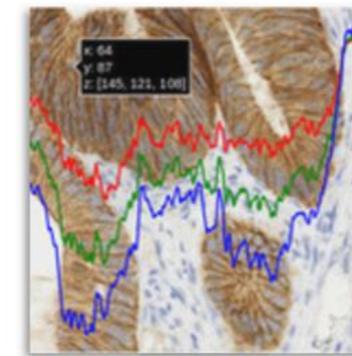
[More Scientific Charts »](#)



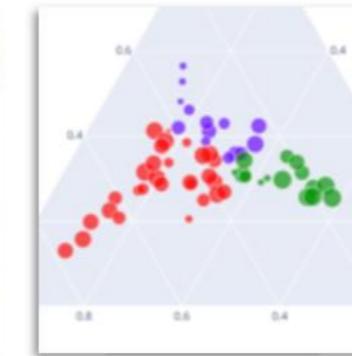
Contour Plots



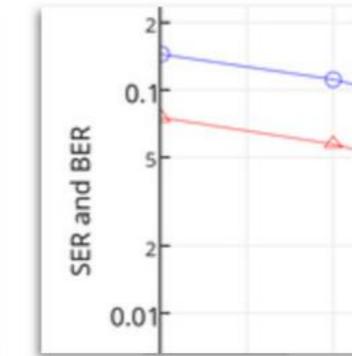
Heatmaps



Imshow



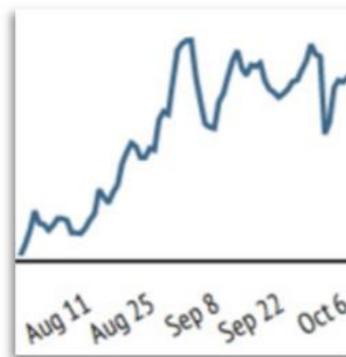
Ternary Plots



Log Plots

Financial Charts

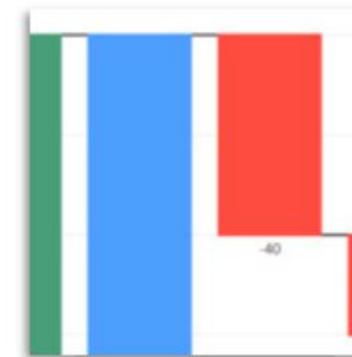
[More Financial Charts »](#)



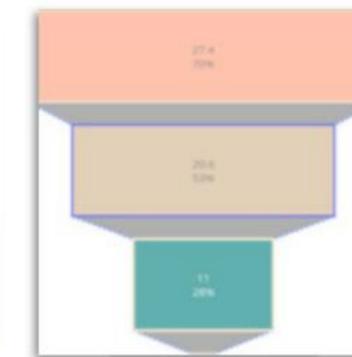
Time Series and Date
Axes



Candlestick Charts



Waterfall Charts



Funnel Chart



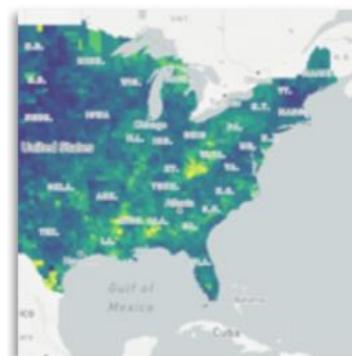
OHLC Charts

<https://plotly.com/python/>

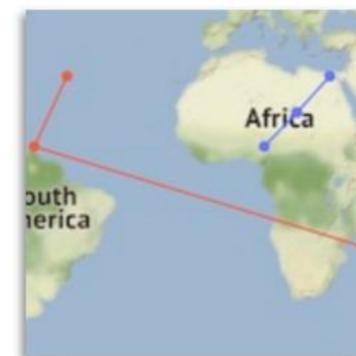


Python Plotly Graphing Library

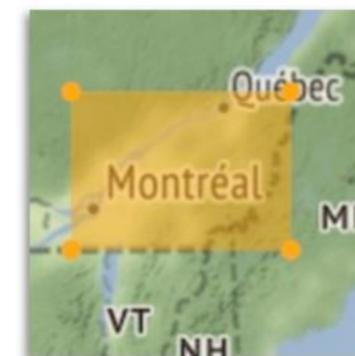
Maps

[More Maps ×](#)

Mapbox Choropleth
Maps



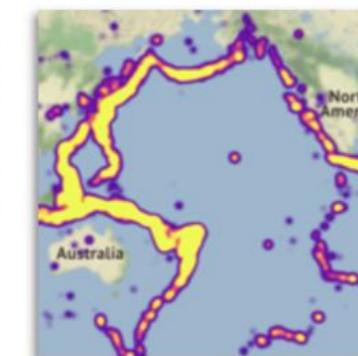
Lines on Mapbox



Filled Area on Maps

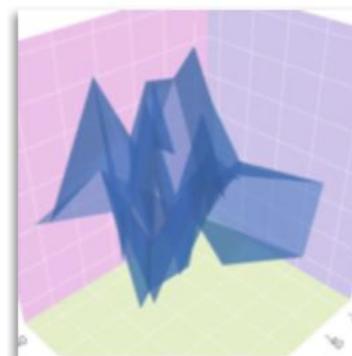


Bubble Maps

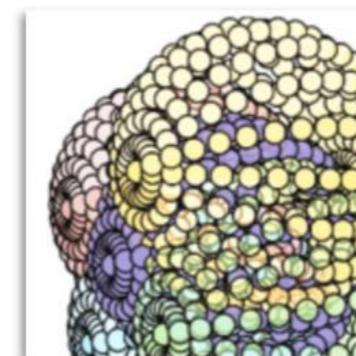


Mapbox Density
Heatmap

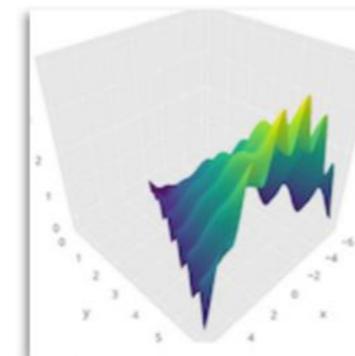
3D Charts

[More 3D Charts ×](#)

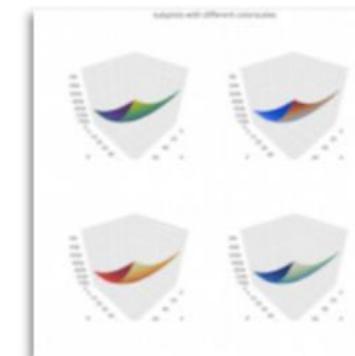
3D Axes



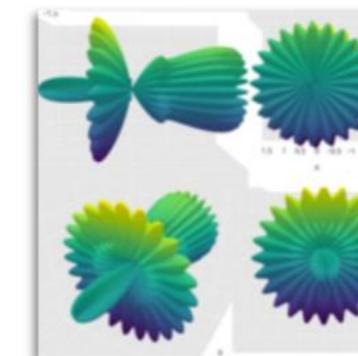
3D Scatter Plots



3D Surface Plots



3D Subplots



3D Camera Controls

<https://plotly.com/python/>



Python Plotly Graphing Library

Subplots



Mixed Subplots



Map Subplots

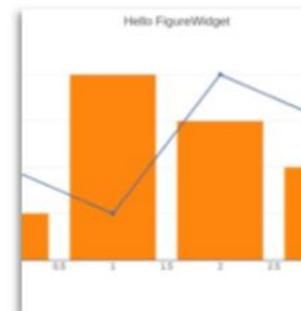


Table and Chart Subplots

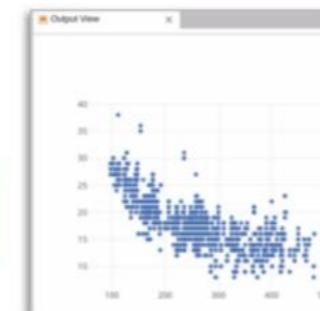


Figure Factory Subplots

Jupyter Widgets Interaction



Plotly FigureWidget Overview



Jupyter Lab with FigureWidget



Interactive Data Analysis with FigureWidget ipywidgets



Click Events

<https://plotly.com/python/>



Python Bokeh

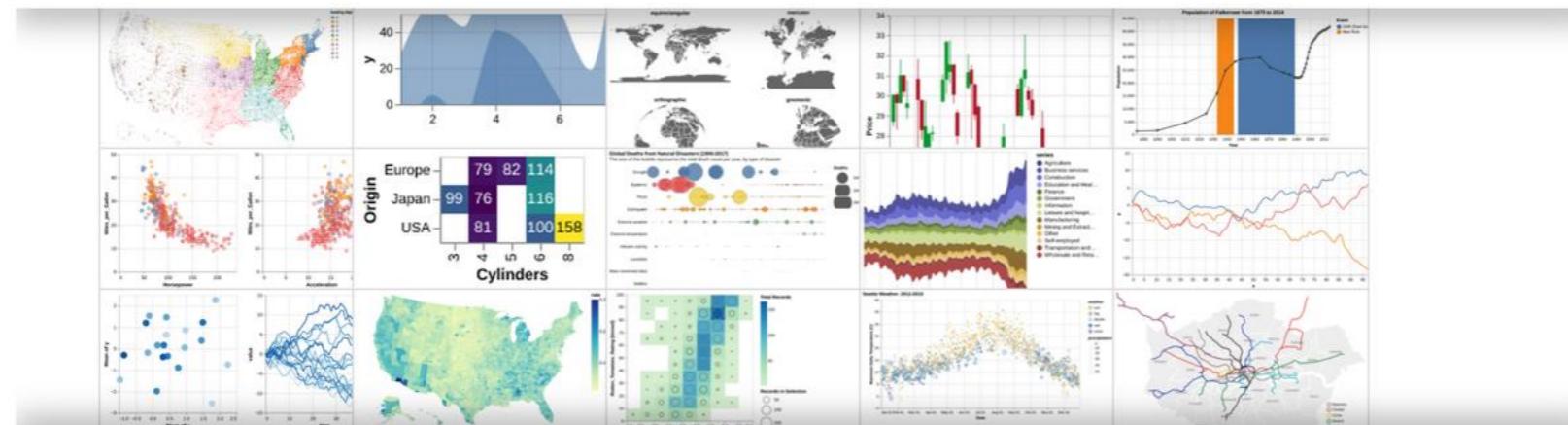




Python Altair



Vega-Altair: Declarative Visualization in Python



Vega-Altair is a declarative statistical visualization library for Python, based on [Vega](#) and [Vega-Lite](#).

The Vega-Altair open source project is not affiliated with Altair Engineering, Inc.

With Vega-Altair, you can spend more time understanding your data and its meaning. Altair's API is simple, friendly and consistent and built on top of the powerful [Vega-Lite](#) visualization grammar. This elegant simplicity produces beautiful and effective visualizations with a minimal amount of code.

You can browse this documentation via the links in the top navigation panel. In addition to reading this documentation page, it can be helpful to also browse the [Vega-Lite documentation](#).

<https://altair-viz.github.io/>

Iris flower data set

setosa



versicolor



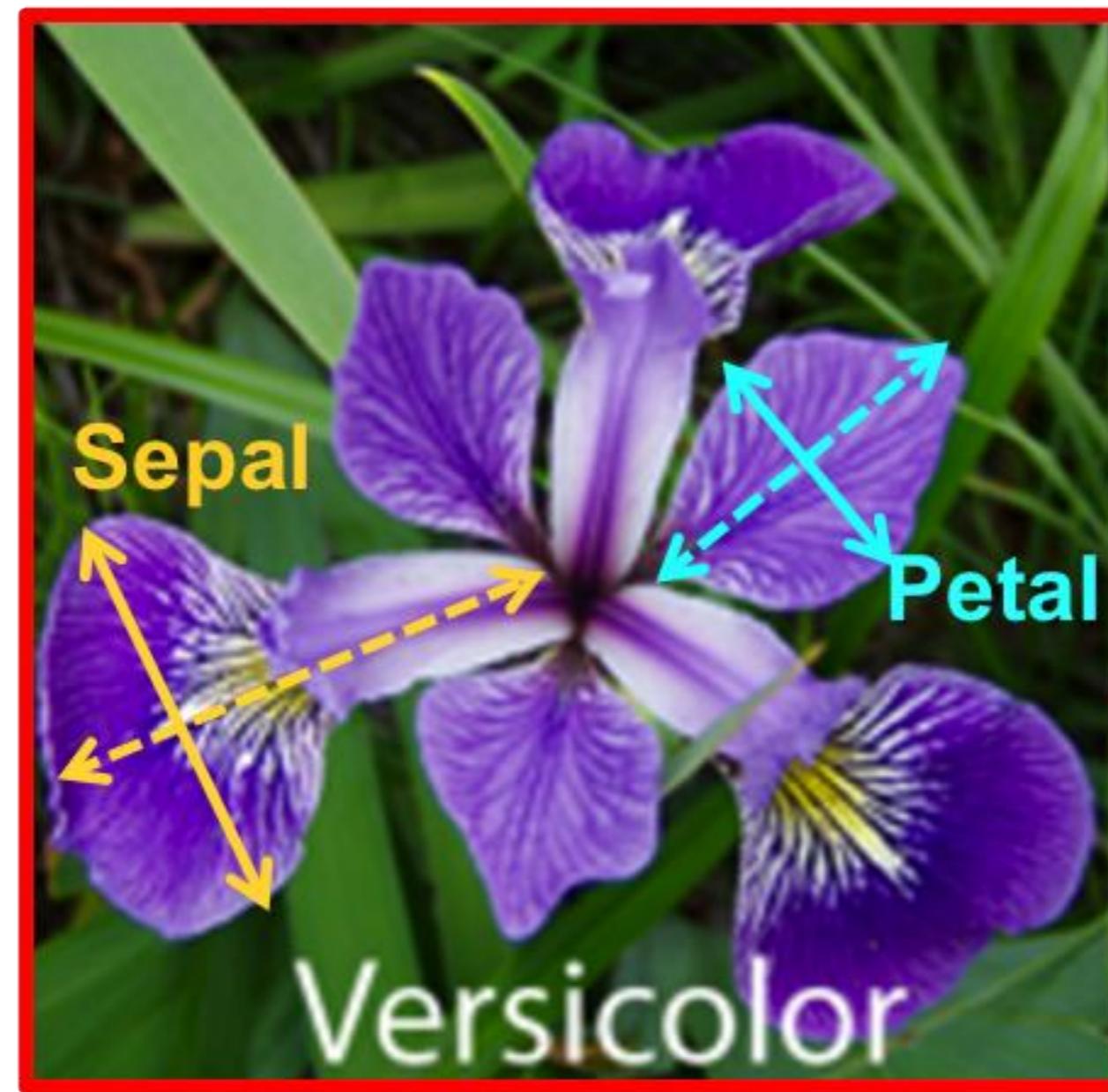
virginica



Source: https://en.wikipedia.org/wiki/Iris_flower_data_set

Source: <http://suruchifialoke.com/2016-10-13-machine-learning-tutorial-iris-classification/>

Iris Classification



Source: <http://suruchifialoke.com/2016-10-13-machine-learning-tutorial-iris-classification/>

iris.data

<https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data>

5.1,3.5,1.4,0.2,Iris-setosa
4.9,3.0,1.4,0.2,Iris-setosa
4.7,3.2,1.3,0.2,Iris-setosa
4.6,3.1,1.5,0.2,Iris-setosa
5.0,3.6,1.4,0.2,Iris-setosa
5.4,3.9,1.7,0.4,Iris-setosa
4.6,3.4,1.4,0.3,Iris-setosa
5.0,3.4,1.5,0.2,Iris-setosa
4.4,2.9,1.4,0.2,Iris-setosa
4.9,3.1,1.5,0.1,Iris-setosa
5.4,3.7,1.5,0.2,Iris-setosa
4.8,3.4,1.6,0.2,Iris-setosa
4.8,3.0,1.4,0.1,Iris-setosa
4.3,3.0,1.1,0.1,Iris-setosa
5.8,4.0,1.2,0.2,Iris-setosa
5.7,4.4,1.5,0.4,Iris-setosa
5.4,3.9,1.3,0.4,Iris-setosa
5.1,3.5,1.4,0.3,Iris-setosa
5.7,3.8,1.7,0.3,Iris-setosa
5.1,3.8,1.5,0.3,Iris-setosa
5.4,3.4,1.7,0.2,Iris-setosa
5.1,3.7,1.5,0.4,Iris-setosa
4.6,3.6,1.0,0.2,Iris-setosa
5.1,3.3,1.7,0.5,Iris-setosa
4.8,3.4,1.9,0.2,Iris-setosa
5.0,3.0,1.6,0.2,Iris-setosa
5.0,3.4,1.6,0.4,Iris-setosa
5.2,3.5,1.5,0.2,Iris-setosa
5.2,3.4,1.4,0.2,Iris-setosa
4.7,3.2,1.6,0.2,Iris-setosa
4.8,3.1,1.6,0.2,Iris-setosa
5.4,3.4,1.5,0.4,Iris-setosa

setosa



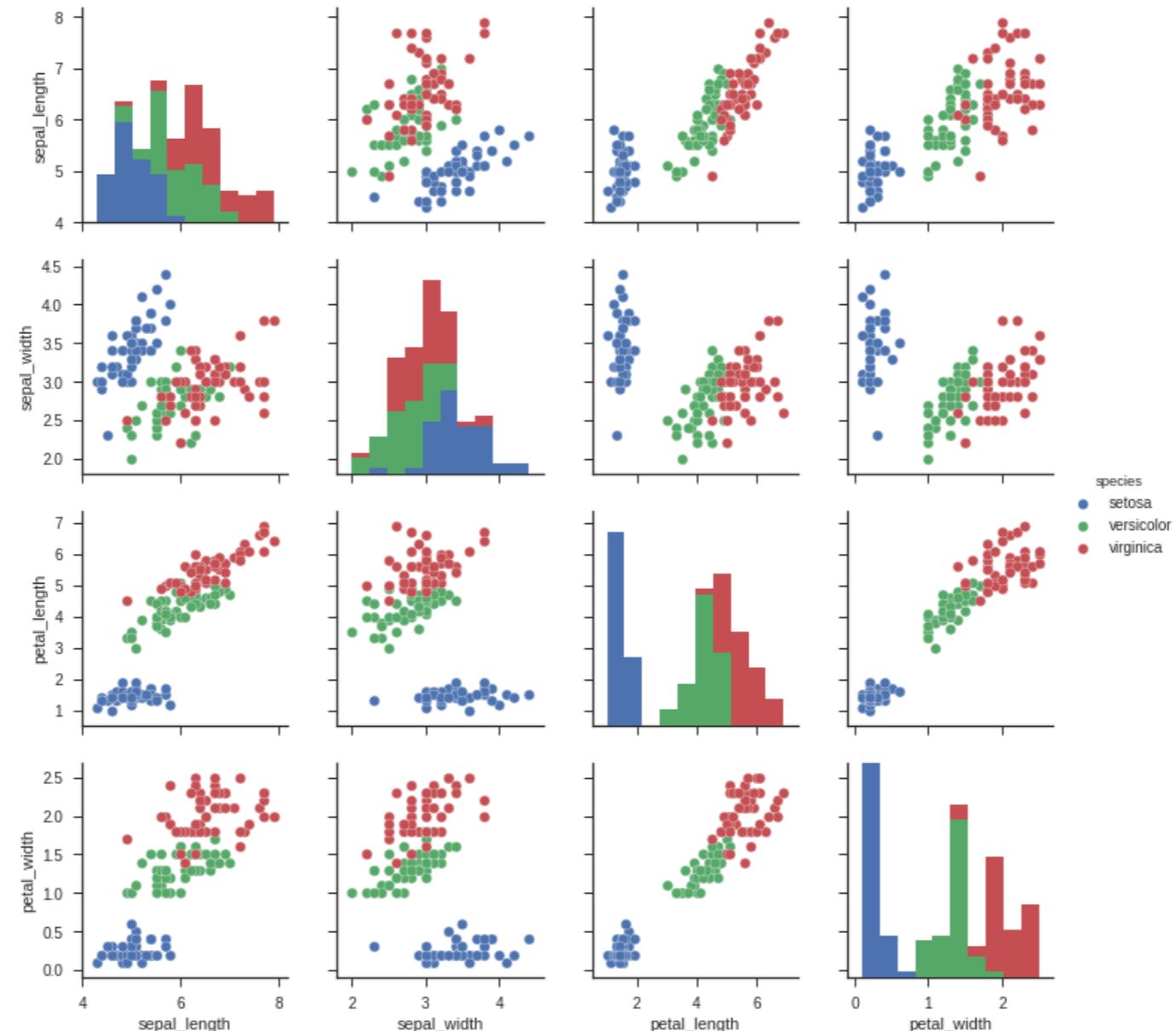
virginica



versicolor



Iris Data Visualization



Source: <https://seaborn.pydata.org/generated/seaborn.pairplot.html>

Data Visualization in Google Colab

The screenshot shows a Google Colab notebook interface. The top bar includes the Colab logo, the file name 'python101.ipynb', and standard menu options: File, Edit, View, Insert, Runtime, Tools, Help, and a status message 'All changes saved'. On the right, there are buttons for Comment, Share, settings, and a large 'A' for accessibility.

The left sidebar features a 'Table of contents' section with various Python-related topics, including 'Python Data Visualization' which is currently selected. Below it is a list of sub-topics under 'Python Data Visualization'.

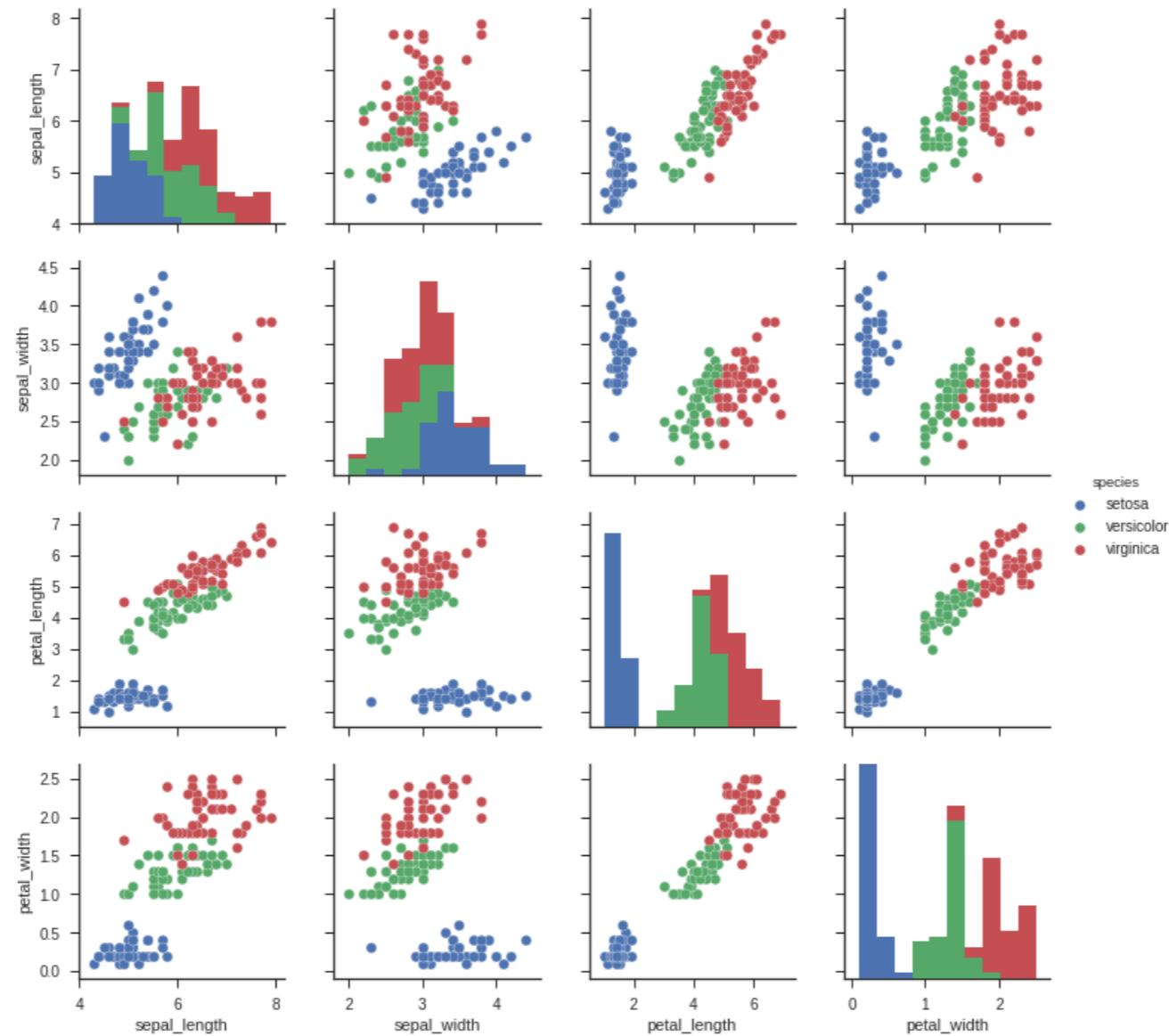
The main workspace contains a code cell [2] and a resulting visualization. The code cell contains:

```
[2]: 1 import seaborn as sns
2 sns.set(style="ticks", color_codes=True)
3 iris = sns.load_dataset("iris")
4 g = sns.pairplot(iris, hue="species")
```

The resulting visualization is a 3x3 grid of plots (pairplot) for the Iris dataset. The diagonal elements show histograms for each feature: 'sepal_length', 'sepal_width', and 'petal'. The off-diagonal elements show scatter plots between pairs of features: (sepal_length, sepal_width), (sepal_length, petal), (sepal_width, petal), (petal, sepal_length), (petal, sepal_width), and (petal, petal). Each scatter plot is colored by the 'species' column, with three distinct colors: blue for 'setosa', orange for 'versicolor', and green for 'virginica'. A legend on the bottom right identifies the species by color.

<https://tinyurl.com/aintpuppython101>

```
import seaborn as sns
sns.set(style="ticks", color_codes=True)
iris = sns.load_dataset("iris")
g = sns.pairplot(iris, hue="species")
```



Source: <https://seaborn.pydata.org/generated/seaborn.pairplot.html>

```
import numpy as np
import pandas as pd
%matplotlib inline
import matplotlib.pyplot as plt
import seaborn as sns
from pandas.plotting import scatter_matrix
```

```
# Import Libraries
import numpy as np
import pandas as pd
%matplotlib inline
import matplotlib.pyplot as plt
import seaborn as sns
from pandas.plotting import scatter_matrix
print('imported')
```

imported

```
url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width', 'class']
df = pd.read_csv(url, names=names)
print(df.head(10))
```

```
# Load dataset
url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width', 'class']
df = pd.read_csv(url, names=names)
print(df.head(10)).
```

	sepal-length	sepal-width	petal-length	petal-width	class
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa
5	5.4	3.9	1.7	0.4	Iris-setosa
6	4.6	3.4	1.4	0.3	Iris-setosa
7	5.0	3.4	1.5	0.2	Iris-setosa
8	4.4	2.9	1.4	0.2	Iris-setosa
9	4.9	3.1	1.5	0.1	Iris-setosa

df.tail(10)

```
print(df.tail(10))
```

	sepal-length	sepal-width	petal-length	petal-width	class
140	6.7	3.1	5.6	2.4	Iris-virginica
141	6.9	3.1	5.1	2.3	Iris-virginica
142	5.8	2.7	5.1	1.9	Iris-virginica
143	6.8	3.2	5.9	2.3	Iris-virginica
144	6.7	3.3	5.7	2.5	Iris-virginica
145	6.7	3.0	5.2	2.3	Iris-virginica
146	6.3	2.5	5.0	1.9	Iris-virginica
147	6.5	3.0	5.2	2.0	Iris-virginica
148	6.2	3.4	5.4	2.3	Iris-virginica
149	5.9	3.0	5.1	1.8	Iris-virginica

df.describe()

```
print(df.describe())
```

	sepal-length	sepal-width	petal-length	petal-width
count	150.000000	150.000000	150.000000	150.000000
mean	5.843333	3.054000	3.758667	1.198667
std	0.828066	0.433594	1.764420	0.763161
min	4.300000	2.000000	1.000000	0.100000
25%	5.100000	2.800000	1.600000	0.300000
50%	5.800000	3.000000	4.350000	1.300000
75%	6.400000	3.300000	5.100000	1.800000
max	7.900000	4.400000	6.900000	2.500000

```
print(df.info())
print(df.shape)
```

```
print(df.info())

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 5 columns):
sepal-length      150 non-null float64
sepal-width       150 non-null float64
petal-length      150 non-null float64
petal-width       150 non-null float64
class             150 non-null object
dtypes: float64(4), object(1)
memory usage: 5.9+ KB
None
```

```
print(df.shape)
```

```
(150, 5)
```

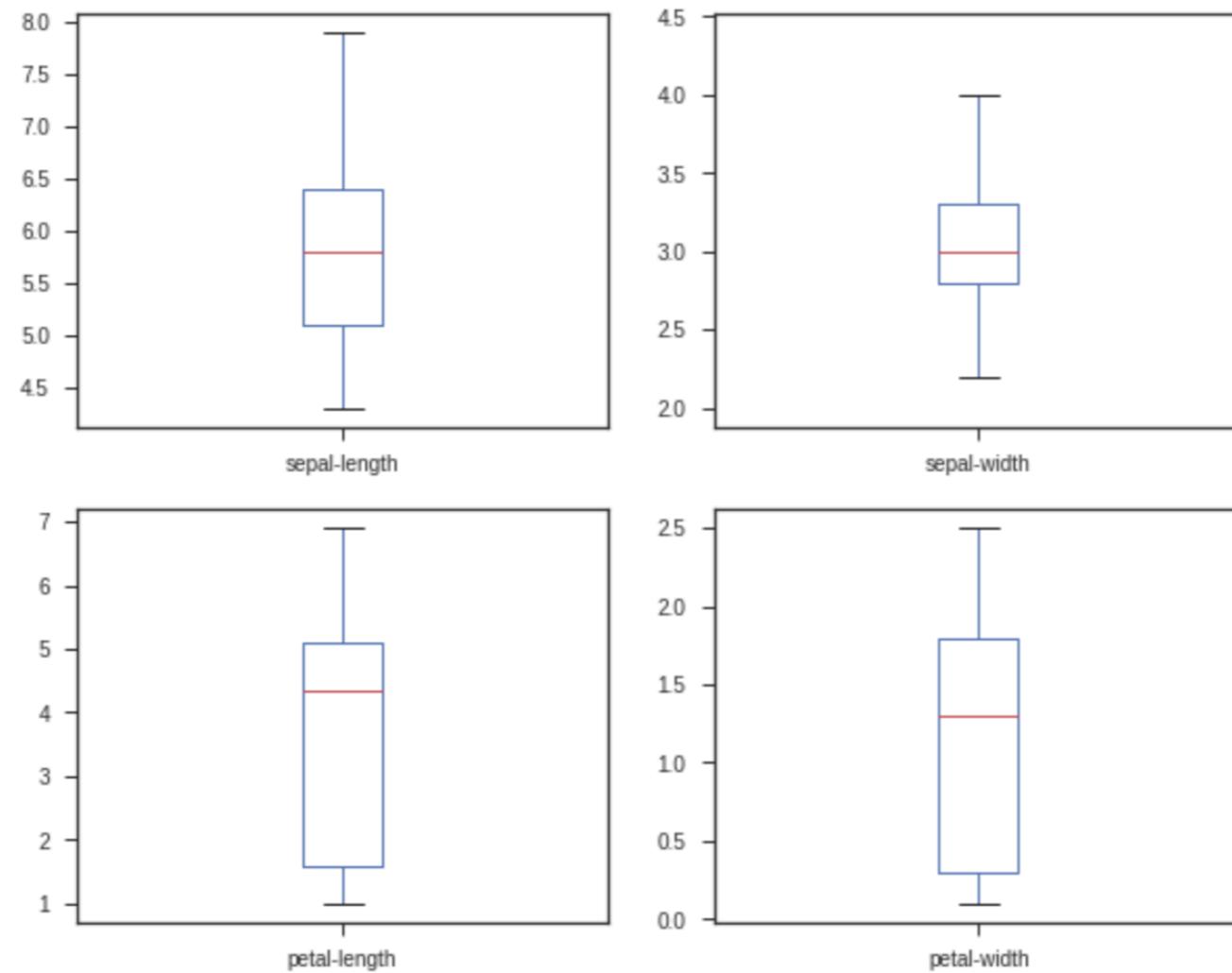
```
df.groupby('class').size()
```

```
print(df.groupby('class').size())
```

```
class
Iris-setosa      50
Iris-versicolor  50
Iris-virginica   50
dtype: int64
```

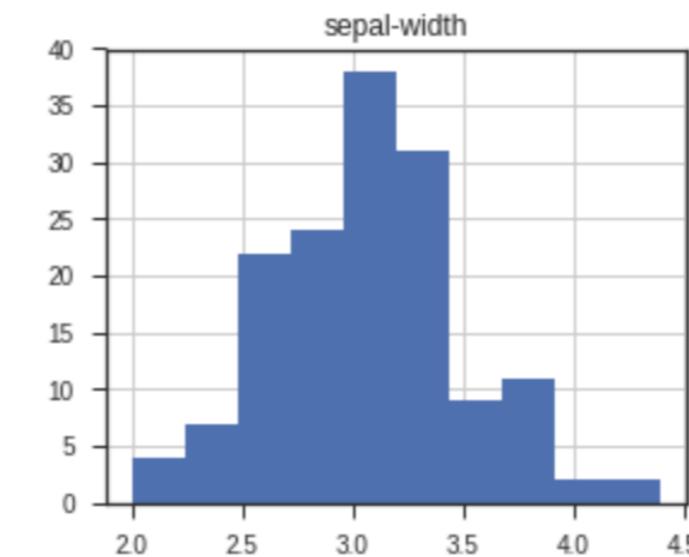
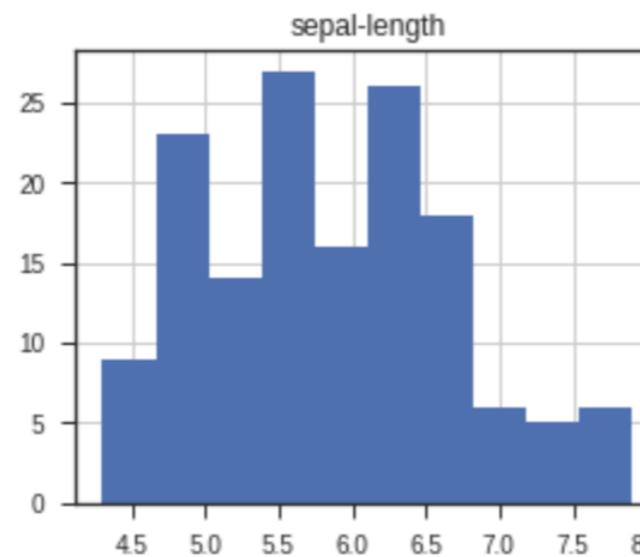
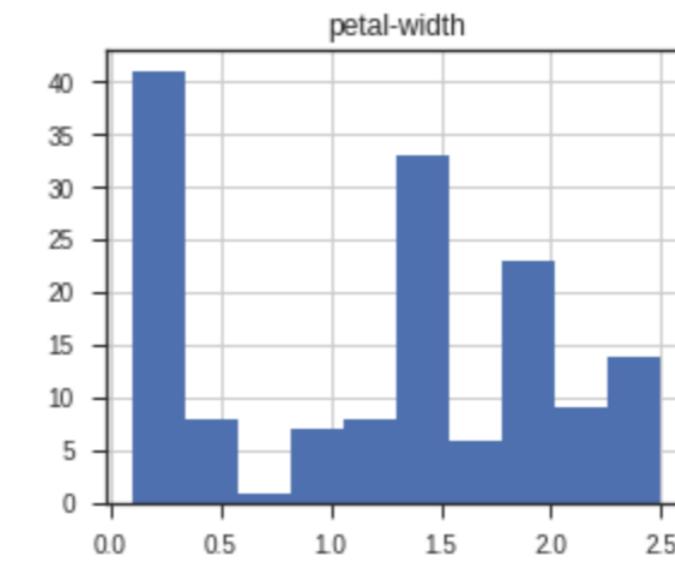
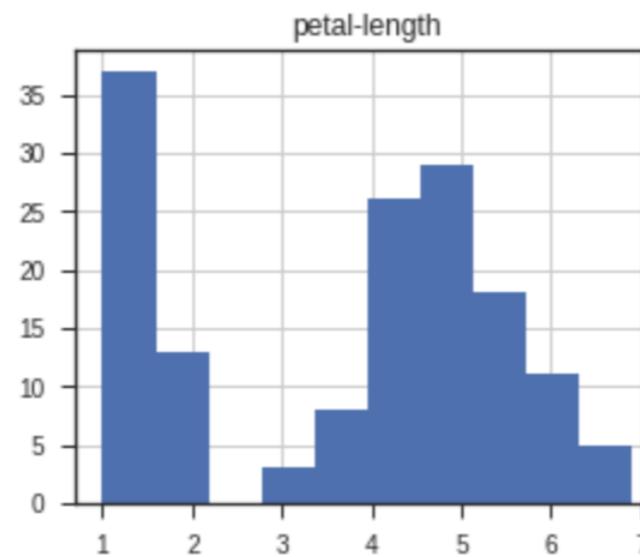
```
plt.rcParams["figure.figsize"] = (10,8)
df.plot(kind='box', subplots=True, layout=(2,2), sharex=False, sharey=False)
plt.show()
```

```
plt.rcParams["figure.figsize"] = (10,8)
df.plot(kind='box', subplots=True, layout=(2,2), sharex=False, sharey=False)
plt.show()
```



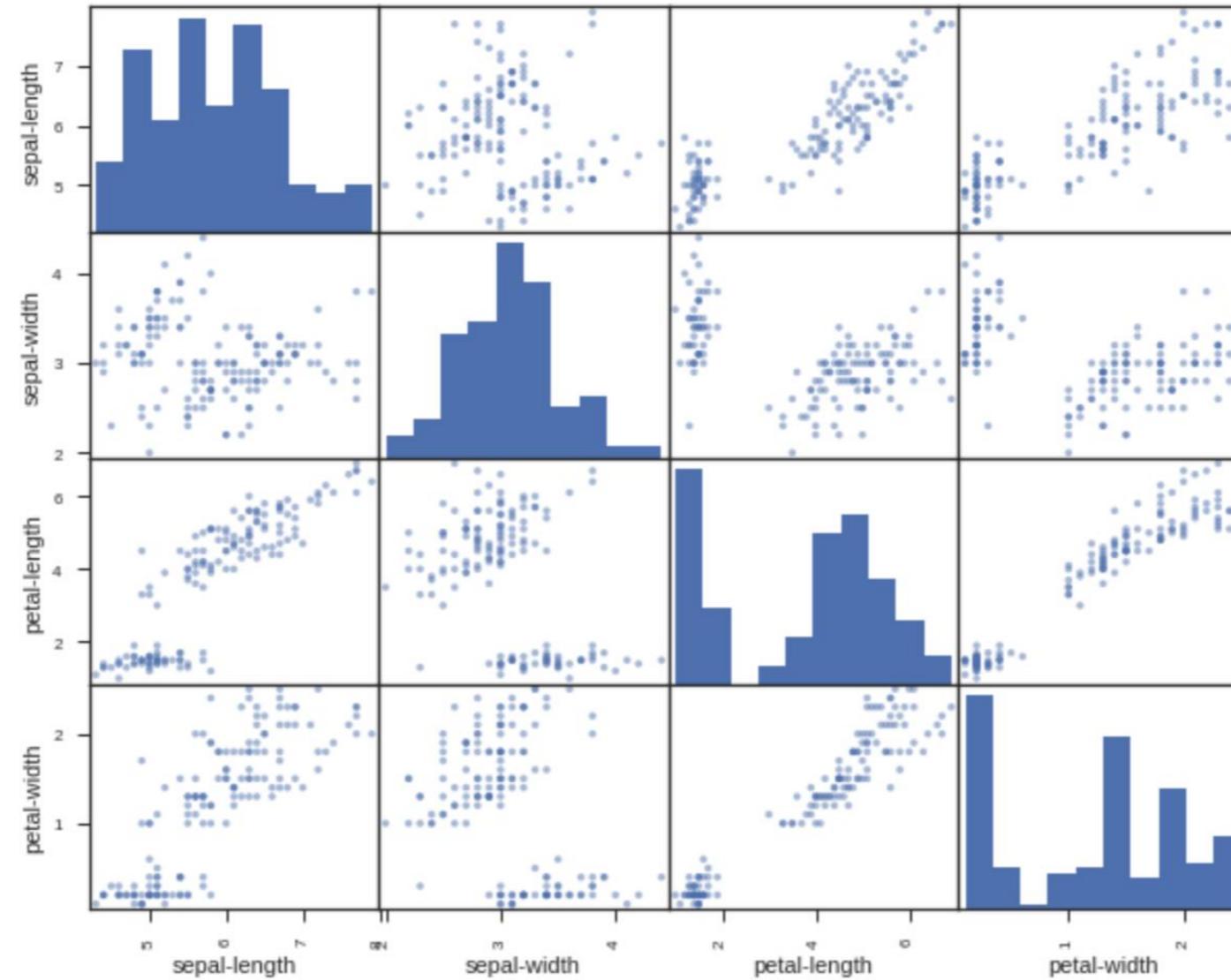
df.hist() plt.show()

```
df.hist()  
plt.show()
```



```
scatter_matrix(df)  
plt.show()
```

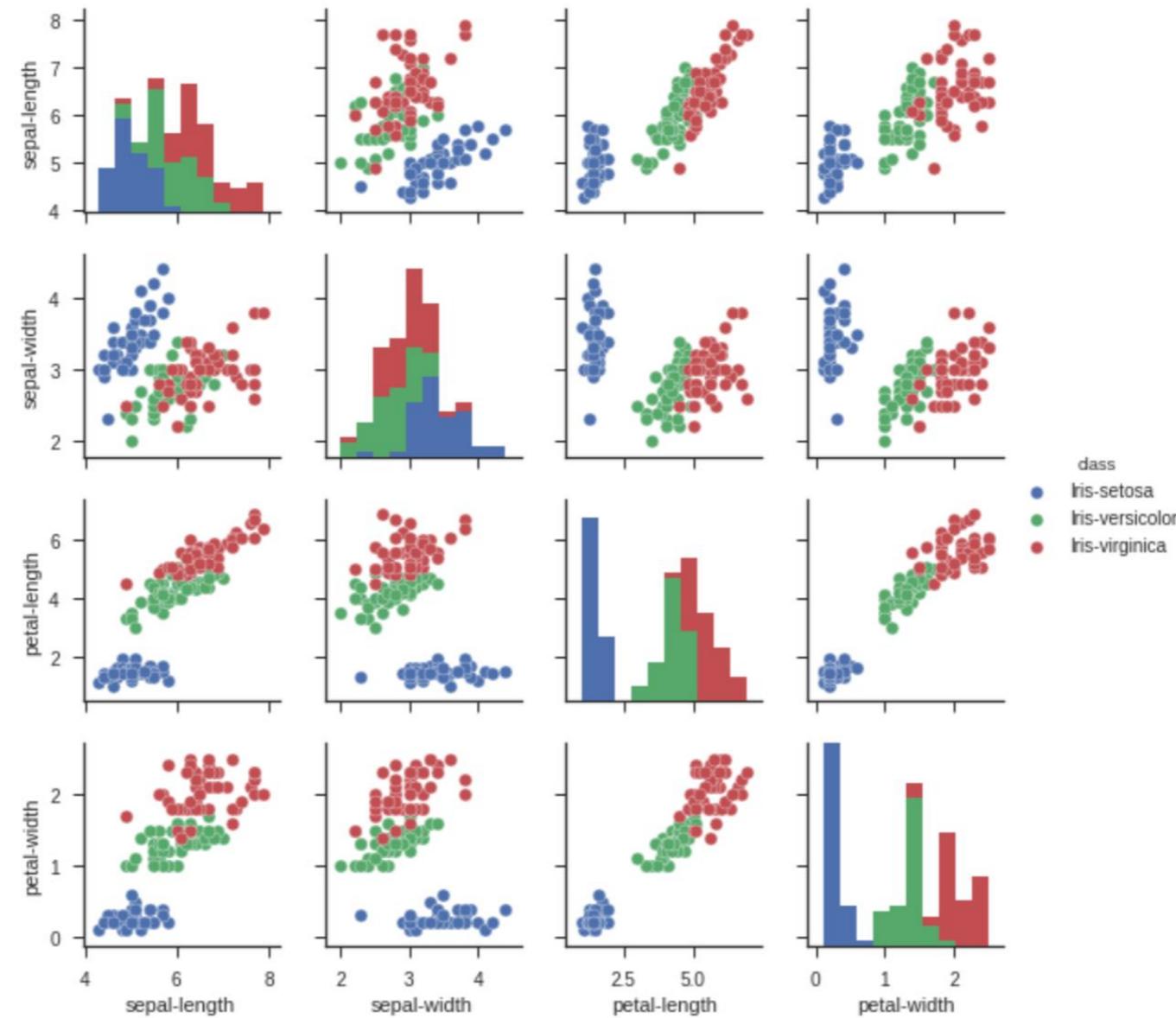
```
scatter_matrix(df)  
plt.show(.)
```



```
sns.pairplot(df, hue="class", size=2)
```

```
sns.pairplot(df, hue="class", size=2)
```

```
<seaborn.axisgrid.PairGrid at 0x7f1d21267390>
```



Wes McKinney (2022), "Python for Data Analysis: Data Wrangling with pandas, NumPy, and Jupyter", 3rd Edition, O'Reilly Media.

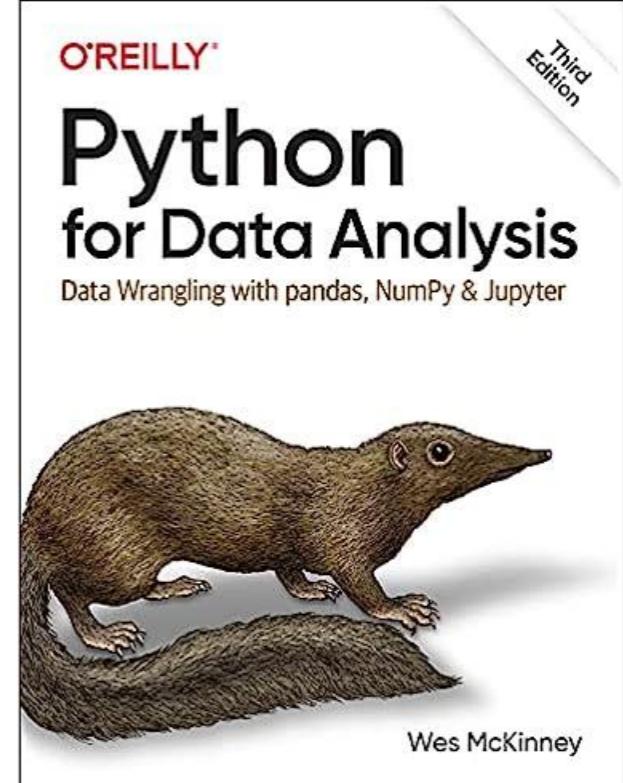
[wesm/pydata-book](https://github.com/wesm/pydata-book) Public

Code Issues 2 Pull requests 1 Actions Projects Wiki Security Insights

3rd-edition 3 branches 0 tags Go to file Code About

File	Description	Time
wesm	Upload cleaner notebook files without internal build toolchai...	3 days ago
datasets	Add fec.parquet	10 months ago
examples	Simplifying datasets	10 months ago
.gitignore	Add gitignore	8 years ago
COPYING	Update COPYING	4 months ago
README.md	Update notebooks in advance of publication	7 months ago
appa.ipynb	Upload cleaner notebook files without internal build toolchai...	3 days ago
appb.ipynb	Upload cleaner notebook files without internal build toolchai...	3 days ago
ch02.ipynb	Upload cleaner notebook files without internal build toolchai...	3 days ago
ch03.ipynb	Upload cleaner notebook files without internal build toolchai...	3 days ago
ch04.ipynb	Upload cleaner notebook files without internal build toolchai...	3 days ago
ch05.ipynb	Upload cleaner notebook files without internal build toolchai...	3 days ago
ch06.ipynb	Upload cleaner notebook files without internal build toolchai...	3 days ago

O'REILLY
Python for Data Analysis
Data Wrangling with pandas, NumPy & Jupyter
Third Edition
Wes McKinney



<https://github.com/wesm/pydata-book>

Wes McKinney (2022), "Python for Data Analysis: Data Wrangling with pandas, NumPy, and Jupyter", 3rd Edition, O'Reilly Media.

3rd-edition ▾ pydata-book / ch04.ipynb Go to file ...

wesm Upload cleaner notebook files without internal build toolchain instru... ... Latest commit f1757b8 3 days ago ⏲ History

3 contributors

1224 lines (1224 sloc) | 21.9 KB

In [1]:

```
import numpy as np
np.random.seed(12345)
import matplotlib.pyplot as plt
plt.rc("figure", figsize=(10, 6))
np.set_printoptions(precision=4, suppress=True)
```

In [2]:

```
import numpy as np

my_arr = np.arange(1_000_000)
my_list = list(range(1_000_000))
```

In [3]:

```
%timeit my_arr2 = my_arr * 2
%timeit my_list2 = [x * 2 for x in my_list]
```

In [4]:

```
import numpy as np
data = np.array([[1.5, -0.1, 3], [0, -3, 6.5]])
data
```

Source: <https://github.com/wesm/pydata-book/blob/3rd-edition/ch04.ipynb>

Python in Google Colab

The screenshot shows the Google Colab interface with the following code examples:

```
# Future Value
pv = 100
r = 0.1
n = 7
fv = pv * ((1 + (r)) ** n)
print(round(fv, 2))

[11] amount = 100
interest = 10 #10% = 0.01 * 10
years = 7
future_value = amount * ((1 + (0.01 * interest)) ** years)
print(round(future_value, 2))

# Python Function def
def getfv(pv, r, n):
    fv = pv * ((1 + (r)) ** n)
    return fv
fv = getfv(100, 0.1, 7.)
print(round(fv, 2))

[13] # Python if else
score = 80
if score >=60 :
    print("Pass")
else:
    print("Fail")

[14] Pass
```

Python in Google Colab

Table of contents X + Code + Text ✓ RAM Disk Editing ^

Python101
Python File Input / Output
OS, IO, files, and Google Drive
Python Try Except
Python Class
Python Programming
Pythong String and Text
Python Numpy
Python Pandas

Python Data Visualization

Machine Learning with scikit-learn
Classification and Prediction
K-Means Clustering
Deep Learning for Financial Time Series Forecasting
Portfolio Optimization and Algorithmic Trading
Investment Portfolio Optimisation with Python
Efficient Frontier Portfolio Optimisation in Python

[2] 1 import seaborn as sns
2 sns.set(style="ticks", color_codes=True)
3 iris = sns.load_dataset("iris")
4 g = sns.pairplot(iris, hue="species")

sepal_length
sepal_width
petal_length
petal_width

species
setosa
versicolor
virginica

Papers with Code State-of-the-Art (SOTA)

Computer Vision



Semantic Segmentation

185 benchmarks

3397 papers with code



Image Classification

390 benchmarks

2778 papers with code



Object Detection

269 benchmarks

2559 papers with code



Contrastive Learning

2 benchmarks

1119 papers with code



Image Generation

208 benchmarks

1097 papers with code

▶ See all 1415 tasks

Natural Language Processing



Language Modelling

458 benchmarks

2248 papers with code



Question Answering

181 benchmarks

1818 papers with code



Machine Translation

78 benchmarks

1721 papers with code



Sentiment Analysis

87 benchmarks

1040 papers with code



Text Generation

242 benchmarks

931 papers with code

▶ See all 664 tasks

<https://paperswithcode.com/sota>

Summary

- NumPy
 - Numerical Python N-dimensional array
- Pandas
 - Data Analytics
- Matplotlib
 - Basic Data Visualization
- Seaborn
 - Advanced Visualization

**THANKS FOR YOUR
ATTENTION!**