

MBAZA NLP COMMUNITY

VIRTUAL TRAINING
22-24.06 | 3 - 5 PM

BASIC DATA WRANGLING WITH PANDAS
TEXT PRE-PROCESSING BASICS
NLP MODELING & ALGORITHMS

Why are we here?

Module 1

Basic Data Wrangling
with Pandas



Any Data Science work
you do in the future

Module 2

Text Pre-Processing
Basics



Any NLP work you do in
the future

Dataset cleaning
challenge in July

Module 3

NLP Modelling &
Algorithms



Your introduction to
Machine Learning

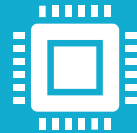
Community activities on
Chatbots and Voice

Mbaza NLP Community



Knowledge sharing

- NLP, e. g. speech recognition and chatbots
- Software architecture design
- DevOps



Access to:

- Developed NLP models
- Datasets
- Infrastructure: Dev/test environments, computing resources



Use case development

- Gain practice & improve your skills
- Practical experiences to show to employers
- Work on projects with rewards

Learning Outcomes of Today



You understand:

- what Jupyter Notebooks are and how to use them
- what the Pandas library does
- the basic Pandas data structures



You can:

- use Pandas to import CSV data
- apply basic methods of inspecting your data
- filter and select data
- combine data from various sources
- apply basic data manipulation methods

Jupyter Notebooks



- Web-based interactive coding environment for Python and other languages
- Also supports [Markdown](#) cells for descriptions and explanations

Code cell

Cell output

Markdown cell

A screenshot of a Jupyter Notebook interface. The top bar shows the Jupyter logo, the name 'matrix', and a 'Logout' button. Below the top bar is a menu bar with 'File', 'Edit', 'View', 'Insert', 'Cell', 'Kernel', 'Widgets', and 'Help'. A toolbar contains icons for saving, adding, deleting, and running cells. The notebook content area shows three code cells. The first cell contains 'import numpy as np'. The second cell contains 'np.__version__' and its output is '1.15.4'. The third cell contains code to create a 4x4 matrix and its output is a 4x4 array of floats. Below the code cells is a Markdown cell containing a note about calculating a 5000x5000 matrix.

```
jupyter matrix
```

```
File Edit View Insert Cell Kernel Widgets Help Trusted Python 3
```

```
In [1]: import numpy as np
```

```
In [2]: np.__version__
```

```
Out[2]: '1.15.4'
```

```
In [3]: n = 4 # n - order square matrix
        h = 0.5
        np.fromfunction(lambda i, j: h / ((i - j) ** 2 + h * h), (n, n), dtype=float)
```

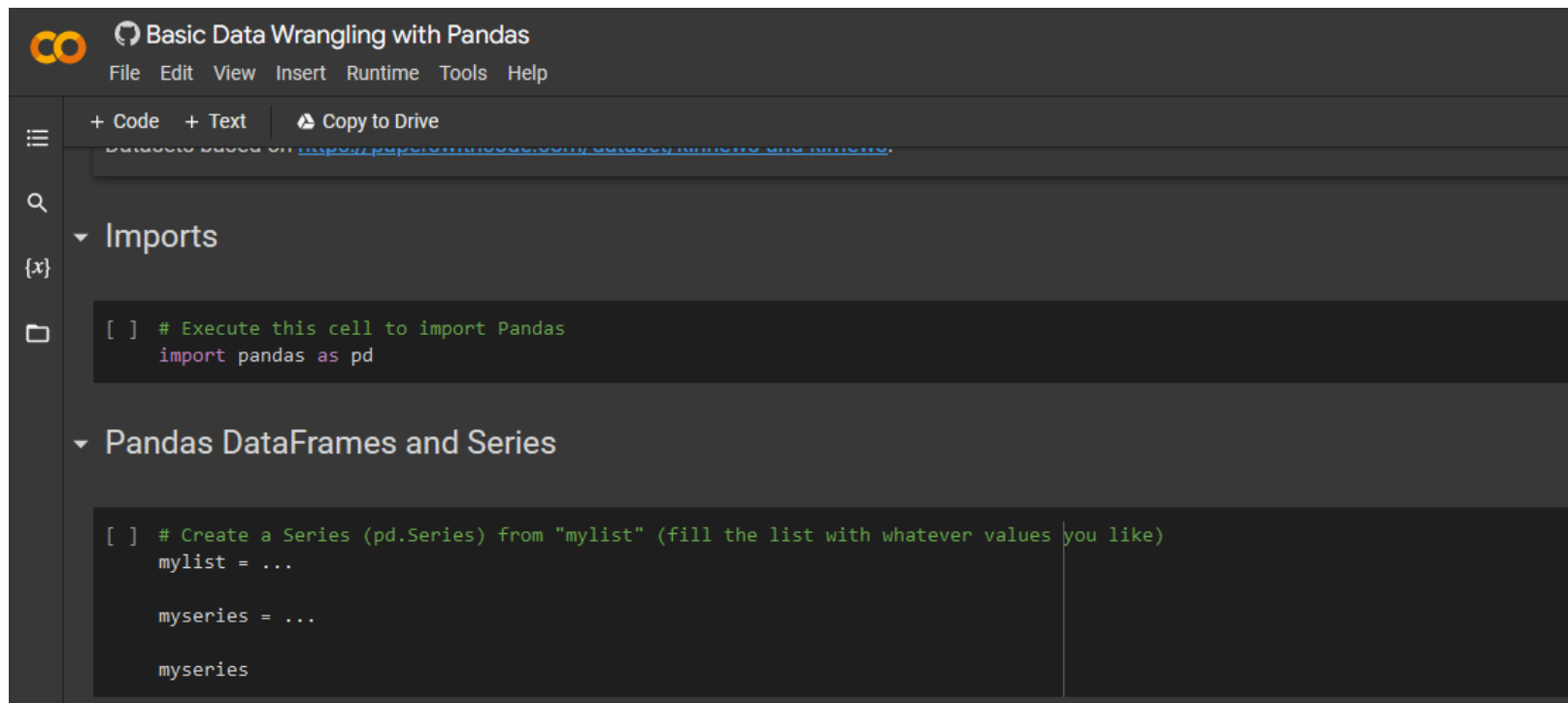
```
Out[3]: array([[2.         , 0.4         , 0.11764706, 0.05405405],
               [0.4         , 2.         , 0.4         , 0.11764706],
               [0.11764706, 0.4         , 2.         , 0.4         ],
               [0.05405405, 0.11764706, 0.4         , 2.         ]])
```

Note. To calculate the matrix size 5000x5000, need to replace n = 4 in the previous input cell on n = 5000 n to push Shift+Enter

Google Colab



- Jupyter Notebook hosted online by Google with some added features
- Your code is executed on Google's servers



The screenshot displays the Google Colab web interface. At the top, the title 'Basic Data Wrangling with Pandas' is visible, along with a menu bar containing 'File', 'Edit', 'View', 'Insert', 'Runtime', 'Tools', and 'Help'. Below the menu, there are tabs for '+ Code', '+ Text', and 'Copy to Drive'. The left sidebar shows a file explorer with a search icon and a folder icon. The main area contains two code cells. The first cell is under the 'Imports' section and contains the code:

```
[ ] # Execute this cell to import Pandas
import pandas as pd
```

. The second cell is under the 'Pandas DataFrames and Series' section and contains the code:

```
[ ] # Create a Series (pd.Series) from "mylist" (fill the list with whatever values you like)
mylist = ...

myseries = ...

myseries
```

let's get started!

Open the Colab link

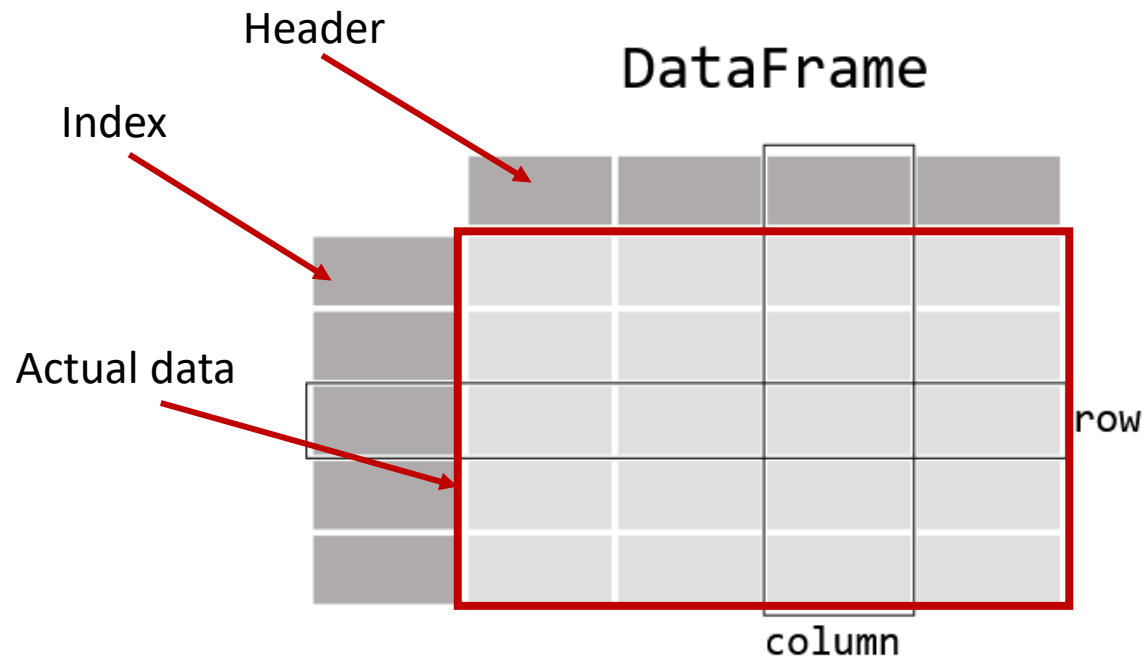
Pandas is the core library for Data Science with Python

- Open-source Python library for data analysis & management
- Provides data structures for table-like data as well as tools for data manipulation (e. g. filtering, merging, or sorting).



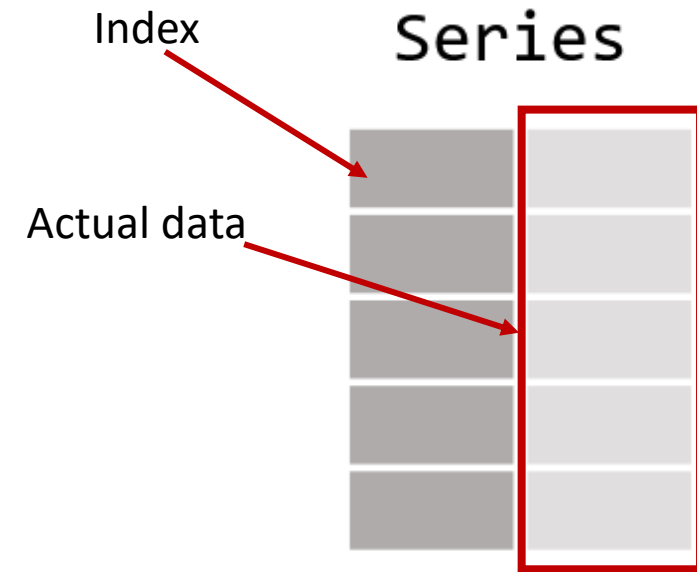
Pandas has two basic data structures

DataFrame: Two-dimensional,
can be thought of as table



```
df = pd.DataFrame({'number': [1, 2, 3],  
                  'squared': [1, 4, 9]})
```

Series: One-dimensional,
more similar to a list



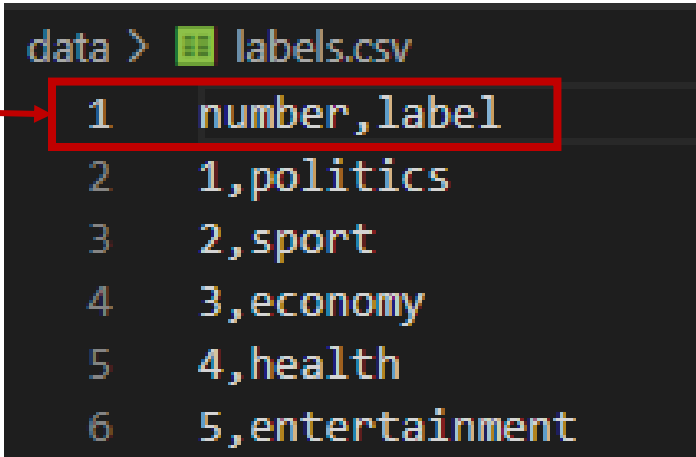
```
s = pd.Series([1, 2, 3])
```

Importing data

- A common file format in data science is **CSV** (comma-separated values)

Header row →

```
data > labels.csv
1 number,label
2 1,politics
3 2,sport
4 3,economy
5 4,health
6 5,entertainment
```



- Pandas makes it easy to import many file formats, including CSVs

```
data = pd.read_csv('filepath')
```

Understanding your data

There are several ways to inspect your data

1. Console output: Putting the variable name holding your data as the last row of the notebook cell prints it out

```
data
```

✓ 13.8s

	label	title
0	2	ikipe y' u rwanda amavubi yahesheje u rwanda a...
1	11	urubyiruko itorero erc giterane cy'ububyutse k...
2	4	rusizi bambaye udupfukamunwa n'ubwo bamwe bata

2. The **shape** attribute of your data variable shows you the number of rows and columns of a DataFrame

```
data.shape
```

✓ 1.8s

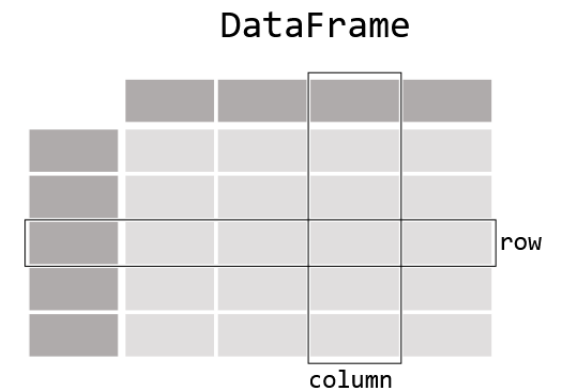
(500, 3)

3. The **head(x)** method returns the first x rows of your data

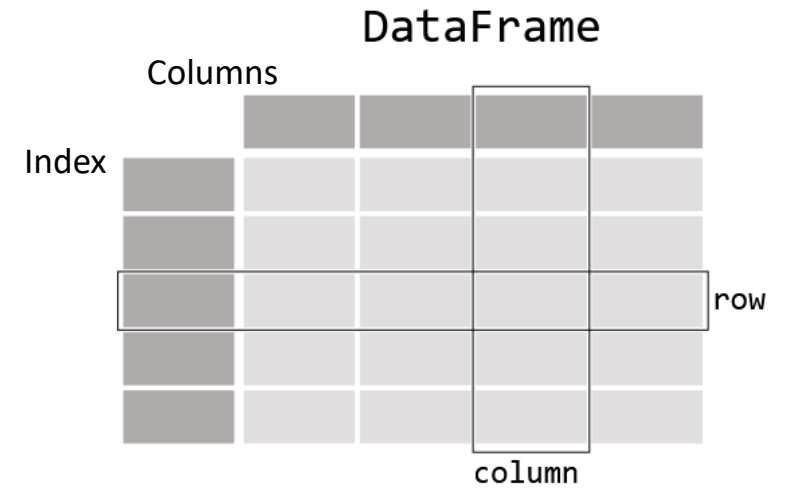
```
data.head(3)
```

✓ 0.6s

	label	title
0	2	ikipe y' u rwanda amavubi yahesheje u rwanda a...
1	11	urubyiruko itorero erc giterane cy'ububyutse k...
2	4	rusizi bambaye udupfukamunwa n'ubwo bamwe bata...



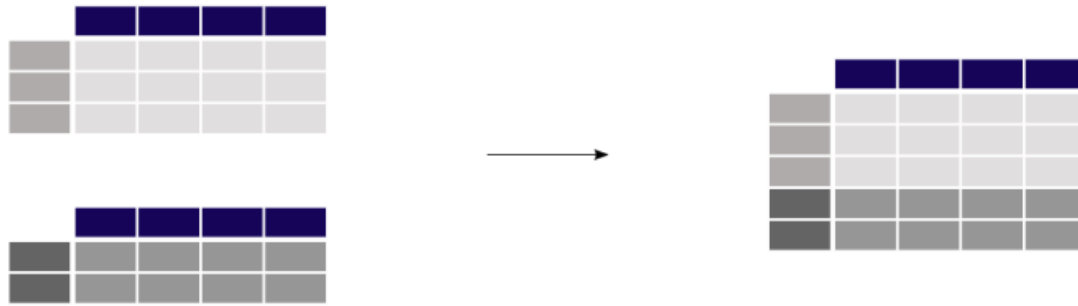
Filtering & selecting data



Task	Code	Output
Select one column of a DataFrame	<code>data["column_name"]</code>	Series
Select rows of a DataFrame by their Index number	<code>data[index_start:index_end]</code>	DataFrame
Select multiple columns of a DataFrame in the order you want them	<code>data[["column A", "column B"]]</code>	DataFrame
Select rows from a specific column	<code>data["column_name"][index_start:index_end]</code>	Series
Select only rows where some condition is met	<code>data[data["column_name"] == 5]</code>	DataFrame

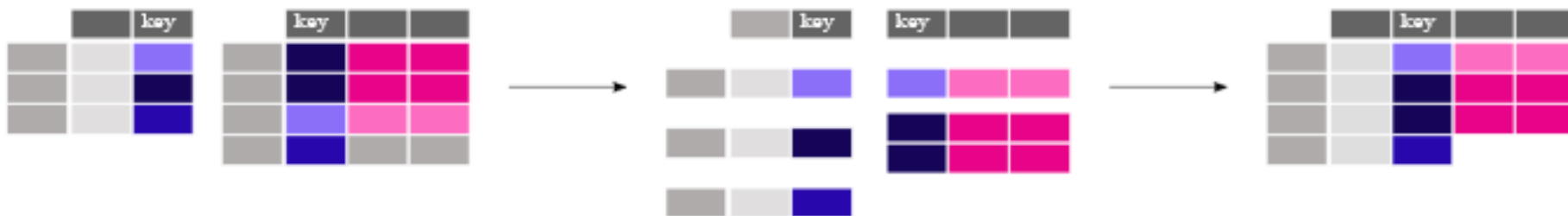
Combining data from multiple source

Concatenating two DataFrames – `pd.concat()`



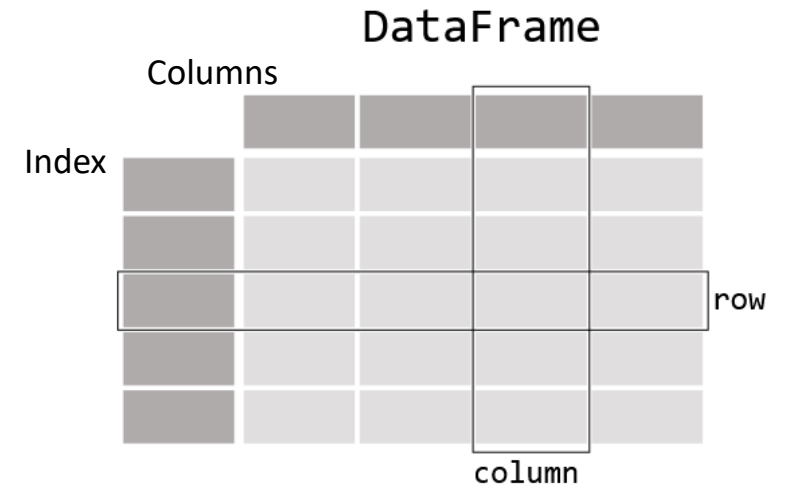
```
pd.concat([data, data2])
```

Join DataFrames using a common identifier – `pd.merge()`



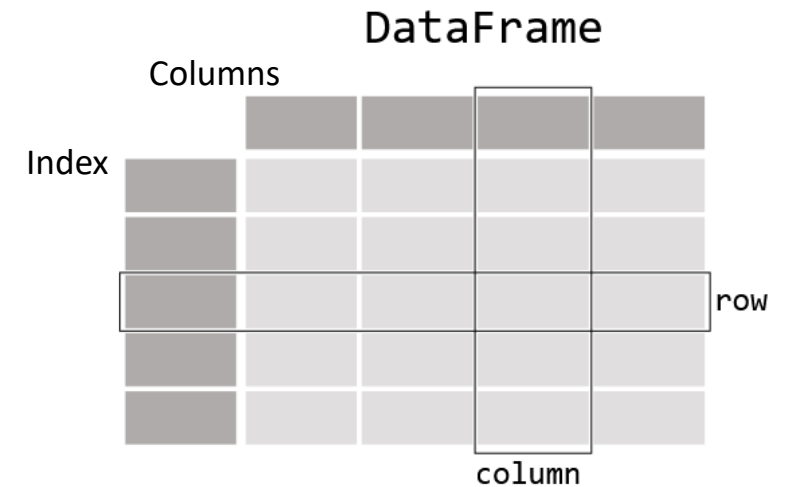
```
pd.merge(all_data, text_labels, on='label')
```

Basic DataFrame operations



Task	Code	Output
Add column	<code>data["column_name"] = pd.Series(...)</code>	
Remove column	<code>data.drop("column_name", inplace=True)</code> <i>or</i> <code>data = data.drop("column_name")</code>	DataFrame
Text functions on columns	<code>data["column_name"].str.len()</code> <code>.contains()</code> <code>.lower()</code> <code>.replace()</code> <code>.split()</code>	Series

Basic DataFrame operations



Task	Code	Output
Aggregation functions on columns	<code>data["column_name"].mean()</code> <code>.median()</code> <code>.std()</code> <code>.min()</code> <code>.max()</code>	float
Remove duplicates	<code>data.drop_duplicates("column_name", inplace=True)</code> <i>or</i> <code>data = data.drop_duplicates("column_name")</code>	DataFrame
Sorting by column	<code>data.sort_values(by='column_name', ascending=True)</code>	DataFrame
Sorting with function	<code>data.sort_values(by='content', key=lambda col: col.str.len())</code>	DataFrame

Report back

Q&A

Summary

- **Jupyter Notebooks** provide a web-based interactive coding environment
- **Pandas** is an open-source Python library for data analysis & management
- It provides two basic data structures: **DataFrame** (two-dimensional) and **Series** (one-dimensional)
- `pd.read_csv()` allows you to import CSV files
- The `shape` attribute and `head()` method help you understand your data
- Pandas offers many methods to select data based on squared brackets:
`data["column_name"][index_start:index_end]`
- Combine data by concatenation (`pd.concat()`) or via common identifier (`pd.merge()`)
- Basic operations: `.sort_values()` for sorting, `.drop_duplicates()` for dup removal
- Aggregation functions: `.mean()`, `.median()`, `.std()`, ...
- Text functions: `.str.len()`, `.str.contains()`, `.str.lower()`, `.str.replace()`, ...

Outlook for Tomorrow



Understand:

- what NLP is good for
- why text pre-processing is crucial for NLP
- what parallel datasets are
- what regular expressions are good for



Know:

- how to cHaNgE cAsInG
- how to remove text parts
- how to handle empty values
- how to use regular expressions in Python
- how to deal with parallel datasets and their challenges

Join the Mbaza NLP Community!

WhatsApp

<https://chat.whatsapp.com/BRlxzsFiZgsLmK5SBT2XUo>



Slack

https://join.slack.com/t/mbazanlpcommunity/shared_invite/zt-19ie5idhj-f0yWfOBgTKzs7VOKCcr_pw



GitHub

<https://github.com/MBAZA-NLP>



Hugging face

<https://huggingface.co/organizations/mbazaNLP/share/mUKyOkYpSRisRpspbfuwUvoQgWyfdiJYqU>

