# Programming for Data Science

Lab01-Getting started
Instructor: Jakramate Bootkrajang

# Outline

- Python distribution

- Programming environments

- Google colaboration

  – Some exercises

- Interesting websites

# About Python

- Python 2.xx
  - Older version of Python
  - Still being used in legacy programs
- Python 3.xx
  - Newer version
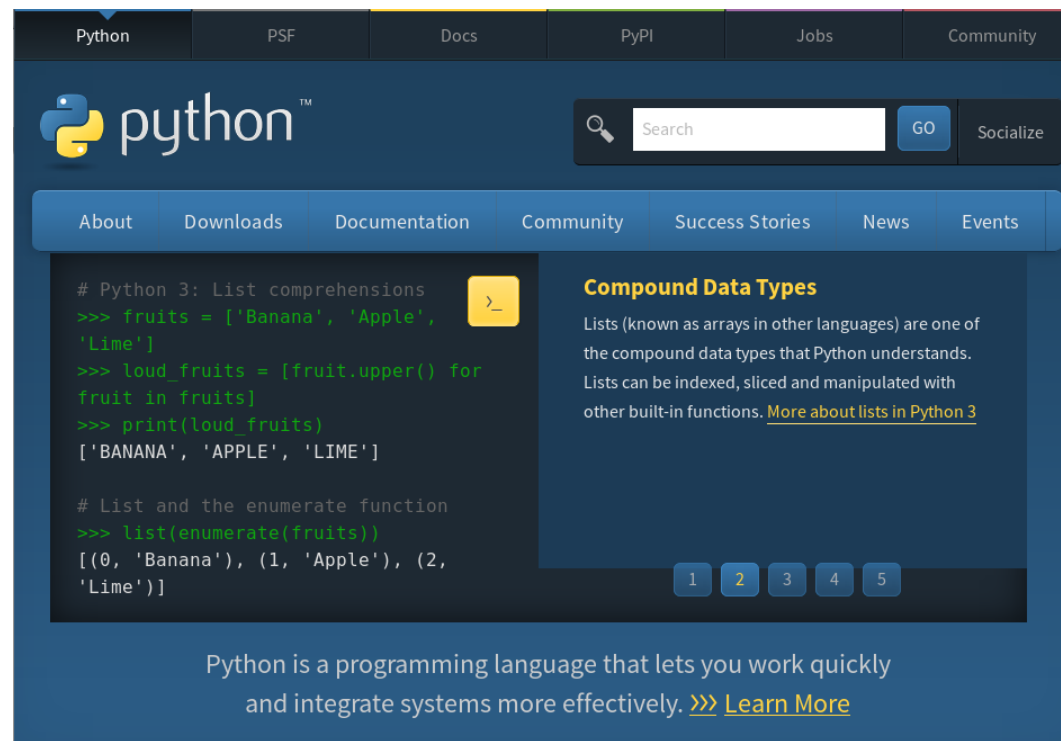  - We will use Python 3.xx in this class

# Python distribution

- A distribution is
  - A set of Python interpreter + additional packages
- The most widely used ones are
  - CPython distribution (standard)
  - Anaconda distribution (CPython + packages for data science)

# CPython distribution

- Standard distribution
- Can be downloaded from www.python.org

# Anaconda distribution

- Standard distribution + appox. 1400 packages

- Those packages are essential for data science, data mining, machine learning, etc.

- There is a smaller version of Anaconda called miniconda

- Can be downloaded from

    – www.anaconda.com

# Snapshot of the webpage

# Benefit of using Anaconda

- Support multiple virtual environments
- Excellent package manager named conda
- Conflicting packages can be easily avoided
  - You can install two versions of Numpy on the same computer, but in different virtual environments.

# Running Python

- On personal computer

  – Requires installation of python distribution

- On cloud computing service

  – Does not require software installation

  – Your codes are with you all the time

  – Usually free

# Cloud computing service

- Microsoft Azure

  - Available for CMU students for free

- Google Colaboration (colab in short)

  - Free

  - Requires Google's account (gmail)

  - Automatically links to Google drive

# Google Colaboration

- Visit https://colab.research.google.com
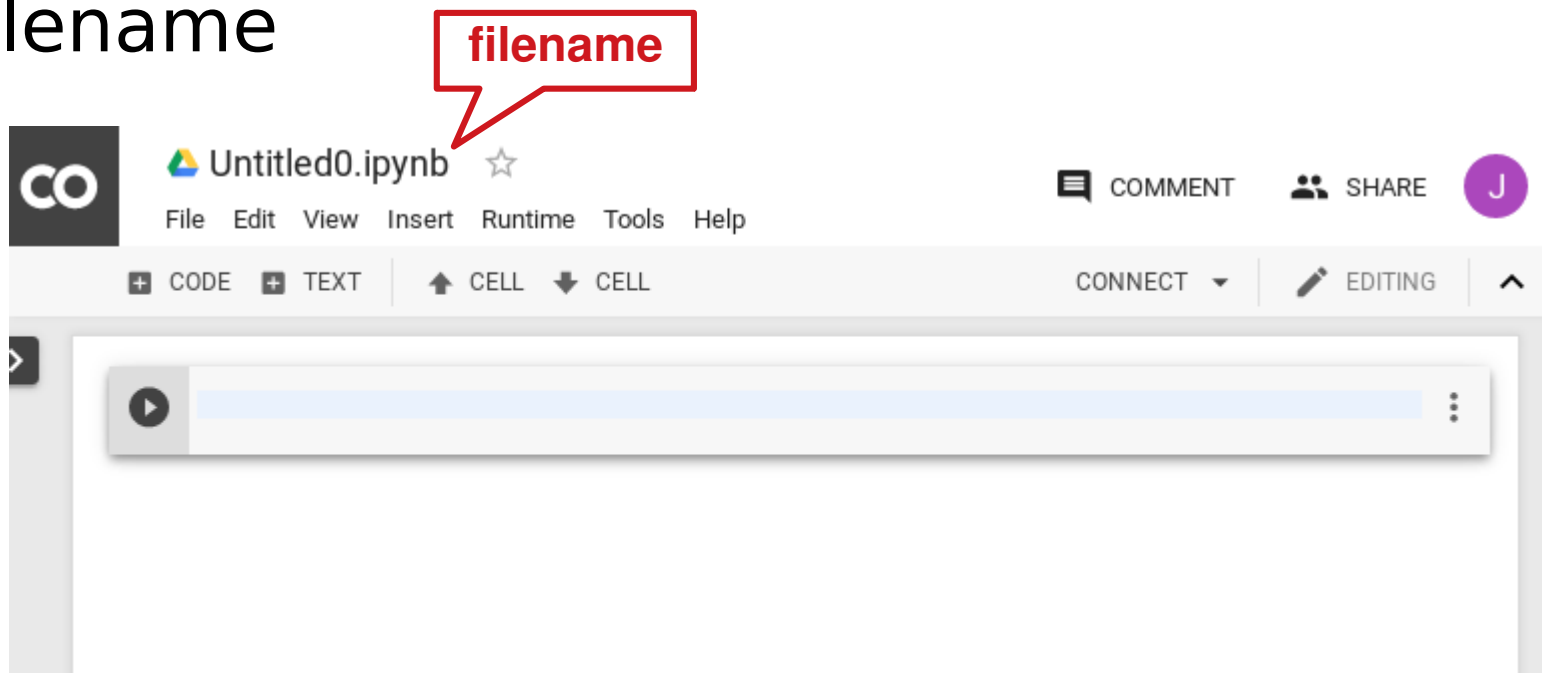- And sign in with your Google's credential

# After signing in

# **Alternatively**

- You can create new file by choosing file menu on the top-left corner

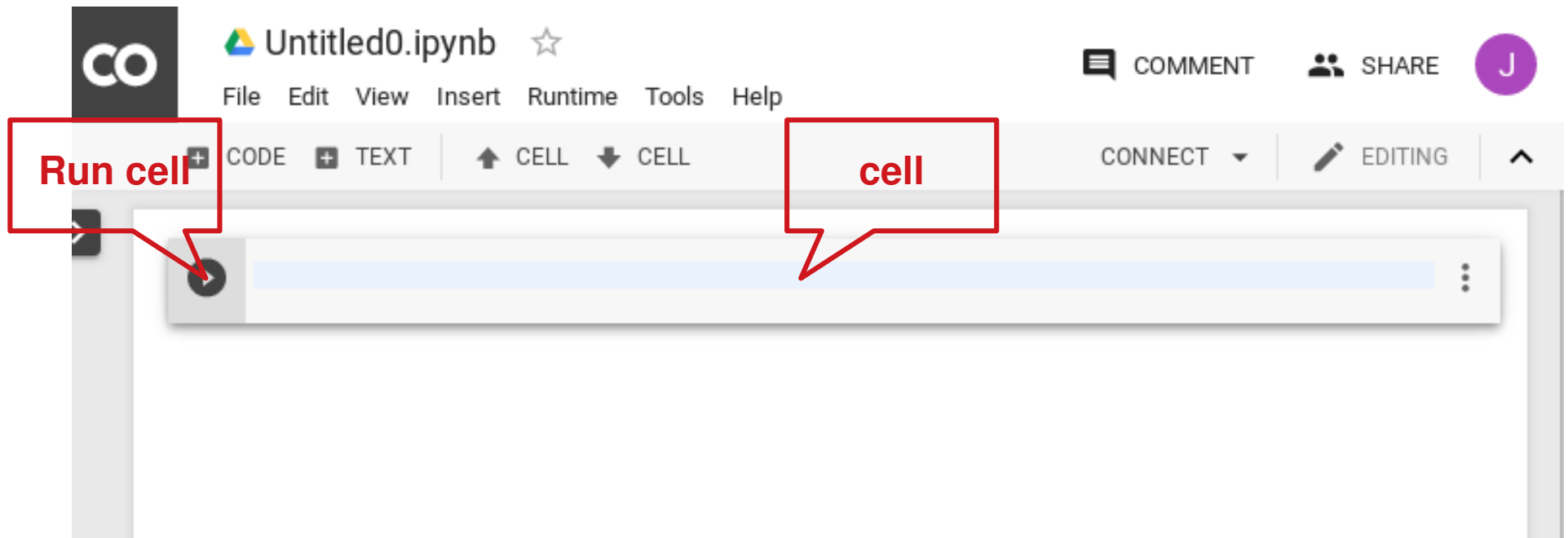  - File → New Python 3 notebook

# Your new Python file

- A Python sourcecode is called a notebook

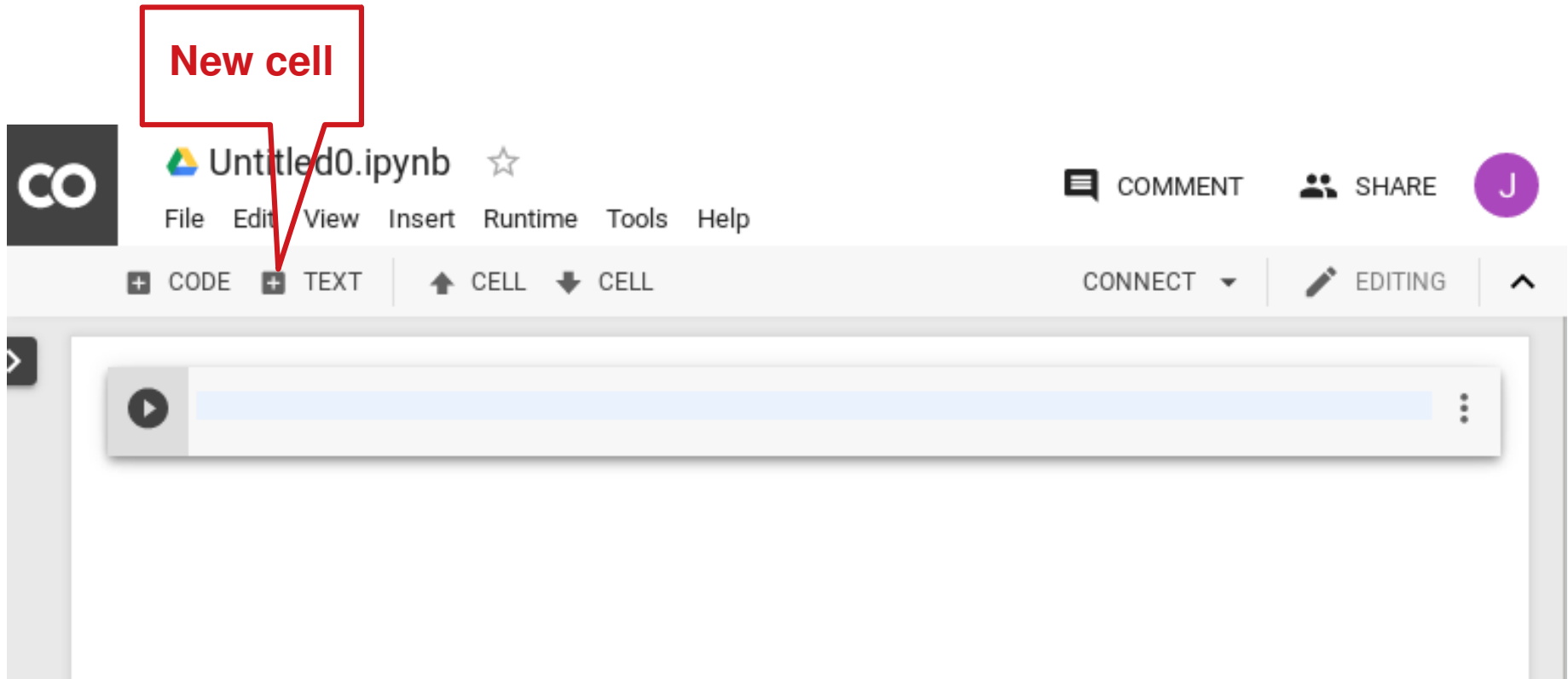- You can rename the file by double click filename

# What's in a notebook ?

# Adding new code cell



New cell

# Adding new text cell



**New cell**

Text cell is useful for adding context to the code

# Exercise 1

- Add one code cell and one text cell

# Deleting cell



Click this dots

# Running a cell

- Click the `play' button in front of cell
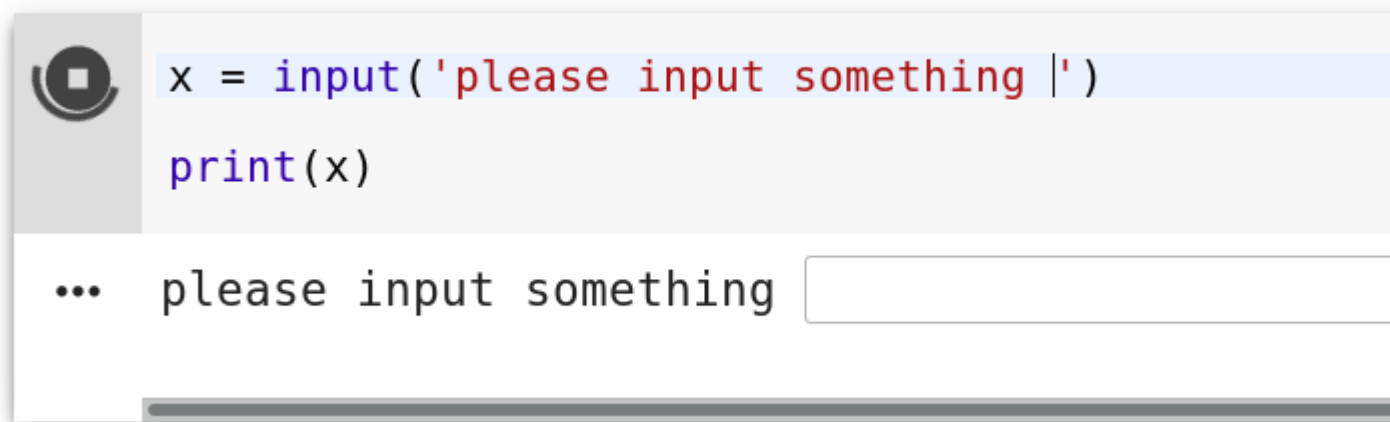- Or press `CTRL-ENTER'
- Let's try running the following code

```
[1]  x = 2
     y = 5

     print(x+y)

  ⤷  7
```

# **Getting input from user**

- You can get input from user in the same way you did with IDLE

```
x = input('please input something |')
print(x)
```

... please input something

- It waits for an input and will continue after you pressed ENTER

# Notes

- Once you've defined variables in a cell (and run it), the variables can be referenced in the subsequent cells.

- You can split your BIG code into multiple cells

- It will be easier to debug your code

# Example

```
[5]  x = 2

[6]  y = 3

▶  print(x+y)

⤷  5
```

# Colab and Google drive

- Files will be saved on your Google drive
- Click Colab icon to open the drive

**Go to Google Drive**

# Code snippets

- A snippet is a short code for doing some specific task

- Colab provides many useful snippets as examples

- They can be reused. (Problem solving using analogy and reduction)

# Let's try some snippet

# Useful websites [1]

- Data science competition / learning hub
- https://www.kaggle.com/

# Useful websites [2]

- Data sharing webset

- https://archive.ics.uci.edu/ml/index.php

# Useful Website [3]

- For asking programming related questions
- Or looking for solutions to problems similar to yours
- https://stackoverflow.com/

# The End

- Thank you and have fun !!

# Homework 0

- Read Overview of Colaboration