# Coding seminar

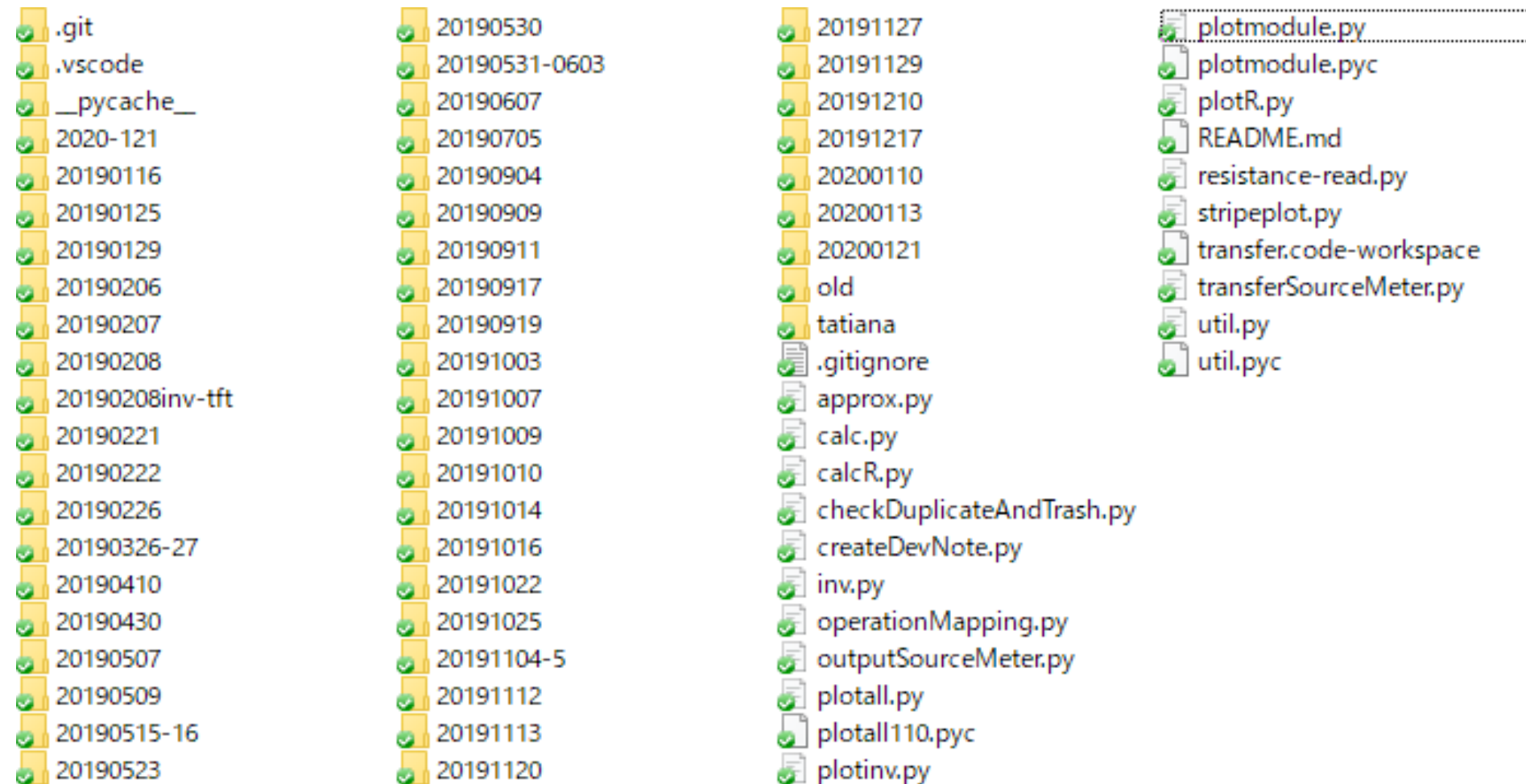## Lesson 4: UNIX commands and Statistics

Ikue Hirata, PhD

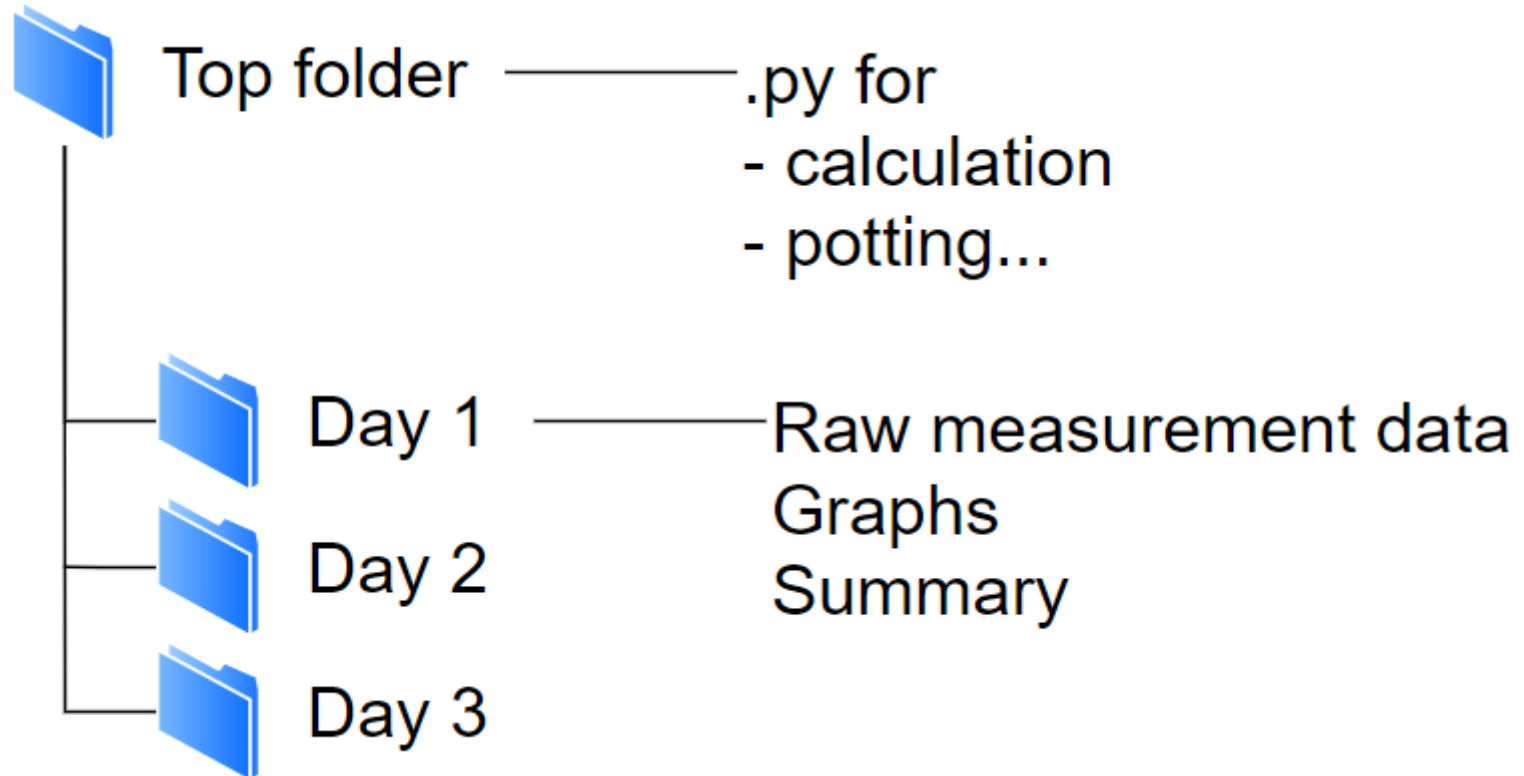# Contents

Folder structure and UNIX commands

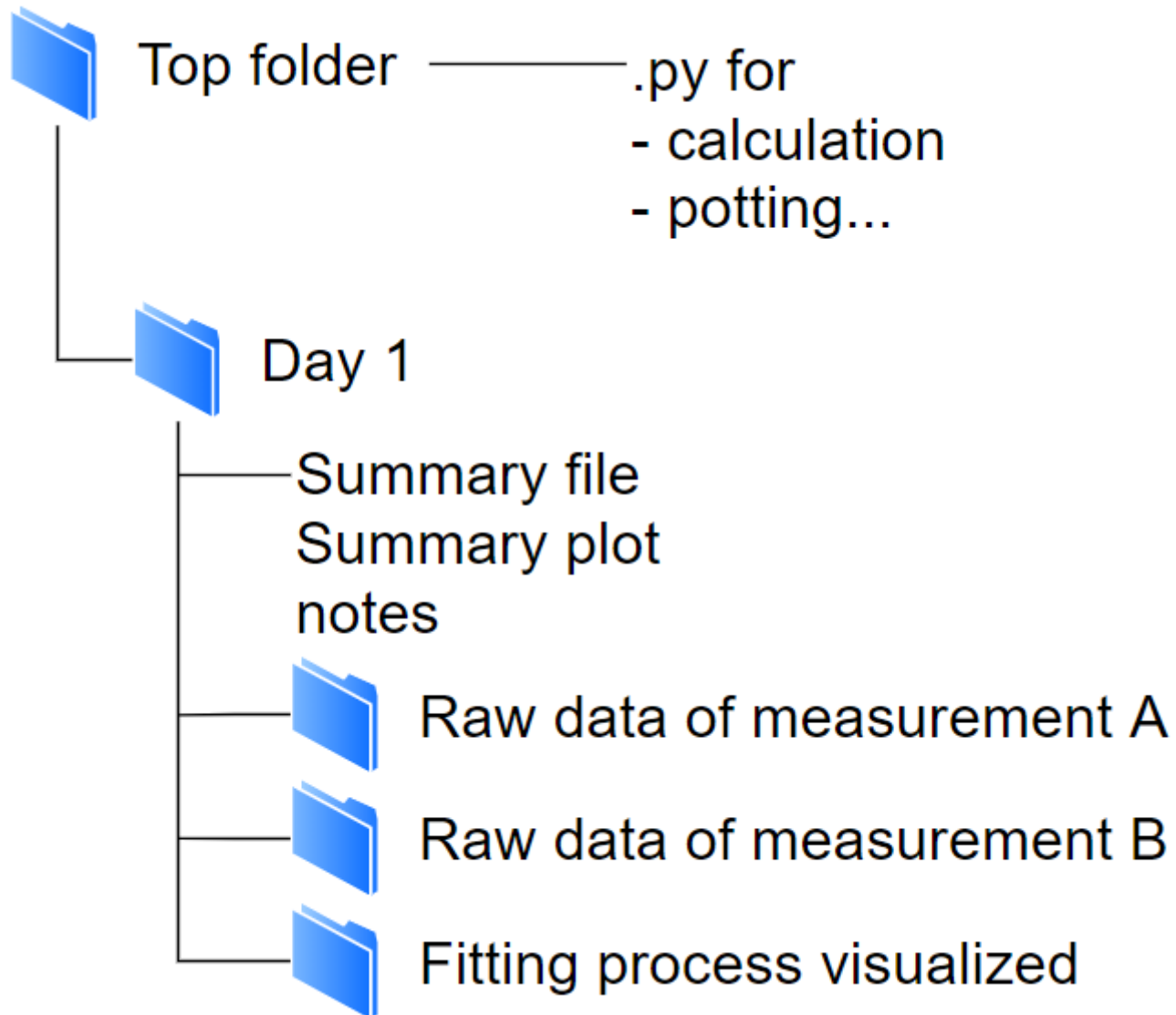Pandas, seaborn

# Getting more practical/systematic

| | | | |
|---|---|---|---|
| .git | 20190530 | 20191127 | plotmodule.py |
| .vscode | 20190531-0603 | 20191129 | plotmodule.pyc |
| __pycache__ | 20190607 | 20191210 | plotR.py |
| 2020-121 | 20190705 | 20191217 | README.md |
| 20190116 | 20190904 | 20200110 | resistance-read.py |
| 20190125 | 20190909 | 20200113 | stripeplot.py |
| 20190129 | 20190911 | 20200121 | transfer.code-workspace |
| 20190206 | 20190917 | old | transferSourceMeter.py |
| 20190207 | 20190919 | tatiana | util.py |
| 20190208 | 20191003 | .gitignore | util.pyc |
| 20190208inv-tft | 20191007 | approx.py | |
| 20190221 | 20191009 | calc.py | |
| 20190222 | 20191010 | calcR.py | |
| 20190226 | 20191014 | checkDuplicateAndTrash.py | |
| 20190326-27 | 20191016 | createDevNote.py | |
| 20190410 | 20191022 | inv.py | |
| 20190430 | 20191025 | operationMapping.py | |
| 20190507 | 20191104-5 | outputSourceMeter.py | |
| 20190509 | 20191112 | plotall.py | |
| 20190515-16 | 20191113 | plotall110.pyc | |
| 20190523 | 20191120 | plotinv.py | |

# Directory structure



Top folder —————— .py for
  - calculation
  - potting...

Day 1 —————— Raw measurement data
  Graphs
  Summary

Day 2

Day 3

# Directory structure

Top folder ——————— .py for
- calculation
- potting...

Day 1
Summary file
Summary plot
notes

Raw data of measurement A

Raw data of measurement B

Fitting process visualized

# Directory structure



Top folder ———— .py for
- calculation
- potting...

**Call codes in parent**

Day 1

**Summary files in current**

Summary file
Summary plot
notes

**Read files in children**

Raw data of measurement A

Raw data of measurement B

Fitting process visualized

# Launch Anaconda Prompt

# Anaconda Prompt



Anaconda Prompt (Anaconda3)

```
(base) C:¥Users¥Ikue>_
```

**Home directory**

# Move to other directory

# cd Move to other directory

Anaconda Prompt (Anaconda3)

```
(base) C:\Users\Ikue>cd Dropbox\PythonCourse

(base) C:\Users\Ikue\Dropbox\PythonCourse>cd Lesson4

(base) C:\Users\Ikue\Dropbox\PythonCourse\Lesson4>
```

# `dir` Show info on the directory



Anaconda Prompt (Anaconda3)

```
(base) C:¥Users¥Ikue¥Dropbox¥PythonCourse¥Lesson4>dir
 ドライブ C のボリューム ラベルは             です
 ボリューム シリアル番号は          です

 C:¥Users¥Ikue¥Dropbox¥PythonCourse¥Lesson4 のディレクトリ

2020/01/23  17:05    <DIR>          .
2020/01/23  17:05    <DIR>          ..
2020/01/23  17:06    <DIR>          dummyfiles1
2020/01/23  18:42               239 showfiles.py
               1 個のファイル               239 バイト
               3 個のディレクトリ  105,405,521,920 バイトの空き領域

(base) C:¥Users¥Ikue¥Dropbox¥PythonCourse¥Lesson4>_
```

# Terminal in VS Code is similar

# (optional) from Jupyter Notebook

**Only in Jupyter**

In [17]:

```
%pwd
```

'C:\\Users\\Ikue\\Dropbox\\PythonCourse\\Lesson4'

In [14]:

```
%cd Lesson4/
```

C:\Users\Ikue\Dropbox\PythonCourse\Lesson4

# Run .py file

**At this level**

Lesson4
- showfiles.py
- dummyfiles1
  - text1.txt
  - text2.txt
  - text3.txt

```
In [14]:    %cd Lesson4/
```

```
C:\Users\Ikue\Dropbox\PythonCourse\Lesson4
```

```
In [20]:    !python showfiles.py
```

```
Files in the current folder
['dummyfiles1', 'showfiles.py']
Files in the child folders
['dummyfiles1\\text1.txt', 'dummyfiles1\\te
xt2.txt', 'dummyfiles1\\text3.txt', 'dummyf
iles1\\text4.txt']
```

# Run .py file



```
In [21]:  %cd dummyfiles1/
          !python ../showfiles.py

          C:\Users\Ikue\Dropbox\PythonCourse\Lesson4
          \dummyfiles1
          Files in the current folder
          ['text1.txt', 'text2.txt', 'text3.txt', 'te
          xt4.txt']
          Files in the child folders
          []
```

**Parent ../**

**At this level**

# In practical case

```
ExperimentFolder\Today> python ../analyze.py
```
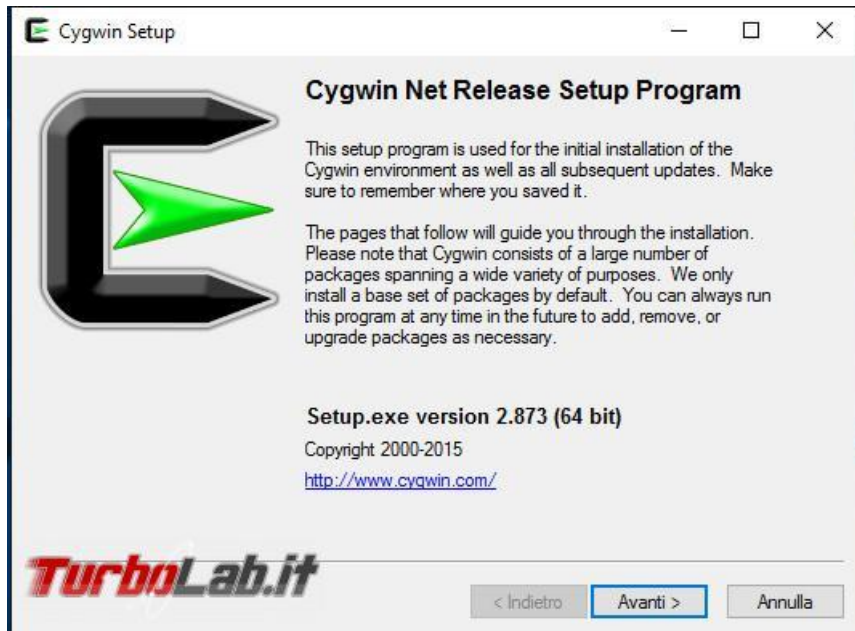
Experiment folder
- analyze.py
- Today
  - analysis result 1
  - analysis result 2
  - data
    - data files

```python
# analyze.py

l = glob.glob("data/*")
for files in l:
    some_calculations()
np.savetxt("result1.csv")
```

# UNIX commands

| Function | UNIX | Python |
|---|---|---|
| show current path | `pwd` | `os.getcwd()` |
| change directory | `cd` | `os.chdir()` |
| copy file(s) | `cp` | `shutil.copyfile()` |
| mv file(s) | `mv` | `os.rename()`<br>`shutil.move()` |
| make directory(ies) | `mkdir` | `os.mkdir()` |
| search words | `grep` | `glob.glob() + find()` |
| replace words | `sed`<br>`awk` | `re.sub() …?` |

# Run UNIX commands on Windows

## Option 1: Cygwin (and other shells)



## Option 2: WSL



## More in Appendix

# SciPy package family



https://www.scipy.org/

# pandas: database handling

**Table Of Contents**

**Search**

[ ] Go

Enter search terms or a module, class or function name.

## pandas: powerful Python data analysis toolkit

**Date**: Nov 09, 2019 **Version**: 0.25.3

**Download documentation**: PDF Version | Zipped HTML

**Useful links**: Binary Installers | Source Repository | Issues & Ideas | Q&A Support | Mailing List

**pandas** is an open source, BSD-licensed library providing high-performance, easy-to-use data structures and data analysis tools for the Python programming language.

See the Package overview for more detail about what's in the library.

- What's new in 0.25.3 (October 31, 2019)
- Installation
- Getting started
  - Package overview
  - 10 minutes to pandas
  - Essential basic functionality
  - Intro to data structures
  - Comparison with other tools
  - Tutorials
- User Guide
  - IO tools (text, CSV, HDF5, …)
  - Indexing and selecting data
  - MultiIndex / advanced indexing
  - Merge, join, and concatenate
  - Reshaping and pivot tables
  - Working with text data
  - Working with missing data
  - Categorical data
  - Nullable integer data type

https://pandas.pydata.org/pandas-docs/stable/index.html

# Benefit of pandas

**numpy**: for numerical array, like x-y

loadtxt() - no missing data

genfromtext() - with missing data
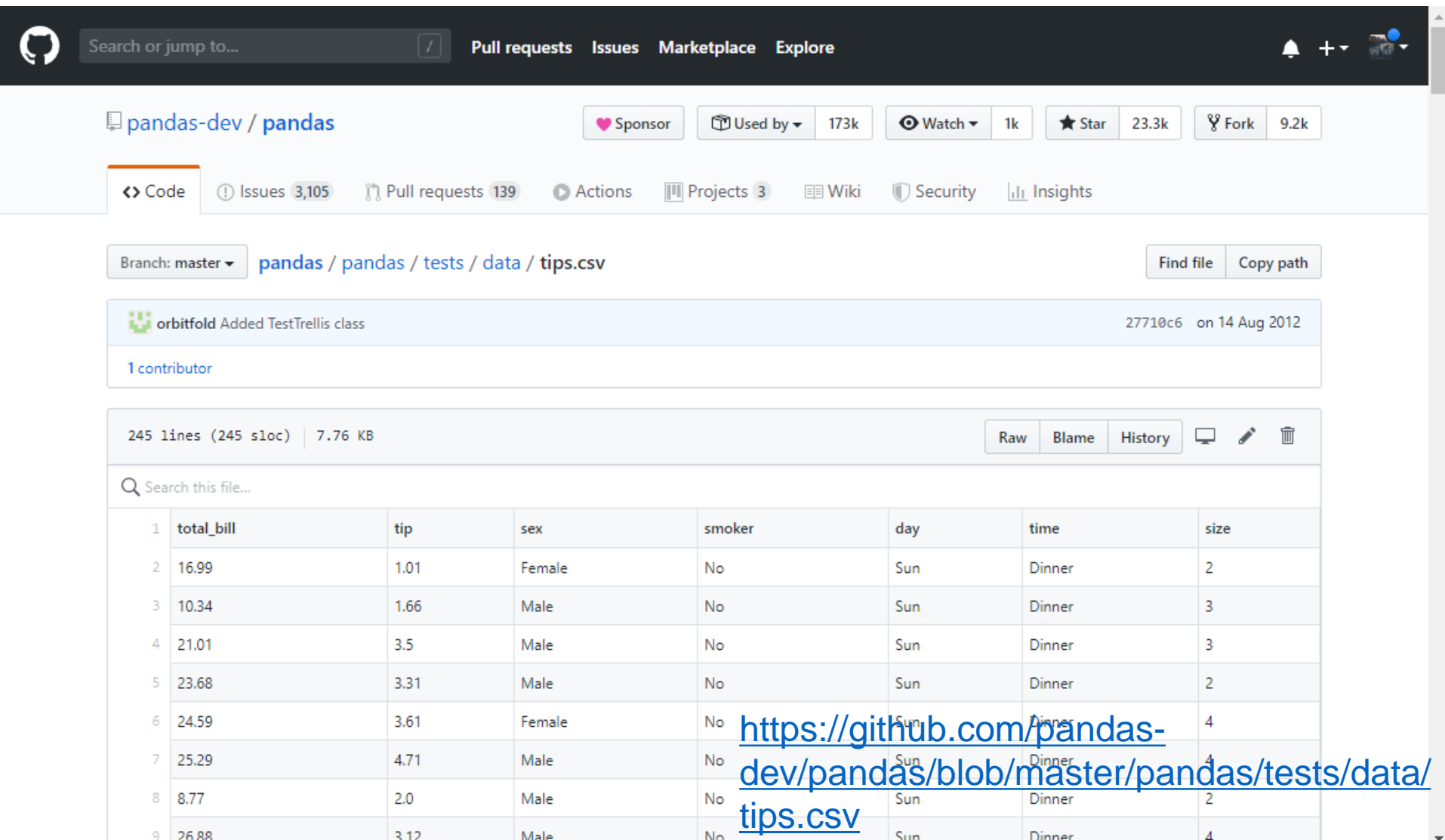
column name -> tuple (cannot access by index)


**pandas**: for data set, including `str`

load_csv() accessible by index and/or label

generic column name

stronger for statistics

# Sample data



https://github.com/pandas-dev/pandas/blob/master/pandas/tests/data/tips.csv

# By Numpy

```
In [31]:  import numpy as np
          dataar = np.genfromtxt("Lesson4/tips.csv",
          delimiter=",", skip_header=1)
          print(dataar)
```

```
[[16.99  1.01   nan ...   nan   nan   2.  ]
 [10.34  1.66   nan ...   nan   nan   3.  ]
 [21.01  3.5    nan ...   nan   nan   3.  ]
 ...
 [22.67  2.     nan ...   nan   nan   2.  ]
 [17.82  1.75   nan ...   nan   nan   2.  ]
 [18.78  3.     nan ...   nan   nan   2.  ]]
```

**Text data lost**
**(needs special treatment)**

# By pandas

```python
import pandas as pd
datadf = pd.read_csv("Lesson4/tips.csv")
print(datadf)
```

```
     total_bill   tip     sex smoker   day    time  size
0         16.99  1.01  Female     No   Sun  Dinner     2
1         10.34  1.66    Male     No   Sun  Dinner     3
2         21.01  3.50    Male     No   Sun  Dinner     3
3         23.68  3.31    Male     No   Sun  Dinner     2
4         24.59  3.61  Female     No   Sun  Dinner     4
..          ...   ...     ...    ...   ...     ...   ...
239       29.03  5.92    Male     No   Sat  Dinner     3
240       27.18  2.00  Female    Yes   Sat  Dinner     2
241       22.67  2.00    Male    Yes   Sat  Dinner     2
242       17.82  1.75    Male     No   Sat  Dinner     2
243       18.78  3.00  Female     No  Thur  Dinner     2

[244 rows x 7 columns]
```

# Partial data by conditions

In [40]:

```python
print(datadf[datadf["sex"]=="Female"])
```

```
     total_bill   tip     sex smoker   day    time  size
0         16.99  1.01  Female     No   Sun  Dinner     2
4         24.59  3.61  Female     No   Sun  Dinner     4
11        35.26  5.00  Female     No   Sun  Dinner     4
14        14.83  3.02  Female     No   Sun  Dinner     2
16        10.33  1.67  Female     No   Sun  Dinner     3
..          ...   ...     ...    ...   ...     ...   ...
226       10.09  2.00  Female    Yes   Fri   Lunch     2
229       22.12  2.88  Female    Yes   Sat  Dinner     2
238       35.83  4.67  Female     No   Sat  Dinner     3
240       27.18  2.00  Female    Yes   Sat  Dinner     2
243       18.78  3.00  Female     No  Thur  Dinner     2

[87 rows x 7 columns]
```

# Some statistics

```
In [34]:   datadf.describe()
```

|        | total_bill | tip        | size       |
|--------|-----------|------------|------------|
| count  | 244.000000 | 244.000000 | 244.000000 |
| mean   | 19.785943  | 2.998279   | 2.569672   |
| std    | 8.902412   | 1.383638   | 0.951100   |
| min    | 3.070000   | 1.000000   | 1.000000   |
| 25%    | 13.347500  | 2.000000   | 2.000000   |
| 50%    | 17.795000  | 2.900000   | 2.000000   |
| 75%    | 24.127500  | 3.562500   | 3.000000   |

# Plotting by seaborn

In [39]:

```python
%matplotlib inline
import seaborn as sns
sns.boxplot(x="day", y="total_bill",
            hue="smoker", data=datadf)
```

<matplotlib.axes._subplots.AxesSubplot at 0x261c4ab3d08>

# More complex plot

In [53]:

```python
sns.catplot(x="day", y="total_bill", col="sex",
row="smoker", hue="time", data=datadf,
kind="violin")
```
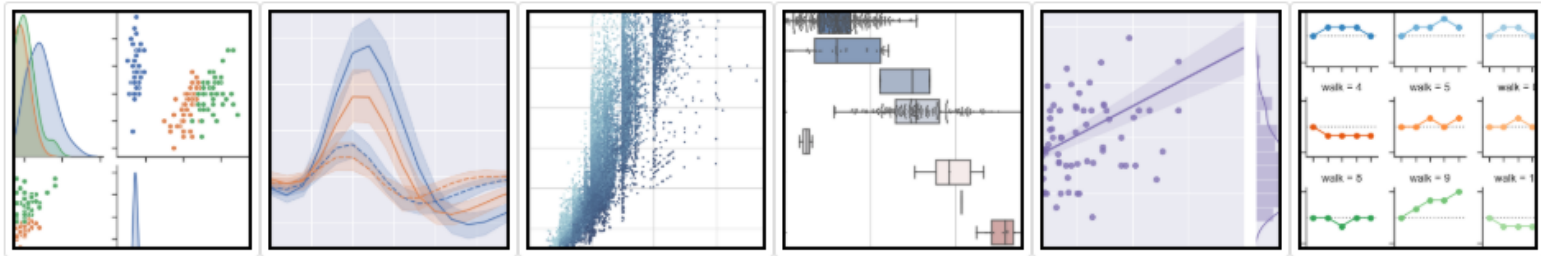
# More about seaborn

# Exercise 1

## Lesson 4 - Exercise 1

Here is a dataset of iris flowers iris.csv and some iris flowers which species names are unknown. Identify the species name of each flowers below.

| Name | SepalLength | SepalWidth | PetalLength | PetalWidth |
|------|-------------|------------|-------------|------------|
| iris1 | 6.5 | 3 | 5 | 1.2 |
| iris2 | 6.5 | 3 | 5 | 2.4 |
| iris3 | 6 | 3 | 2 | 0.5 |

Hint:

1. Import `iris.csv` to a `pandas.DataFrame`.
2. Add the unknown flowers data to the `DataFrame`.
3. Plot each parameters by `seaborn.pairplot` to show the characteristics of the species.

# Exercise 2 (difficult!)

## Lesson 4 - Exercise 2

Simple substitution cipher is a method of encrypting. Each plaintext alphabet is replaced by other alphabet ([Wikipedia](#)). `Lesson4-files/cipher.txt` (shown below) is an English text encrpyted by simple substitution cipher. Decrypt this and obtain the encryption table.

```
NRMZFYKRKRMRSOUUCNISRACWDRRCRWDOMACWDRDOMYARHOFERGWARUIMKEWIRJR
SCAIYCWIERRCIFWIBRTWIHINRGMWIFHDTGWDCIFEJRSCIMRERRETDOWCIFOXXRO
SWIARBRTSROUUCLMITTDOWCIFOSROMYWDIERBRTYOSRMIWIXXIERWDRNERUJREW
IWDRIXGMGIMIBWDRNOMCTDIDOJRWDRNOZREWCIBWDREWOWRWIYRBRMYWDRNOMYG
MWDROHWGIMEIBOUUNRMOMYREXRHGOUUCIBXSGMHRETDGHDGWGEMIWXSFYRMWWIH
DOUURMKRIMRZFYKREACWDRSREFUWFISWDOWSROEIMURWOXSGMHRDOJRWDRHSRYG
WIBHIMPFRSGMKOMYDIUYGMKDGEEWOWRWDRNROMETGUUOUTOCEARHIMEGYRSRYDI
MREWOMYDRTGUUARXSOGERYACRJRSCAIYCARHOFERWDRJFUKOSOSROUTOCEWOLRM
ACTDOWOWDGMKERRNEWIAROMYACTDOWHINREIBGWOMYGMWDRTISUYWDRSROSRIMU
CWDRJFUKOSBISWDRBRTBGMYOXUOHRWDRSRIMUCTDRMWDRNOMCDOJRMIKSIFMYWI
SREWIMOMRXSGMHRIBWDRXSRERMWWGNRTDINGWGEMIWTRUUWIMONRMRJRSXSROHD
REOMCWDGMKRUERAFWXROHROMYKIIYBOGWDOMYWIAIWDDRGENIEWDIEWGUROMYRG
WDRSGBDRDOYLRXWGWTIFUYDOJRYRXSGJRYDGNIBSRXFWOWGIMOMYLGMKYINNOMC
OWGNRWSGWWRMACNGHIUIMOHDGOJRUUGWSOMEUOWRYACWKMOSSGIWWTDRPSGMHRD
WWXETTTKFWRMARSKISK
```

1. Frequency analysis is a powerful method to help break a simple substitution cipher. In this method, each alphabet is replaced according to the frequency of the appearance in the text. It is known that in any given stretch of written language, certain letters and combinations of letters occur with varying frequencies ([Wikipedia](#)).
   A. Count the number of the appearance of each alphabet in a plaintext file `Lesson4-files/kant.txt` and plot it in a histogram.
   B. Count the number of the appearance of the combination of any given two alphabets and plot top 20 combinations in a histogram.
   C. Count the number of the appearance of the combination of any given three alphabets and plot top 20 combinations in a histogram.
2. Using the results above, decipher the text and obtain the encryption table.

# Study by yourself

Pandas official document

https://pandas.pydata.org/pandas-docs/stable/index.html

Pandas unofficial tutorial

https://www.learndatasci.com/tutorials/python-pandas-tutorial-complete-introduction-for-beginners/
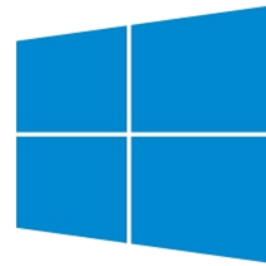
seaborn official document

https://seaborn.pydata.org/

# Appendix for Windows users

Terminals and shells

# Why commands are different?

| DOS commands | Linux command |
|---|---|
| <command> /? | man <command> or command --help |
| cd | cd |
| chdir | pwd |
| cls | clear |
| copy | cp |
| date | date |
| del | rm |
| dir | ls |
| echo | echo |
| edit | vim (or other editor) |
| exit | exit |
| fc | diff |
| find | grep |
| format | mke2fs or mformat |
| mem | free |
| mkdir | mkdir |
| more | more or even less |
| move | mv |
| ren | mv |
| time | date |

# Because of kernel



Applications

Kernel

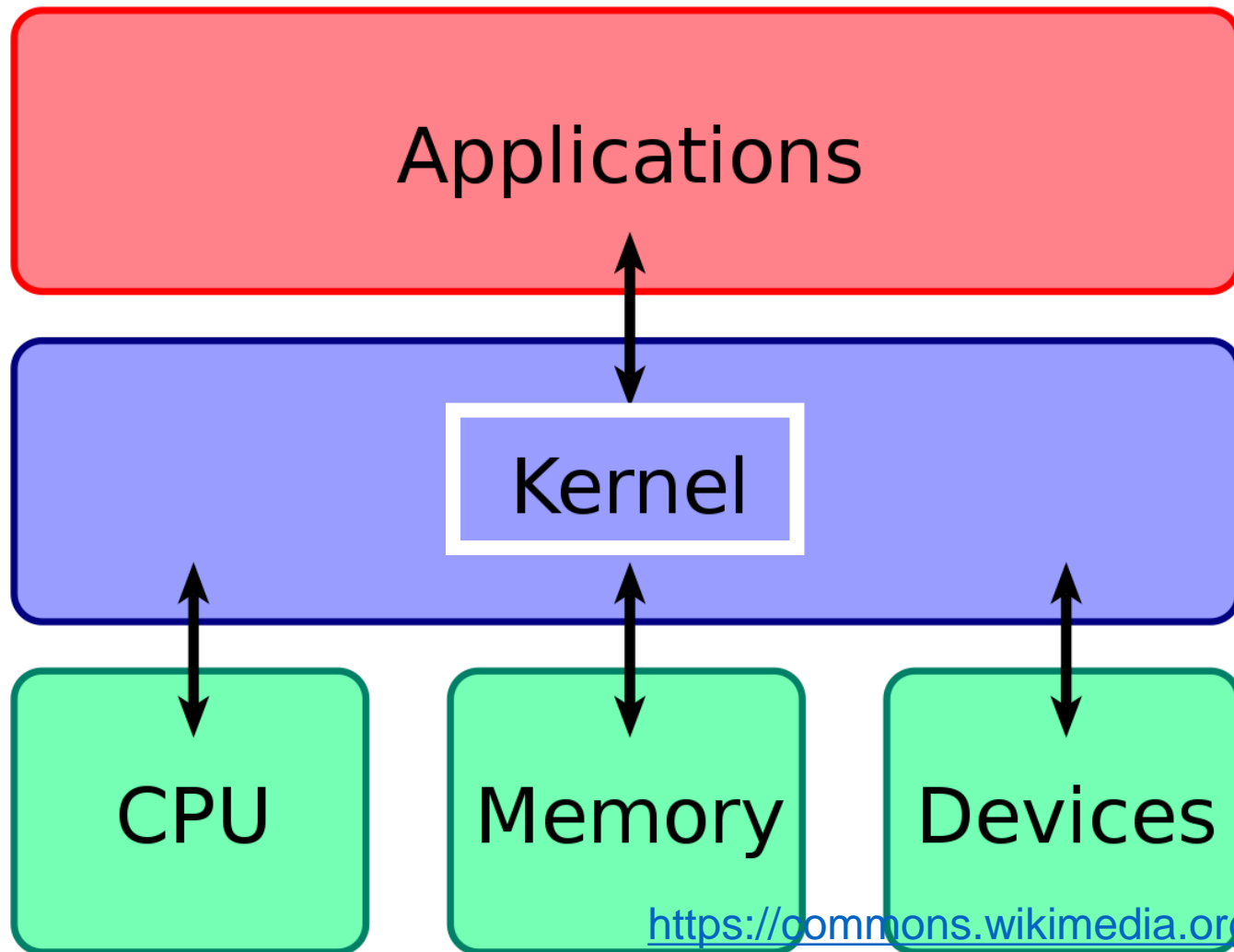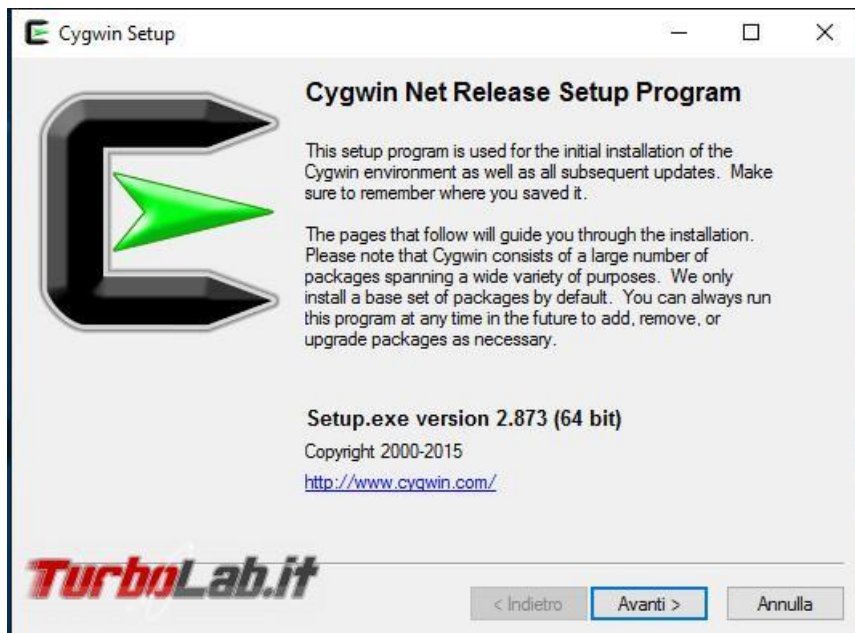CPU    Memory    Devices
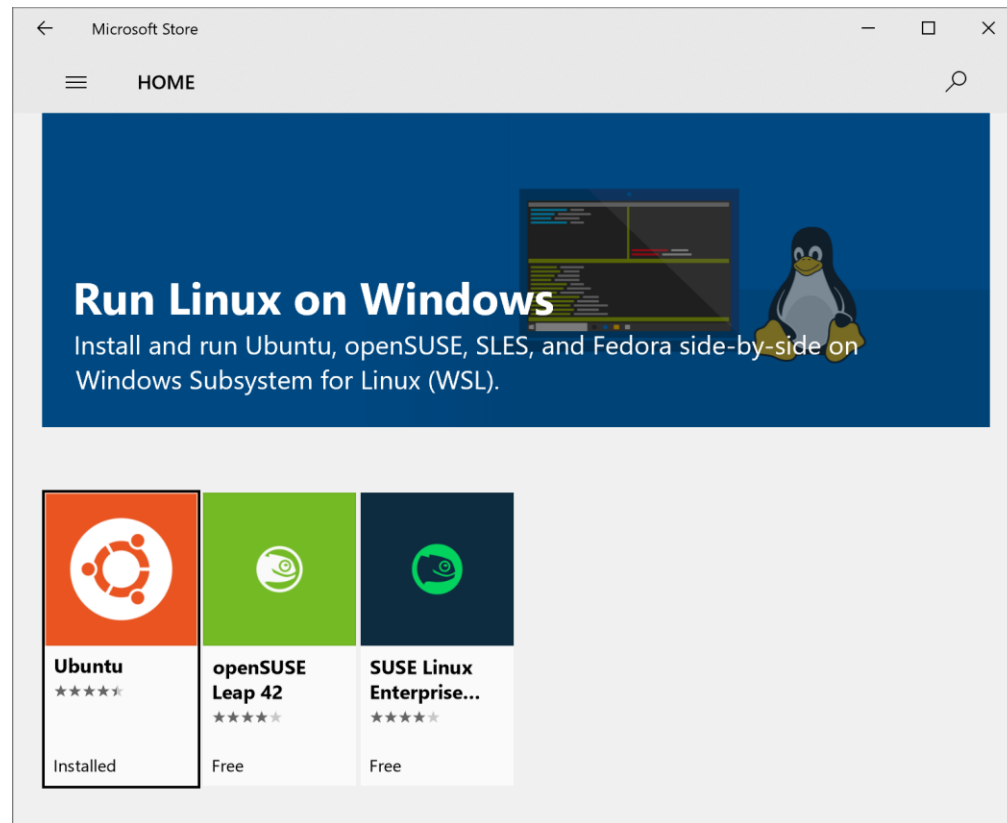
# Install UNIX-like environment

## Option 1: Cygwin (and other shells)



## Option 2: WSL

# Keywords to learn further

About CUI:    -   stdin/stdout
              -   redirect
              -   pipe
              -   alias
              -   PATH

About OS:     -   kernel
              -   terminal
              -   shell
              -   Linux/UNIX

# Have fun! :)