

Data Science

Week 2

Getting Start with Data Manipulation

Outline

1. Data Collection
2. Data Loading and File Formats

Data from Public (Open) Dataset

- If you've ever worked on a personal data science project, you've probably spent a lot of time browsing the [internet](#) looking for interesting datasets to analyze.
- A [dataset](#), or data set, is simply a collection of data.
- The simplest and most common format for datasets you'll find online is a spreadsheet or [CSV format](#) — a single file organized as a table of rows and columns.
- But some datasets will be stored in [other formats](#), and they don't have to be just one file. Sometimes a dataset may be a zip file or folder containing multiple data tables with related data.
- **When you use the open dataset, you must follow the requirements of using the dataset (copyright, citation, acknowledge, and so on).**

Example of Open Dataset

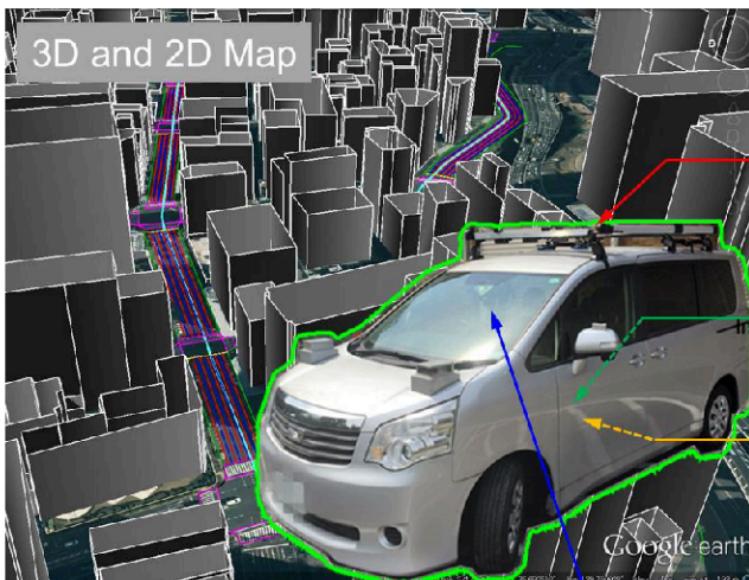
- Data science community:
 - E.g. Kaggle, <https://www.kdnuggets.com/datasets/index.html>
- Dataset platform
 - E.g. Google, <https://cloud.google.com/bigquery/public-data/>
- Dataset of Governmental institution
 - E.g. NASA, earth data: <https://earthdata.nasa.gov/>, data related space: <https://pds.nasa.gov/datasearch/data-search/>
- Specific website for benchmark
 - COCO dataset: <https://cocodataset.org/#home>
- And so on...

Building a Dataset from Scratch

- In my experience, we could not find the open dataset for most of projects, and we need to collect data by ourselves.
- Data collection is the process of gathering quantitative and qualitative information on specific variables. Good data collection requires a clear process to ensure the data you collect is clean, consistent, and reliable.
- Establishing that process involves taking stock of your objectives (limitations if any), identifying your data requirements, and finally organizing a data collection plan that synthesizes the most important aspects of your project.
- The following 5 examples (3 projects conducted by the instructor of this class) demonstrate how to perform the data collection for the projects with specific objectives.

Example 1: Multi-sensor based vehicle self-localization project

- **Objective:** using GNSS receiver, IMU (Inertial Measurement Unit) sensor, Speedometer, Camera, Map information for vehicle positioning



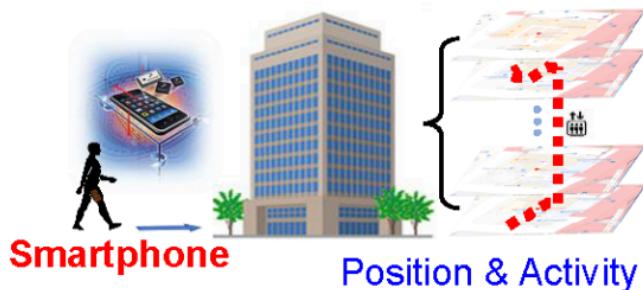
Example 1: Multi-sensor based vehicle self-localization project

- **Limitation:** this project needs the highly accurate map information for vehicle positioning, an accurate map at a specific place of Tokyo has been provided.
- This project has to collect other sensor data when vehicle drives at the place with map information.
- There is NO open dataset can be used for this research.

Example 1: Multi-sensor based vehicle self-localization project

- Data requirements
 - Synchronization of multi-sensor data
 - Use one Computer to control all sensors, and the local time of the computer is used for synchronization
 - Resolution of each data
 - It is necessary to consider the communication bandwidth between computer and each sensor
 - If it is possible, choose the highest resolution for each sensor data in the collection, and adjust the resolution in the later steps of the project
- Do not forget to collect ground truth data
 - Evaluation needs ground truth (compare the estimated vehicle position and the ground truth position in the evaluation).
 - In this project, the vehicle lane-position (ground truth) is recorded by camera.

Example 2: Smartphone Lifelog project



Smartphone-based Lifelog automatically annotates the users' daily experience (position & activity) from multisensory streams on smartphones.

- **Objective:** using multisensory streams (GPS receiver, accelerometer, gyroscope, barometer, compass sensor, WiFi receiver) to calculate the position and activity of user.
- After survey, we found that there is NO open dataset can be used for this research.

Example 2: Smartphone Lifelog project

- Data requirements
 - Synchronization of multi-sensor data
 - Smartphone controls all sensors, and the local time of the smartphone is used for synchronization
 - Resolution of each data
 - If it is possible, choose the highest resolution for each sensor data in the collection, and adjust the resolution in the later steps of the project
 - **Do not forget to collect ground truth data**
 - Evaluation for positioning needs the position ground truth.
 - Evaluation for activity recognition needs to record the activity ground truth.
 - Learning process (activity recognition/classification task) need the activity ground truth.
 - **We need to record the position and activity of user in the experiment (e.g. record synchronized user's voice-over which tells the ground truth).**

Example 3: Stock performance prediction project

- **Objective:** Using historical data of stocks to predict the future of stock performance
- After survey, we found that there is NO open dataset (completed data) can be used for this research.
- We found that **Thomson Reuters database (not free)** can provide the data, and it is possible to use Datastream to access the financial data.
- We prepare a simple program to acquire the data automatically by **batch processing**.

15.89	+5.34%	543.23	300,000
45.34	-7.89%	254.23	320,000
17.34	+5.97%	921.56	430,000
34.89	+2.13%	564.23	120,000
16.45	+6.43%	765.80	900,000
23.67	-11.6%	120.84	300,000
34.64	+23.1%	893.23	120,000
43.69	+5.56%	128.98	320,000
12.78	-3.67%	432.12	750,000
5.44	+11.3%	765.23	150,000
		432.24	120,000
		200,000	

Example 4: Data collection from internet-Wikipedia

- Wikipedia is one of modern humanity's most impressive creations. In just a few years, anonymous contributors working for free could create the greatest source of online knowledge the world.
- Not only is Wikipedia the best place to get information for writing your college papers, but it's also an extremely rich source of data that can fuel numerous data science projects from natural language processing to supervised machine learning.

Example 4: Data collection from internet-Wikipedia

- Using Wikipedia API and obtain the data from Wikipedia by python.

```
install wikipedia

import wikipedia

wikipedia.set_lang("en")
page = wikipedia.page('Ritsumeikan', auto_suggest=False)
print(page.content)
```

There are many useful functions `wikipedia.*`, please check the instruction of wikipedia lib when you use it.

Example 5: Data collection from internet-Website

- “Web scrapers” automatically collect information and data that's usually only accessible by visiting a website in a browser.
- By doing this autonomously, web scraping scripts open up a world of possibilities in data mining, data analysis, statistical analysis, and much more.

The copyright issue is existing when you use Web scrapers, please do NOT violate the **copyright laws** when you use Web scrapers

Example 5: Data collection from internet-Website

Collect data from website using python

```
install beautifulSoup4
```

```
from urllib import request  
  
from bs4 import BeautifulSoup  
  
# link of Yanlei Gu's Google Scholar  
response = request.urlopen('https://scholar.google.co.jp/citations?user=2vLFO00AAAAJ&hl=en')  
contents = BeautifulSoup(response, "html.parser")  
print (contents)
```

Outline

1. Data Collection
2. Data Loading and File Formats

Getting Started with pandas

- Pandas adopts many coding idioms from NumPy, the biggest difference is that
 - Pandas is designed for working with tabular or heterogeneous (different types of) data.
 - NumPy, by contrast, is best suited for working with homogeneous numerical array data.

```
import pandas as pd
```

Data Loading function in pandas

Function	Description
<code>read_csv</code>	Load delimited data from a file, URL, or file-like object; use comma as default delimiter
<code>read_table</code>	Load delimited data from a file, URL, or file-like object; use tab (' \t ') as default delimiter
<code>read_fwf</code>	Read data in fixed-width column format (i.e., no delimiters)
<code>read_clipboard</code>	Version of <code>read_table</code> that reads data from the clipboard; useful for converting tables from web pages
<code>read_excel</code>	Read tabular data from an Excel XLS or XLSX file
<code>read_hdf</code>	Read HDF5 files written by pandas
<code>read_html</code>	Read all tables found in the given HTML document
<code>read_json</code>	Read data from a JSON (JavaScript Object Notation) string representation
<code>read_msgpack</code>	Read pandas data encoded using the MessagePack binary format
<code>read_pickle</code>	Read an arbitrary object stored in Python pickle format
<code>read_sas</code>	Read a SAS dataset stored in one of the SAS system's custom storage formats
<code>read_sql</code>	Read the results of a SQL query (using SQLAlchemy) as a pandas DataFrame
<code>read_stata</code>	Read a dataset from Stata file format
<code>read_feather</code>	Read the Feather binary file format

Comma-separated (CSV) text file

- 1) Use `read_csv` to read it into a DataFrame
- 2) A file will not always have a header row, we should use `header=None` to avoid losing the first row in this case
- 3) Give the column name (header) manually
- 4) Use one column to be the index of the returned DataFrame

Files used for programming

ex1.csv

	A	B	C	D	E
1	a	b	c	d	message
2	1	2	3	4	hello
3	5	6	7	8	world
4	9	10	11	12	foo

ex2.csv

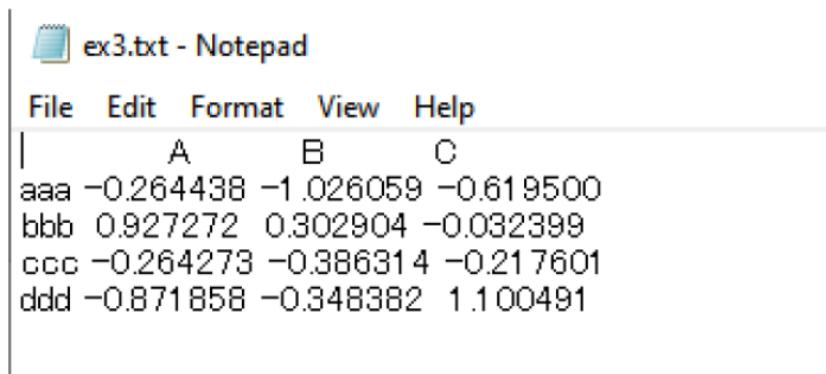
	A	B	C	D	E
1	1	2	3	4	hello
2	5	6	7	8	world
3	9	10	11	12	foo

Check the code by yourself

```
# (1) read standard CSV
# Code Example 1
print ("-----")
df = pd.read_csv('ch06/ex1.csv')
print (df)
print ("-----")
# (2 )A file will not always have a header row. Consider this file:
df2 = pd.read_csv('ch06/ex2.csv')
print (df2)
print ("-----")
# should use header=None to avoid lossing the first row
df3 = pd.read_csv('ch06/ex2.csv', header=None)
print (df3)
print ("-----")
# (3) give the column name manually
df4 = pd.read_csv('ch06/ex2.csv', names=['a', 'b', 'c', 'd', 'message'])
print (df4)
print ("-----")
# (4) use message column to be the index of the returned DataFrame
names = ['a', 'b', 'c', 'd', 'message']
df5 = pd.read_csv('ch06/ex2.csv', names=names, index_col='message')
print (df5)
print ("-----")
```

Whitespace-separated text file

- If the data are separated by the whitespace, we can pass a regular expression as a delimiter for `read_csv`. This can be expressed by the regular expression `\s+`:



The screenshot shows a Notepad window titled "ex3.txt - Notepad". The menu bar includes File, Edit, Format, View, and Help. The content of the file is as follows:

	A	B	C
aaa	-0.264438	-1.026059	-0.619500
bbb	0.927272	0.302904	-0.032399
ccc	-0.264273	-0.386314	-0.217601
ddd	-0.871858	-0.348382	1.100491

```
# read other type Text
# Code Example 2
#In some cases, a table might not have a fixed delimiter, using
whitespace or some
#other pattern to separate fields. Consider a text file that looks
like this:
list(open('ch06/ex3.txt'))
result = pd.read_csv('ch06/ex3.txt')
print ("-----")
print (result)
print ("-----")
print (result.iloc[0,[0]])

result1 = pd.read_csv('ch06/ex3.txt', sep='\s+')
print ("-----")
print (result1)
print ("-----")
print (result1.iloc[0,[0]])
print ("-----")
print (result1.loc[['aaa','A']])
print ("-----")
```

File with missing data

ex5.csv

	A	B	C	D	E	F
1	somethin\ a	b	c	d	4	message
2	one	1	2	3	all	
3	two	5	6		8	world
4	three	9	10	11	12	foo

```
# read a file with missing data
# Code Example 3
result = pd.read_csv('ch06/ex5.csv')
print ("-----")
print (result)
print ("-----")
```

JSON Data

- Short for JavaScript Object Notation
- JSON has become one of the standard formats for sending data by HTTP request between web browsers and other applications.
- It is a much more free-form data format than a tabular text form like CSV.

Load JSON Data

JSON_example.json

```
1 [ {"a": 1, "b": 2, "c": 3},  
2  {"a": 4, "b": 5, "c": 6},  
3   {"a": 7, "b": 8, "c": 9}]
```

- The `pandas.read_json` can automatically convert JSON datasets in specific arrangements into a Series or DataFrame.
- The default options for `pandas.read_json` assume that each object in the JSON array is a row in the table

Check the code by yourself

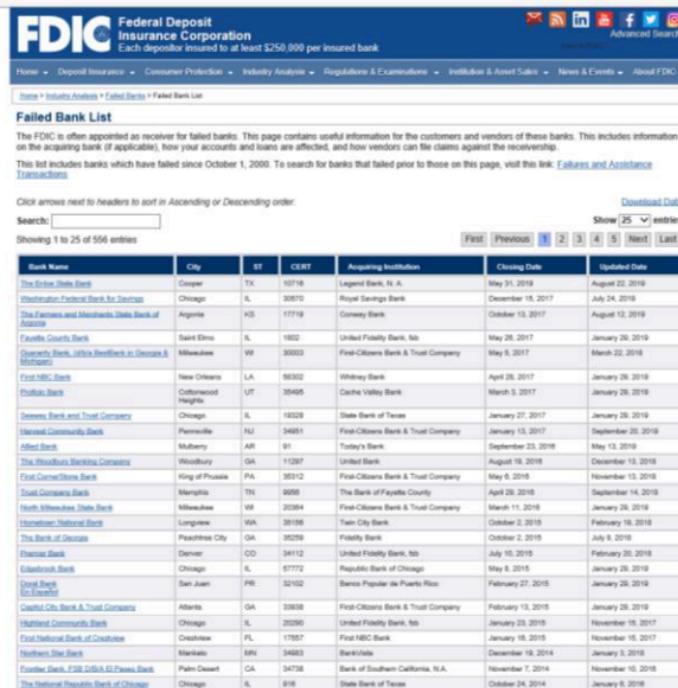
```
# read a json file  
data = pd.read_json('ch06/JSON_example.json')  
print ("-----")  
print (data)  
print ("-----")  
print (data.to_json())  
print ("-----") open
```

XML and HTML files

Extensible Markup Language and
Hypertext Markup Language

- Python has many libraries for reading and writing data in HTML and XML formats.
- Examples include lxml, BeautifulSoup, and html5lib.
- pandas has a built-in function, `read_html`, which uses libraries like lxml and BeautifulSoup to automatically parse tables out of HTML files as DataFrame objects.
- First, you must install some additional libraries used by `read_html`

Downloaded an HTML file from the United States FDIC government agency showing bank failures.



The screenshot shows the FDIC website's "Failed Bank List" page. At the top, there's a navigation bar with links like Home, Deposit Insurance, Consumer Protection, Industry Analysis, Regulations & Examinations, Institutions & Asset Sales, News & Events, and About FDIC. Below the navigation is a search bar with options for Advanced Search and a blue "Search" button. The main content area has a header "Failed Bank List" with a sub-instruction: "The FDIC is often appointed as receiver for failed banks. This page contains useful information for the customers and vendors of these banks. This includes information on the acquiring bank (if applicable), how your accounts and loans are affected, and how vendors can file claims against the receivership." It also notes that the list includes banks which have failed since October 1, 2000. A link to "Failure and Assistance Transactions" is provided. Below this, there's a note about sorting: "Click arrows next to headers to sort in Ascending or Descending order." A search input field is present. The table itself has columns for Bank Name, City, ST, CERT, Acquiring Institution, Closing Date, and Updated Date. The table lists numerous failed banks across various states, such as The First State Bank, Washington Federal Bank for Savings, The Farmers and Merchants State Bank of Arkansas, Evansville County Bank, Germany Bank (with a link to its receivership), First NBC Bank, FirstBank, Geesens Bank and Trust Company, Universal Community Bank, Allentown Bank, The Woodbury Banking Company, First Contractors Bank, First Commerce Bank, First Commerce Bank, First Midwest Bank, First Milwaukee, First Milwaukee, First National Bank, The Bank of Georgia, Elmeras Bank, EdensTrust Bank, First Bank On Earth, Coastal City Bank A Trust Company, Midland Community Bank, First National Bank, The Bank of Chicago, Northern Star Bank, Escobar Bank, FSB DABA II Passes Bank, and The National Peoples Bank of Chicago. The table includes a "Download Data" link, a "Show 25 entries" dropdown, and a navigation footer with links for First, Previous, Next, and Last.

Bank Name	City	ST	CERT	Acquiring Institution	Closing Date	Updated Date
The First State Bank	Casper	WY	10718	Legend Bank, N.A.	May 31, 2018	August 22, 2018
Washington Federal Bank for Savings	Chicago	IL	30870	Royal Savings Bank	December 15, 2017	July 24, 2018
The Farmers and Merchants State Bank of Arkansas	Argenta	AR	17719	Conway Bank	October 13, 2017	August 12, 2018
Evansville County Bank	Saint Albans	IL	1802	United Fidelity Bank, N.A.	May 26, 2017	January 26, 2018
Germany Bank (with link to its receivership)	Milwaukee	WI	30003	First-Citizens Bank & Trust Company	May 8, 2017	March 22, 2018
First NBC Bank	New Orleans	LA	30032	Whitney Bank	April 26, 2017	January 26, 2018
FirstBank	Cottonwood	UT	35495	Cache Valley Bank	March 3, 2017	January 26, 2018
Geesens Bank and Trust Company	Chicago	IL	18328	State Bank of Texas	January 27, 2017	January 26, 2018
Universal Community Bank	Pennsville	NJ	34681	First-Citizens Bank & Trust Company	January 13, 2017	September 20, 2018
Allentown Bank	Mulberry	AR	91	Totally's Bank	September 23, 2016	May 13, 2018
The Woodbury Banking Company	Woodbury	GA	11287	United Bank	August 18, 2016	December 13, 2018
First Contractors Bank	King of Prussia	PA	35112	First-Citizens Bank & Trust Company	May 8, 2016	November 13, 2018
First Commerce Bank	Memphis	TN	9999	The Bank of Fayette County	April 26, 2016	September 14, 2018
First Midwest Bank	Milwaukee	WI	30384	First-Citizens Bank & Trust Company	March 11, 2016	January 26, 2018
First Milwaukee, D&L Bank	Milwaukee	WI	30384	First-Citizens Bank & Trust Company	January 26, 2016	January 26, 2018
Hannover National Bank	Lagrange	GA	36198	Team City Bank	October 2, 2015	February 18, 2018
The Bank of Georgia	Peachtree City	GA	35259	Fidelity Bank	October 2, 2015	July 8, 2018
Elmeras Bank	Denver	CO	34112	United Fidelity Bank, N.A.	July 10, 2015	February 20, 2018
EdensTrust Bank	Chicago	IL	67772	Republic Bank of Chicago	May 8, 2015	January 26, 2018
FirstBank On Earth	San Juan	PR	32152	Banco Popular de Puerto Rico	February 27, 2015	January 26, 2018
Coastal City Bank A Trust Company	Atlanta	GA	33038	First-Citizens Bank & Trust Company	February 13, 2015	January 26, 2018
Midland Community Bank	Chicago	IL	20280	United Fidelity Bank, N.A.	January 23, 2015	November 15, 2017
First National Bank of Crestview	Crestview	FL	17587	First NBC Bank	January 18, 2015	November 15, 2017
Northern Star Bank	Minnetonka	MN	34683	BankHoltz	December 18, 2014	January 3, 2018
Escobar Bank, FSB DABA II Passes Bank	Palm Desert	CA	34708	Bank of Southern California, N.A.	November 7, 2014	November 10, 2018
The National Peoples Bank of Chicago	Chicago	IL	918	State Bank of Texas	October 24, 2014	January 8, 2018



Load <table> data of HTML

- The `pandas.read_html` function has a number of options, but by default it searches for and attempts to parse all `tabular data contained within <table> tags`.
- The result is a list of `DataFrame` objects.

```
# read a HTML file
tables = pd.read_html('ch06/FDICFailedBankList.html')
print(tables)
```

Check the code by yourself

HDF5 Format

- The “HDF” in HDF5 stands for hierarchical data format.
- HDF5 is a well-regarded file format intended for storing large quantities of scientific array data.
- Each HDF5 file can store multiple datasets and supporting metadata.
- Compared with simpler formats, HDF5 supports compression modes, enabling data with repeated patterns to be stored more efficiently.
- HDF5 can be a good choice for working with very large datasets that don’t fit into memory, as you can efficiently read and write small sections of much larger arrays.

Load HDF5 File

Check the code by yourself

```
# read a HDF file  
frame = pd.DataFrame({'a': np.random.randn(100)})  
frame.to_hdf('ch06/mydata.h5', key='frame', mode='w')  
frame2 = pd.read_hdf('ch06/mydata.h5')
```

mydata.h5

```
1 #HDF  
2 SUB  
3 NUL NUL NUL NUL BS BS NUL EOT NUL DLE NUL  
4 NUL  
5 NUL  
6 NUL NUL NUL NUL NUL NUL NUL SO NUL DLE NUL NUL NUL NUL NUL NUL DC1 NUL NUL  
7 àmEö!8ç@Ù'JÖÉÄ!GSEéWçñ?ít@moå?^SÅJñDéçþNAKxSEžñç é:qºETB'?'jBündçöçGSÉðP«çñçy;ñWýzDCLæççvNAKÙ'úåçí,R CAN<çÜ?æíDC  
8 ÁízDbfoîòð?g-'Ii»+?é"jcôçä?)vSEESX?pRS<áSåÈçOSOFsÁéSORù?RSETB(p" è?|ES°_i]SOH@Q1í(Eúç°[vêézÙ?LuMR4á?uåÁçUSüç  
9 öçDC3-DC1tSíççYççç;ESYñjaNUD NUD  
10 NUD SOH NUD NUD
```

Microsoft Excel Files

- pandas also supports reading tabular data stored in Excel 2003 (and higher) files using `pandas.read_excel` function.
- Internally these tools use the add-on packages xlrd and openpyxl to read XLS and XLSX files, respectively.
- You may need to install these manually with pip or conda.

Interacting with Web APIs

- Many websites have public APIs providing data feeds via JSON or some other format.
- There are a number of ways to access these APIs from Python.
- One easy-to-use method is the `requests` package.

Example for Web APIs

- To find the last 30 GitHub issues for pandas on GitHub, we can make a GET HTTP request using the add-on requests library:

Check the code by yourself

```
# find the last 30 GitHub issues for pandas on GitHub
url = 'https://api.github.com/repos/pandas-dev/pandas/issues'
resp = requests.get(url)
print (resp)
data = resp.json()

print (data[0]['title'])
issues = pd.DataFrame(data, columns=['number', 'title'])
print (issues)
```

Example for Web APIs

- The Response object's json method will return a dictionary containing JSON parsed into native Python objects.
- Each element in data is a dictionary containing all of the data found on a GitHub issue page (except for the comments). We can pass data directly to DataFrame and extract fields of interest

Image file and image processing lib

- Image Processing is a field of knowledge that falls in Computer Vision.
- The premises of machine learning were first laid down by computer vision theory, applying a whole set of techniques to process and analyze imagery data to extract valuable information that computers and machines may use for a wide range of applications.
- Some of the commonly used image processing techniques leveraging a very popular Computer Vision library, OpenCV. We can use the functions of [OpenCV lib](#) load image files.

Reference

- Devin Pickell, Structured vs Unstructured Data – What's the Difference?, <https://learn.g2.com/structured-vs-unstructured-data>.
- 21 Places to Find Free Datasets for Data Science Projects, <https://www.dataquest.io/blog/free-datasets-for-projects/>
- Will Koehrsen, Wikipedia Data Science: Working with the World's Largest Encyclopedia, <https://towardsdatascience.com/wikipedia-data-science-working-with-the-worlds-largest-encyclopedia-c08efbac5f5c>
- McKinney, Wes. Python for data analysis: Data wrangling with Pandas, NumPy, and IPython. " O'Reilly Media, Inc.", 2012.
- Provost, Foster, and Tom Fawcett. Data Science for Business: What you need to know about data mining and data-analytic thinking. O'Reilly Media, Inc., 2013.

Data Science

Week 3

Statistics for Data Science