



THE HONG KONG
POLYTECHNIC UNIVERSITY
香港理工大學



PolyU 理大商學院
Business School
Innovation-driven Education and Scholarship

School of
**ACCOUNTING
& FINANCE**
會計及金融學院

Week 6: Introduction to Scientific Computing, Web Protocols, REST API, Accessing and Processing Data

AF3214 Python Programming for Accounting and Finance

Vincent Y. Zhuang, Ph.D.
vincent.zhuang@polyu.edu.hk

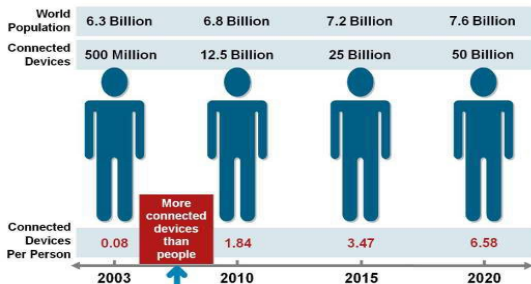
School of Accounting and Finance
The Hong Kong Polytechnic University

R508, 8:30 am – 11:20 am, Tuesdays, Semester 2, AY 2024-25

Why We need Big Data

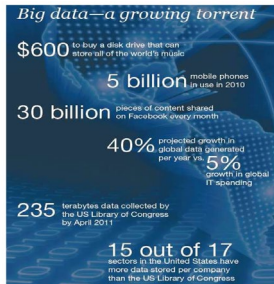
- Facebook generates 10TB of data daily
- Twitter generates 7TB of data daily
- IBM claims 90% of today's stored data was generated in just the last two years.....

Figure 1. The Internet of Things Was "Born" Between 2008 and 2009



Source: Cisco IBSG, April 2011 <https://www.worldometers.info/world-population/>

- Mobile Devices
- Microphones
- Readers/Scanners
- Science facilities
- Programs/ Software
- Social Media
- Cameras



← Examples of data generation points

Big Data Characteristics

- Walmart handles more than 1 million customer transactions every hour.
- Facebook handles 40 billion posts/photos/videos from its user base.
- Decoding the human genome originally took 10 years to process; now it can be achieved in one week.

Three Characteristics Of Big Data – V3s

Volume

- Data quantity

Velocity

- Data Speed

Variety

- Data Types

Volume – Data quantity

- A typical PC might have had 10 GB of storage in 2000.
- Today, Facebook has 4 petabytes (4,096 TB) of new data every day.
- Boeing 737 generates 240 TB of flight data during a single flight across The US. (Automatic Dependent Surveillance-Broadcast (ADS-B): <https://www.flightradar24.com/>)

Velocity – Data speed

- Clickstreams and ad impressions capture user behavior at millions of events per second.
- High-frequency trading algorithms reflect market changes in microseconds.
- Machine to machine processes exchange data between billions of devices.
- On-line gaming systems support millions of concurrent users, each producing multiple inputs per second.

Variety – Data Types

- Big Data isn't just numbers, dates, and strings. Big Data also include **location** data, **3D** data, **audio** and **video**, and **unstructured** text, including **log files** and **social media**.
- Traditional database systems were designed to address smaller volumes of structured data, fewer updates or a predictable, consistent data structure.
- Big Data analysis includes different types of data.

Why We need Big Data

- Growth of Big Data is needed
 - Increase of storage capacities
 - Increase of processing power
 - Availability of data (different data types)

Big Data Analytics

- Examining large amount of data
- Appropriate information (about data)
- Identification of hidden patterns, unknown correlations
- Better business decisions: strategic and operational
- Effective marketing, customer satisfaction, increased revenue

Applications of Big Data Analytics

Smarter
Healthcare



Homeland
Security



Traffic Control



Manufacturing



Multi-channel
sales



Telecom



Trading
Analytics



Search
Quality

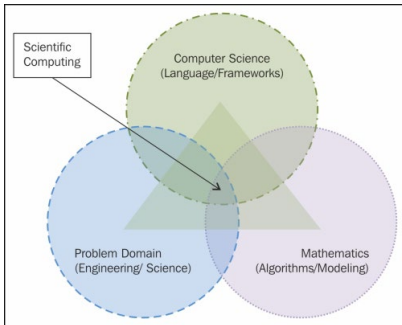


Benefits of Big Data

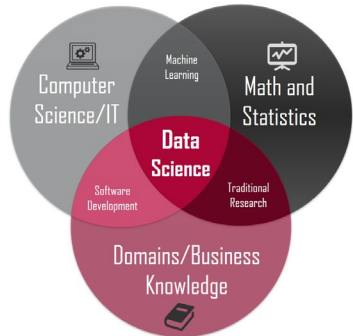
- Real-time big data isn't just a process for storing data in a data warehouse. It's about the ability to make better decisions and take meaningful actions at the right time.
- Technologies give you the scale and flexibility to store data before you know how you are going to process it.

What are Scientific Computing and Data Science?

- "Scientific computing is the collection of tools, techniques and theories required to solve on a computer the mathematical models of problems in science and engineering."
- --Gene H. Golub and James M. Ortega



Scientific computing can be described as an interdisciplinary field, as presented in the above diagram.



Data science is the field of study that uses mathematics, programming, and domain/business knowledge to extract meaningful insights from data. 8

Cornerstones of The Scientific Method



Replication – Numerical calculations should be able to rerun the simulations and replicate the results upon request. Others should also be able to **perform the same calculations** and **obtain the same results**, given the information about their methods.



Reproducibility – Results obtained from an experiment, or an observational study or in a statistical analysis should be achieved again with a high degree of reliability when the study is replicated.

In other words, “**Reproducibility**” refers to instances in which the original researcher's data and computer codes are used to regenerate the results, while “**Replicability**” refers to instances in which a researcher collects new data to arrive at the same scientific findings as a previous study.

FinTech



FinTech? or TechFin?

- TechFin: which refers to technology firms entering the financial field
- It's basically trying to tell you the philosophical way what FinTech is and make it really kind of like it's all driven by machines
- Tech is great...but what is **not** tech? Is excel sheet tech?
- Must be something beyond our imagination. Think in a different way.
- Before embracing, better understand...

FinTech - Cont'd

My view of different parts of Tech that can help you think in the financial industry.

- **Tech** boosts quality service or strengthen one part of traditional financial industry
 - Peer-to-Peer lending on platforms (P2P)
 - **Big Data** → better risk control
 - P2P may not be such a big industry to start with, but it's kinda like a natural supplement to the traditional financial industry
 - Insurance (fraud prevention):
 - **Artificial Intelligence, Machine Learning**
- **Tech** creates a new sector but not revolutionary
 - Say, robotic investment advisory based on AI and ML
- **Blockchain** is significantly different
 - Revolutionary ideas confronts fundamental (and well-studied!) finance/economics principles
- (Almost) no **tech** in ICO.....

Examples of Applications in Accounting and Finance

Risk
Analytics

Real-Time
Analytics

Fraud
Detection

Algorithmic
Trading

Customer
Data
Management

Consumer
Analytics

Providing
Personalized
Services

Where do we Start?



Three Common Data Types in Data Analysis

Cross-Sectional Data:

Consists of a sample data taken in a single time period with many subjects (such as individuals, firms, countries). **Ordering** of the data does **not matter** (↗, ↘, randomized order)

Time Series Data:

Consists of observations on a variable or several variables over time, or a sequence of data points indexed in time order. **Ordering matters**, typically presented in chronological order.

Panel (or Longitudinal) Data:

Consists of a time series for each cross-sectional member in the data. Essentially **combining above two** and **observe the changes over a time series**.

Types of Data - Cont'd

What kind of data is there in Accounting and Finance?



Quantitative Data

numeric

- Security Prices
- Income Statements
- Balance Sheets
- Forecasts
- Many Others



Qualitative Data

non-numeric

- News
- Management Discussion and Analysis in Financial Statements (MD&A)
- Many Others

Types of Data - Cont'd

Quantitative and qualitative data can be gathered from the *same* data unit depending on whether the variable of interest is numerical or categorical

Data unit	Numeric variable	= <u>Quantitative data</u>	Categorical variable	= <u>Qualitative data</u>
A person	"How many children do you have?"	3 children	"In which country were your children born?"	Australia
	"How much do you earn?"	\$60,000 p.a.	"What is your occupation?"	Photographer
	"How many hours do you work?"	38 hours per week	"Do you work full-time or part-time?"	Full-time
A house	"How many square metres is the house?"	200 square metres	"In which city or town is the house located?"	Sydney
A business	"How many workers are currently employed?"	264 employees	"What is the industry of the business?"	Retail
A farm	"How many milk cows are located on the farm?"	36 cows	"What is the main activity of the farm?"	Dairy

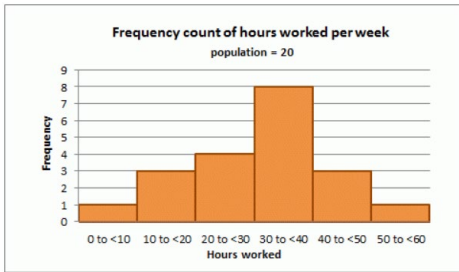
Quantitative data: Quantity, measures of values or counts and are expressed as numbers

Qualitative data: Quality, measures of 'types' and may be represented by a name, symbol, or a number code.

Types of Data - Cont'd

The number of times an observation occurs (frequency) for a data item (variable) can be shown for both quantitative and qualitative data.

Quantitative data

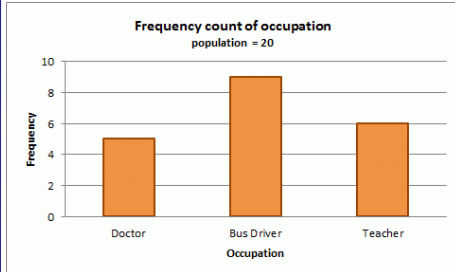


Column graph showing the frequency counts of hours worked per week for 20 people.

Hours worked:

0 to 10—**1**, 10 to 20—**3**,
20 to 30—**4**, 30 to 40—**8**,
40 to 50—**3**, 50 to 60—**1**.

Qualitative data

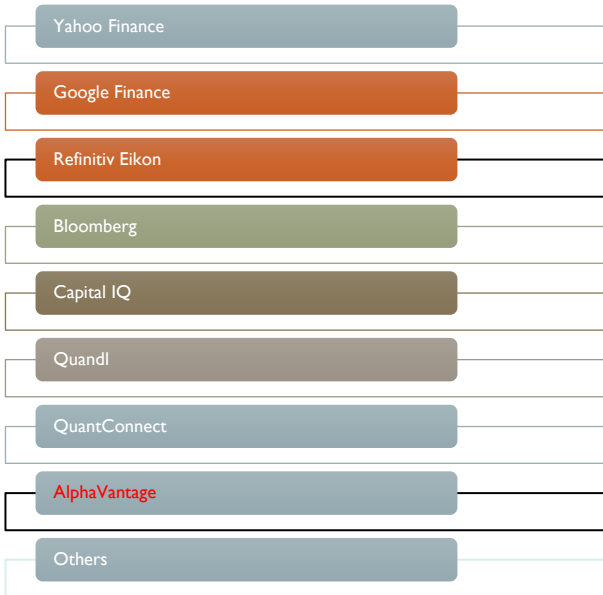


Column graph showing the frequency counts of occupation for 20 people.

Doctor - **5**
Bus driver - **9**
Teacher - **6**

It is important to identify whether the data are quantitative or qualitative as this affects the statistics that can be produced.

Data Sources of Stock Price



Data Cleaning and Manipulating

You will be spending majority of your time **getting** data and **cleaning** and preparing to run through your model or program.

Understanding the structures of data is important!

Types of Structured data

HTML - Hypertext Markup Language, a standardized system for tagging text files to achieve font, colour, graphic, and hyperlink effects on World Wide Web pages.

XML - eXtensible Markup Language (**XML**) is a markup language that **defines** a set of rules for encoding documents in a format that is both human- and machine-readable. The World Wide Web Consortium's **XML 1.0** Specification and several other related specifications—all of them free open standards—**define XML**. It was designed to store and transport data.

CSV - A comma-separated values file is a delimited text file that uses a comma to separate values. A CSV file stores tabular data in plain text. Each line of the file is a data record. Each record consists of one or more fields, separated by commas.

JSON - In computing, JavaScript Object Notation is an open-standard file format that uses human-readable text to transmit data objects consisting of attribute-value pairs and array data types.

Calculating Return from Price Data

The reason for using returns versus prices is normalization. This allows us to measure all variables in a comparable metric.

It is a common practice to use log returns in finance.

$$P_t = P_0(1 + r) = P_0 e^R \Rightarrow \ln \frac{P_t}{P_0} = \ln(1 + r) = \ln e^R = R$$

Time	Price	r	log
0	100	-	-
1	120	20.0%	18.2%
2	100	-16.7%	-18.2%

Total Return: 3.3%? 0%

P_0 – Initial price of the stock

P_t – Price of the stock at the end of the period

R – continuously compounded rate over the period

r – simple return of the stock over time t

Descriptive/Summary Statistics

Mean/Average - The mean (or average) is the most popular and well-known measure of central tendency.

Mode - The mode is the most frequent score in our data set.

Median - The median is the middle score for a set of data that has been arranged in order of magnitude.

Descriptive/Summary Statistics - Cont'd

Variance – The averaged of the squared differences from the mean. first calculate the difference between each point and the mean; then, square and average the results.

Standard Deviation (Volatility) – Measure of how spread out numbers are. The square root of variance by figuring out the variation between each data point relative to the mean. If the points are further (closer) from the mean, there is a higher (lower) deviation.

$$S^2 = \frac{\sum (x_i - \bar{x})^2}{n - 1}$$

S^2 = sample variance

x_i = the value of the one observation

\bar{x} = the mean value of all observations

n = the number of observations

$$\left(\begin{array}{c} \text{Standard} \\ \text{Deviation} \end{array} \right) = \sqrt{\text{Variance}}$$

$$\sigma = \sqrt{\sigma^2}$$

For traders and analysts, these two concepts are of paramount importance as they are used to measure security and market volatility, which in turn plays a large role in creating a profitable trading strategy.

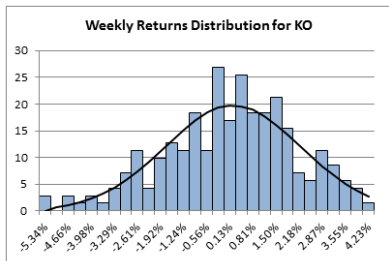
Correlation Matrix

- Fundamental tool for stock market investors. It describes how closely the returns of the assets in a portfolio are correlated.
- A correlation matrix is a table showing correlation coefficients between variables.
- The line of 1.00s going from the top left to the bottom right is the main *diagonal*, which shows that each variable always perfectly correlates with itself.
- Usually, in statistics, we measure four types of correlations: Pearson correlation, Kendall rank correlation, Spearman correlation, and Point-Biserial correlation.

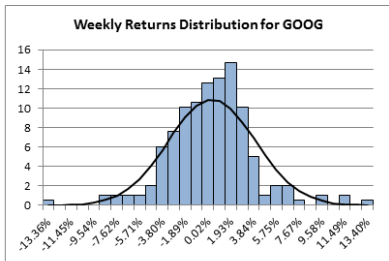
	◊ iPad ▾	◊ iPhone ▾	◊ iPod ▾	◊ Nokia mobile phone ▾	◊ Other mobile phone (not Nokia and not iPhone) ▾	◊ Mac computer - desktop ▾	◊ Mac computer - laptop ▾	◊ PC (non-Mac) ▾	◊ Laptop computer (non-Mac) ▾	◊ None of these ▾
iPad	1.000	.219	.143	.004	-.128	.081	.132	-.007	.009	-.027
iPhone	.219	1.000	.248	-.165	-.332	.161	.191	-.158	.152	-.053
iPod	.143	.248	1.000	-.004	-.054	.186	.238	-.012	.058	-.069
Nokia mobile phone	.004	-.165	-.004	1.000	-.423	.047	-.070	.067	-.059	-.106
Other mobile phone (not Nokia and not iPhone)	-.128	-.332	-.054	-.423	1.000	-.082	-.152	.148	.111	-.101
Mac computer - desktop	.081	.161	.186	.047	-.082	1.000	.203	-.127	.016	-.034
Mac computer - laptop	.132	.191	.238	-.070	-.152	.203	1.000	-.151	-.104	-.034
PC (non-Mac)	-.007	-.158	-.012	.067	.148	-.127	-.151	1.000	-.104	-.159
Laptop computer (non-Mac)	.009	.152	.058	-.059	.111	.016	-.104	-.104	1.000	-.148
None of these	-.027	-.053	-.069	-.106	-.101	-.034	-.034	-.159	-.148	1.000

Visualizing Data

Histograms - used to graphically summarize and display the distribution of a data set.



Mean	0.002196
Standard Deviation	0.020183
Skew	-0.21837
Kurtosis	-0.00766

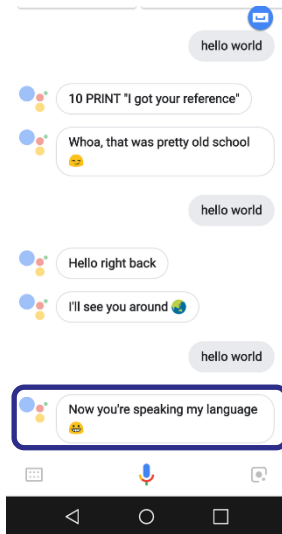


Mean	0.00363
Standard Deviation	0.036469
Skew	0.40915
Kurtosis	2.94533

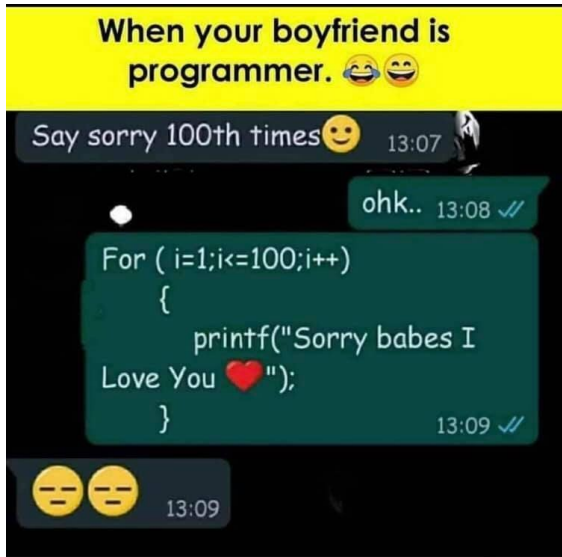
The best is yet to come...



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



Why Python?



It has a large community of users, easy to find help and documentation.

Extensive ecosystem of scientific libraries and environments

Why Python?

Good support for Parallel processing with processes and thread, interprocess communication, GPU computing

Readily available and suitable for use on high-performance computing clusters

```
public class Main {  
    public static void main(String[] args) {  
        System.out.println("hello world");  
    }  
}
```

Why Python?

The majority of Dropbox code is written in Python, and the initial Dropbox product was almost entirely written in Python.

```
print('hello world')
```

Guido van Rossum, the creator of Python, used to work at Dropbox.

Python - Cont'd

Jupyter Notebook – An HTML-based notebook environment for Python. It is a cell-based environment with great interactivity, where calculations can be organized and documented in a structured way.

Spyder – is an IDE for scientific computing with Python. It is a powerful code editor, with syntax high-lighting, dynamic code introspection and integration with Python debugger.

IDE: An integrated development environment (IDE) is a software application that provides comprehensive facilities to computer programmers for software development.

Visual Studio Code - <https://code.visualstudio.com/>

An integrated development environment, supports for debugging, syntax highlighting, intelligent code completion, snippets, code refactoring, and embedded Git.

In the StackOverflow 2021 Developer Survey, Visual Studio Code was ranked the most popular developer environment tool, with 70% of 82,000 respondents reporting that they use it.[^]

[^] <https://insights.stackoverflow.com/survey/2021#section-most-popular-technologies-integrated-development-environment>

Python - Cont'd

What is a Package?

A package contains all the files you need for a module.

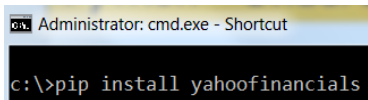
Modules are Python code libraries you can include in your project.

What is PIP?

PIP is a package manager for Python packages/module.

Download a Package

Open the command line interface and tell PIP to download the package you want. Navigate your command line to the location of Python's script directory, and type the following:



```
Administrator: cmd.exe - Shortcut
c:\>pip install yahoofinancials
```

```
1 import urllib
2 import urllib.request
3 from urllib.request import urlopen
4 import json
5 import time
6 from datetime import datetime
7 import requests
8 import csv
9 from urllib.parse import urlparse
10 import datetime
11 import sys
12 import io
13 import pandas as pd
```



Using a Package

Python - Understanding Data Structure

yahoofinancials

Financial data module used for pulling both fundamental and technical data from Yahoo Finance

Tesla



The module returns stock, forex, cryptocurrency, mutual fund, ETF, commodity futures, and US Treasury financial data from Yahoo Finance.

<https://pypi.org/project/yahoofinancials/>

```
"incomeStatementHistory": {
  "TSLA": [
    {
      "2020-12-31": {
        "researchDevelopment": 1491000000,
        "effectOfAccountingChanges": null,
        "incomeBeforeTax": 1154000000,
        "minorityInterest": 1454000000,
        "netIncome": 721000000,
        "sellingGeneralAdministrative": 3188000000,
        "grossProfit": 6630000000,
        "ebit": 1951000000,
        "operatingIncome": 1951000000,
        "otherOperatingExpenses": null,
        "interestExpense": -784000000,
        "extraordinaryItems": null,
        "nonRecurring": null,
        "otherItems": null,
        "incomeTaxExpense": 292000000,
        "totalRevenue": 31536000000,
        "totalOperatingExpenses": 29585000000,
        "costOfRevenue": 24906000000,
        "totalOtherIncomeExpenseNet": -797000000,
        "discontinuedOperations": null,
        "netIncomeFromContinuingOps": 862000000,
        "netIncomeApplicableToCommonShares": 690000000
      }
    }
  ],
},
```

```
demo.py > ...
1 from yahoofinancials import YahooFinancials
2 import json
3 yahoo_financials = YahooFinancials('TSLA')
4 json_file = yahoo_financials.get_financial_stmts('annual', 'income')
5 with open('data.json', 'w') as f:
6     json.dump(json_file, f)
```


Example of HTML Data from SEC EDGAR

How to read it in Python?

Python Modules

- Pandas - *Pandas* is a library providing high-performance, easy-to-use data structures and data analysis tools for Python programming language.
- NumPy - *NumPy* is the fundamental package for scientific computing with Python. It offers comprehensive mathematical functions, random number generators, linear algebra routines, Fourier transforms, and more.

Pandas Tutorial:

<https://www.kaggle.com/learn/pandas>

Numpy Library Lookup:

<https://cs231n.github.io/python-numpy-tutorial/#numpy>



NumPy



SciPy



Scikit

matplotlib



Python Modules - Cont'd

- Scikit-learn is a free machine learning library for Python programming. It features various classification, regression and clustering algorithms.
- Matplotlib is a Python 2D plotting library which produces publication quality figures in a variety of hardcopy formats and interactive environments across platforms.

Scikit-learn Tutorial:

<https://www.datacamp.com/community/tutorials/scikit-learn-python>

Matplotlib Tutorial:

<https://www.datacamp.com/community/tutorials/matplotlib-tutorial-python>



NumPy



SciPy



Scikit



matplotlib

Let's move to Jupyter Notebook Demo

What are Internet and WWW?

- On August 6th 1991 the very first website went online. Marked the birth of the world wide web and technological revolution.
- It was the beginning of one of the biggest achievements of the 20th century - the world wide web: <http://info.cern.ch/hypertext/WWW/TheProject.html>
- The man behind it is computer scientist Sir Tim Berners-Lee. In 1980, he worked as an independent contractor for CERN - The European Council for Nuclear Research.

World Wide Web

The WorldWideWeb (W3) is a wide-area [hypermedia](#) information retrieval initiative aiming to give universal access to a large universe of documents.

Everything there is online about W3 is linked directly or indirectly to this document, including an [executive summary](#) of the project, [Mailing lists](#), [Policy](#), November's [W3 news](#), [Frequently Asked Questions](#).

[What's out there?](#)

Pointers to the world's online information, [subjects](#), [W3 servers](#), etc.

[Help](#)

on the browser you are using

[Software Products](#)

A list of W3 project components and their current state. (e.g. [Line Mode](#), [X11 Viola](#), [NeXTStep](#), [Servers](#), [Tools](#), [Mail robot](#), [Library](#))

[Technical](#)

Details of protocols, formats, program internals etc

[Bibliography](#)

Paper documentation on W3 and references.

[People](#)

A list of some people involved in the project.

[History](#)

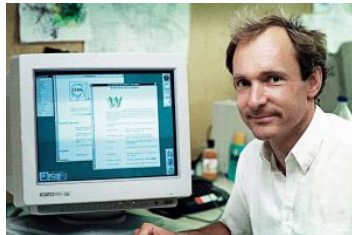
A summary of the history of the project.

[How can I help?](#)

If you would like to support the web..

[Getting code](#)

Getting the code by [anonymous FTP](#), etc.



What are Internet and WWW?

What's the difference between the Web and Internet?



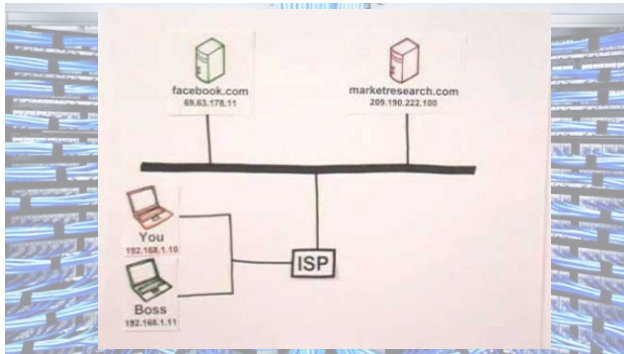
Web = Internet?



but they actually mean different things.

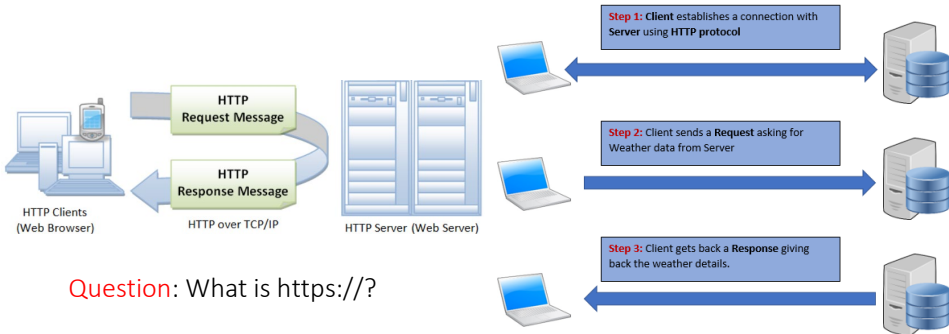
How does the Internet Work?

- We usually don't have to think too much about what is happening on our browser when we are surfing the internet.
- However it is important to understand the mechanisms of the Internet when we are surfing the internet for data.
- We need to understand what is happening at the browser level, as well as the network level.



HTTP Protocol

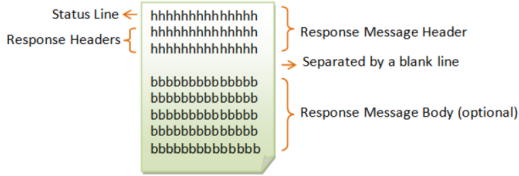
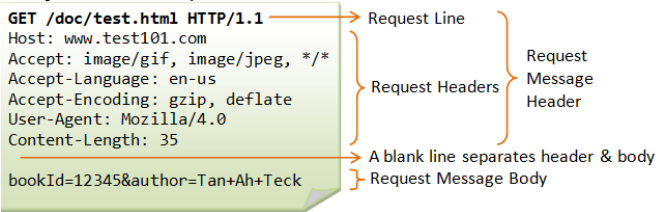
- HTTP (**H**yper**T**ext **T**ransfer **P**rotocol): The standard protocol for transferring web pages and content across the internet.
- When you browse a web page, the URL might be preceded by <http://>. This is telling the web browser to use HTTP to transfer the data.
- Web Browser is the Client and the request are sent to the servers.



HTTP Request

- HTTP Request is a packet of information that one computer sends to another computer to communicate something.
- A HTTP Request contains following parts.
 - Request Line
 - Headers
 - Optional Body of the Request

HTTP request message



HTTP Response Message

HTTP response message

The response headers are in the form: name:value pairs

HTTP Request - Cont'd

- It specifies the method token:

➤ GET

➤ POST

➤ PUT

➤ HEAD

➤ DELETE

➤ TRACE

➤ OPTIONS

```
GET /doc/test.html HTTP/1.1
```

```
Host: www.test101.com
```

```
Accept: image/gif, image/jpeg, /*
```

```
Accept-Language: en-us
```

```
Accept-Encoding: gzip, deflate
```

```
User-Agent: Mozilla/4.0
```

```
Content-Length: 35
```

```
bookId=12345&author=Tan+Ah+Teck
```

Request Line

Request Headers

Request
Message
Header

A blank line separates header & body

Request Message Body

```
HTTP/1.1 200 OK
```

```
Date: Sun, 08 Feb xxxx 01:11:12 GMT
```

```
Server: Apache/1.3.29 (Win32)
```

```
Last-Modified: Sat, 07 Feb xxxx
```

```
ETag: "0-23-4024c3a5"
```

```
Accept-Ranges: bytes
```

```
Content-Length: 35
```

```
Connection: close
```

```
Content-Type: text/html
```

```
<h1>My Home page</h1>
```

Status Line

Response Headers

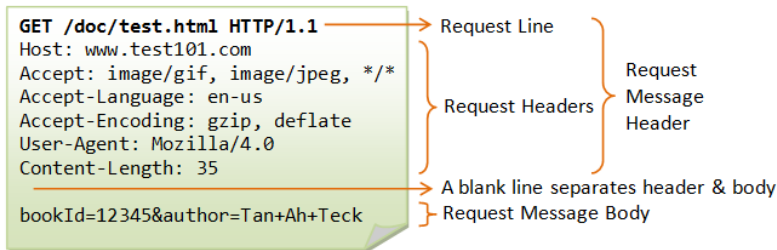
Response
Message
Header

A blank line separates header & body

Response Message Body

HTTP Request - Cont'd

- **GET** is issued to request data from a specified resource.
- It is one of the most common HTTP Methods.
- It can be **cached** and **remains** in browser history
- It should not be used when dealing with sensitive data.
- There is a length restriction – max 2,048 characters.



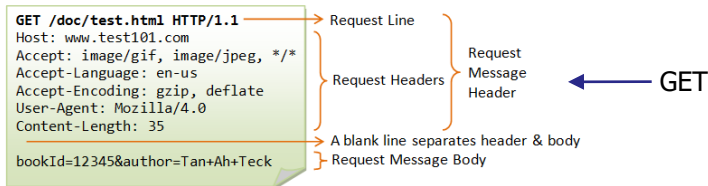
HTTP Request - Cont'd

- If you are using the GET method, you are limited to a maximum of 2,048 characters.
- However, the POST method is not limited by the size because they are transferred in the header and not in the URL.
- **POST** is used to send data to a server to create/update a resource.
- The data sent to the server with POST is stored in the request body.
- POST requests are never cached
- POST requests do not remain in the browser history.
- Has no restrictions on data length.

The HTTP Method	Path to the source on Web Server	Protocol Version Browser supports
The Request Headers	Post /RegisterDao.jsp HTTP/1.1	
	Host: www.javatpoint.com	
	User-Agent: Mozilla/5.0	
	Accept: text/xml,text/html,text/plain,image/jpeg	
	Accept-Language: en-us,en	
	Accept-Encoding: gzip,deflate	
	Accept-Charset: ISO-8859-1,utf-8	
	Keep-Alive:300	
	Connection:keep-alive	
	User=ravi&pass=java	Message body

HTTP Request - Cont'd

- HTTP Request may contain zero or more Request **Headers**.
- In between the Request Line and the Message body is considered as Header.
- An HTTP header consists of its case-insensitive name followed by a colon (:), then by its value.
- Header is used to pass additional information about the request to the server.



HTTP Request - Cont'd

- Request message body is the part of the HTTP request where additional content can be sent to the server.
- Request message body tries to send additional information required by the server to process current request properly.

(a) Request message		(b) Response message	
GET /test/hi-there.txt HTTP/1.0	Start line	HTTP/1.0 200 OK	
Accept: text/* Accept-Language: en,fr	Headers	Content-type: text/plain Content-length: 19	
	Body	Hi! I'm a message!	

Question: What is text/*?

General

Request URL: https://af.polyu.edu.hk/
Request Method: GET
Status Code: 200 OK
Remote Address: 158.132.48.76:443
Referrer Policy: strict-origin-when-cross-origin

Response Headers View source

Cache-Control: private
Content-Length: 67975
Content-Type: text/html; charset=utf-8
Date: Wed, 15 Sep 2021 07:27:00 GMT
Strict-Transport-Security: max-age=10886400; preload
X-Content-Type-Options: nosniff
X-Frame-Options: sameorigin
X-XSS-Protection: 1; mode=block

Request Headers View parsed

GET / HTTP/1.1
Host: af.polyu.edu.hk
Connection: keep-alive
sec-ch-ua: "Google Chrome";v="93", " Not;A Brand";v="99", "Chromium";v="93"
sec-ch-ua-mobile: ?0
sec-ch-ua-platform: "Windows"
Upgrade-Insecure-Requests: 1
DNT: 1
User-Agent: Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/93.0.4577.63 Safari/537.36
Accept: text/html,application/xhtml+xml,application/xml;q=0.9,image/avif,image/webp,image/apng,*/*;q=0.8,application/signed-exchange;v=b3;q=0.9
Sec-Fetch-Site: same-origin
Sec-Fetch-Mode: navigate
Sec-Fetch-User: ?1
Sec-Fetch-Dest: document
Referer: https://af.polyu.edu.hk/about/message-from-head/
Accept-Encoding: gzip, deflate, br
Accept-Language: en-US,en;q=0.9,zh-CN;q=0.8,zh-TW;q=0.7,zh;q=0.6
Cookie: _ga_M8886P9Rj=GS1.1.1630911293.2.1.1630911509.0; _gcl_eu=1.1.692502002.1630914032; _fbp=fb.2.16309140328001.662757297; _ga=GA1.3.644692572.1630908296; cookieNotice=accepted; s_pers=%20v8%3D01631604419059%7C1726212419059%3B%20v8_s%3D0More%2520than%25207%2520days%7C1631606219059%3B%20c19%3Dpr%253Apure%2520portal%253Apersons%253Aview%7C1631606219064%3B%20v6%3D01631604415392%7C1631606219071%3B; AMCV_406368F454EC41940A4C

HTTP Response

After receiving and interpreting a request message, a server responds with an HTTP response message.

- The status line consists of 3 parts
 - ❖ HTTP Protocol Version
 - ❖ Status Code
 - ❖ Reason Phrase: is intended to give a short textual description of the Status Code.

Examples of status line are:

```
HTTP/1.1 200 OK
HTTP/1.0 404 Not Found
HTTP/1.1 403 Forbidden
HTTP/1.1 500 Internal
                Server Error
HTTP/1.1 502 Bad Gateway
```

HTTP/1.1 200 OK

Date: Sun, 08 Feb xxxx 01:11:12 GMT
Server: Apache/1.3.29 (Win32)
Last-Modified: Sat, 07 Feb xxxx
ETag: "0-23-4024c3a5"
Accept-Ranges: bytes
Content-Length: 35
Connection: close
Content-Type: text/html

<h1>My Home page</h1>

→ Status Line

→ Response Headers

} Response
Message
Header

→ A blank line separates header & body

} Response Message Body

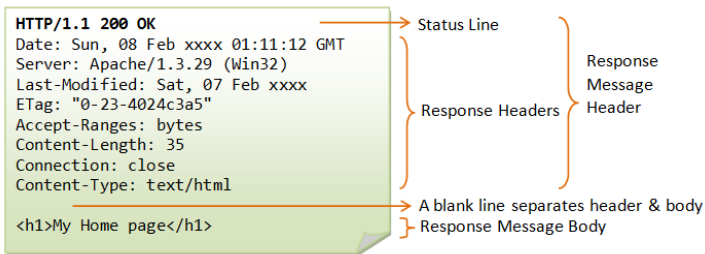
HTTP Response: Status Line

100	Continue	409	Conflict
101	Switching Protocols	410	Gone
102	Processing	411	Length Required
2XX Success		412	Precondition Failed
200	OK	413	Payload Too Large
201	Created	414	Request-URI Too Long
202	Accepted	415	Unsupported Media Type
203	Non-authoritative Information	416	Requested Range Not Satisfiable
204	No Content	417	Expectation Failed
205	Reset Content	418	I'm a teapot
206	Partial Content	421	Misdirected Request
207	Multi-Status	422	Unprocessable Entity
208	Already Reported	423	Locked
226	IM Used	424	Failed Dependency
3XX Redirectional		426	Upgrade Required
300	Multiple Choices	428	Precondition Required
301	Moved Permanently	429	Too Many Requests
302	Found	431	Request Header Fields Too Large
303	See Other	444	Connection Closed Without Response
304	Not Modified	451	Unavailable For Legal Reasons
305	Use Proxy	499	Client Closed Request
307	Temporary Redirect	5XX Server Error	
308	Permanent Redirect	500	Internal Server Error
4XX Client Error		501	Not Implemented
400	Bad Request	502	Bad Gateway
401	Unauthorized	503	Service Unavailable
402	Payment Required	504	Gateway Timeout
403	Forbidden	505	HTTP Version Not Supported
404	Not Found	506	Variant Also Negotiates
405	Method Not Allowed	507	Insufficient Storage
406	Not Acceptable	508	Loop Detected
407	Proxy Authentication Required	510	Not Extended
408	Request Timeout	511	Network Authentication Required
		599	Network Connect Timeout Error

More reading: https://en.wikipedia.org/wiki/List_of_HTTP_status_codes

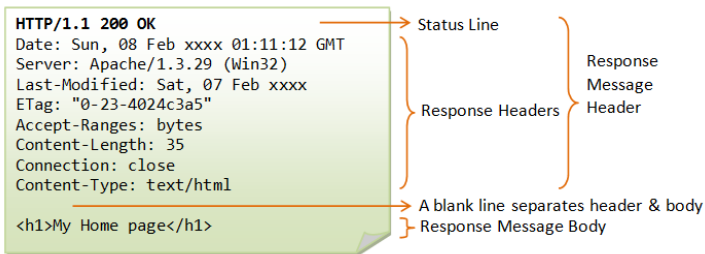
HTTP Response: Response Header

- Just like the Request header, it can contain zero or more lines.
- Very uncommon to have zero headers in the response.
- In Response header, there is a header named content-type. It is used to inform the client that body of response is a certain type.



HTTP Response: Response Body

- Contains the information requested.
- In terms of Web Services, the information requested by a client is referred to as a resource.



Request and Python

It is possible to use **Python** to send these HTTP Request and receive the HTTP Response.

There are many libraries that will help you achieve this goal, but we will concentrate on simplicity to start off.

Python Request library is best for making simple web requests.

<https://pypi.org/project/requests/>

We will have a try later on Jupyter Notebook.

Beautiful Soup

Processing HTML With Beautiful Soup 4

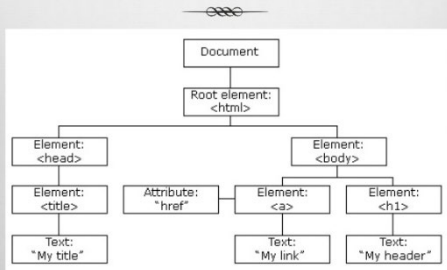
- Beautiful Soup is a python library for pulling data out of HTML and XML files.
- <https://www.crummy.com/software/BeautifulSoup/bs4/doc/>



What is HTML?

- HTML (**H**yper**T**ext **M**arkup **L**anguage) is the standard markup language for document designed to be displayed in a web browser.
- HTML elements are the building blocks of HTML pages. They provide a means to create structured documents.

HTML DOM Tree of Objects



My First Heading

My first paragraph.

Document Object Model (DOM): a cross-platform and language-independent interface that treats an XML or HTML document as a tree structure wherein each node is an object representing a part of the document

SEC EDGAR

- EDGAR (Electronic Data Gathering, Analysis and Retrieval system) performs automated collection, validation, indexing, acceptance, and forwarding of submission by companies and others who are required by law to file forms with the U.S. SEC (U.S. Securities and Exchange Commission).
- The database contains rich information which is freely available to the public via the internet.

Fair access

- Current max request rate: 10 requests/second.

To ensure everyone has equitable access to SEC EDGAR content, please use efficient scripting. Download only what you need and please moderate requests to minimize server load.

SEC reserves the right to limit request rates to preserve fair access for all users. See our [Internet Security Policy](#) for our current **rate request limit**.

The SEC does not allow botnets or automated tools to crawl the site. Any request that has been identified as part of a botnet or an automated tool outside of the acceptable policy will be managed to ensure fair access for all users.

Please declare your **user agent** in request headers:

Sample Declared Bot Request Headers:

User-Agent:	Sample Company Name AdminContact@<sample company domain>.com
Accept-Encoding:	gzip, deflate
Host:	www.sec.gov

We do not offer technical support for developing or debugging scripted processes.

More reading: <https://www.sec.gov/os/accessing-edgar-data>

SEC Filing Types (selected)

- **10-K:**

A 10-K is a comprehensive report filed **annually** by a publicly traded company about its financial performance in the year and is required by the SEC.

- **10-Q:**

A 10-Q is a comprehensive report filed **quarterly** by a publicly traded company about its financial performance in the quarter and is required by the SEC.

- **8-K:**

A 8-K report is of unscheduled material events or corporate changes at a company that could be of importance to the shareholders or the SEC. Such as acquisitions, bankruptcy, the resignation of directors, or changes in the fiscal year.

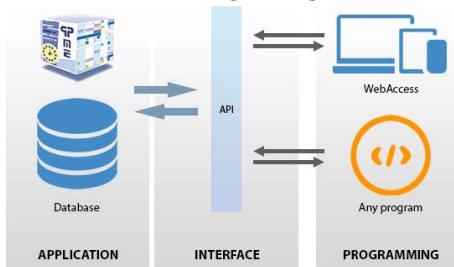
There are many different filings types and can be referenced below.

API / REST API

- An **API** is a “*Application Programming Interface*”.
- An API is a set of **rules** and **protocols** for building and integrating applications.
- Basically an API specifies how software components should interact.
- A good API makes it easier to develop a program by providing all building blocks.
- You can think of an API as a **mediator** between the users or clients and the resources or web services.
- Usually for organizations to share resources and information while maintaining security, control, and authentication - determining who gets access to what.

HKMA Open API

<https://apidocs.hkma.gov.hk/>



API / REST API - Cont'd

What is API?

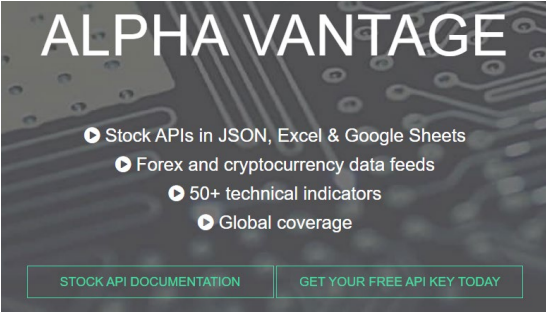


STOCK DATA API

An example:

ALPHA VANTAGE

<https://www.alphavantage.co>

A promotional banner for Alpha Vantage with a dark background featuring a circuit board pattern. The text is white and green. At the bottom, there are two green buttons with white text.

ALPHA VANTAGE

- ▶ Stock APIs in JSON, Excel & Google Sheets
- ▶ Forex and cryptocurrency data feeds
- ▶ 50+ technical indicators
- ▶ Global coverage

STOCK API DOCUMENTATION GET YOUR FREE API KEY TODAY

API / REST API - Cont'd

- **REST** is **R**epresentational **S**tate **T**ransfer – architectural style for distributed hypermedia systems.
- REST is a set of architectural constraints, not a protocol or a standard. API developers can implement REST in a variety of ways.
- It is a set of rules that developers follow when they create their API.
- It is architectural popular way to structure the backend of a web application.
- In a nutshell, REST determines how the API looks like.
- Web applications are usually divided up into two separate developments team.
 - Front-end – Concentrates on the user experience, the UI, and aspects that people see and interact with.
 - Backend – Server-side development.

An example: InteractiveBrokers

<https://www.interactivebrokers.com/en/index.php?f=45185>

The End

JUPYTER NOTEBOOK