

For eg → email $\xrightarrow{\text{spam}}$ $\xrightarrow{\text{not spam}}$.

A classification

→ Python Data Structures

- String
- List
- Tuple
- Dictionary

→ `>>> stuff = 'Hello world'` String

→ `>>> type(stuff)` → type of datastructure
← class 'str' >

→ `>>> dir(stuff)` → All the available methods of data structure
['capitalize', 'casefold', 'center'...]

NOTE → In dir() ignore the ones with underscores they are used by Python.

In although hidden is still counted as character.

→ T-Files

→ `>>> stuff = 'X\my'`

→ `>>> print(stuff)`

X

Y

→ `>>> len(stuff)`

3

→ `>>> fhand = open('mbox-short.txt')`

→ `>>> inp = fhand.read()`

→ `>>> print(len(inp))`

94626

• read() file ko lambi string mein convert keda hai

Dictionary

Count

counts = dict()

names = ['a', 'b', 'c', 'd']

for name in names:

If name not in counts
counts[name] = 1

else:

counts[name] += 1



counts[name] = counts.get(name, 0) + 1

In this case

- get() returns the value to key (name)
if it doesn't exists it returns default(0)

Tuples

dict.items() → Return List of tuples of key value pair

Eg: dict = {'a': 2, 'b': 5, 'c': 1, 'd': 10}

list = [] or list()

for key, val in dict.items()

new tup = (val, key)] → Value first
list.append(new tup)] Then key

list = sorted(list, reverse=True)] → sort based
on value.

Short version imagine as 'forall'

`>>> point (selected ([(v,k) for k,v in dict.items()]))`

\Rightarrow Tuples are comparable

>>> (0, +, 2) < (5, +, 2) → checks first ignores
Done rest i.e

`>>> (0, 1, 2000) < (0, 3, +)` → checks first
True

>>> ('Jones', 'Sally') \in ('Jones', 'Sam') Now second
True 1 < 3 \rightarrow True

⇒ PYTHON TO ACCESS WEB DATA.

→ Regular Expressions / Regex

↖ → ^{matches} beginning of line

\$ → matches the end of line

• → matches any character also called wild card

`\s` → matches whitespace

\S → matches any non-whitespace character

* \rightarrow ~~repeals~~ matches character zero or more time

*? → same as * but non-greedy

+ → Repetition of a character one or more times

+? → same as + but non-greedy

[aeiou] → Matches single character in listed set.

$[^XYZ]$ → Matches single character not in listed set

$[a-z \cup 0-9]$ → The set of characters can include range

library → import ee

`El. search()` → to see if string matches a reg exp.

El. findall (ℓ) \rightarrow extract portion of sking that matches regexp.

Eg -

hand = open ('abc.txt')
for line in hand:
 print line

```
line = line. skip()
```

`[if for line. find('Eom: ') >= 0] ← In index is
point (line) greater than 0`

1

if re.search('From:', line); → It returns
point (line). True/False.

if line.startswith('From:'): → if re.search('From:', line)

Regex used is $[0-9]^*$ ^{any digit} _{one or more times.}

>>> x = 'My 2 fav numbers are 19 and 42'
>>> y = re.findall('([0-9]+)', x)
>>> print(y)
['2', '19', '42']

* and + uses GREEDY MATCHING i.e.
the push ff outward and find largest
possible string

>>> x = 'From: Using the: character'
>>> y = re.findall('^F.+:', x)
>>> print(y)
['From : Using the:']

and not
['From:']
because its greedy.

First F ^{One or more char}
 $\wedge F . + :$
first last character

>>> y = re.findall('^.+?:', x) ^{not greedy.}
>>> print(y)
['From:']

'@([n])*' ^{One or more.}
match non
blank character

Eg abc@d.com
d.com gets
extracted because
of ()

\$ → It stands for something in regex
but \\$ → It means real dollar sign.

Socket → It is an endpoint of a bi-directional inter-process communication flow across an Internet Protocol based computer network, such as the internet.

Using in python

just make connection

```
import socket  
mysock = socket.socket(socket.AF_INET,  
                      socket.SOCK_STREAM)  
mysock.connect(('data.pyc.org', 80))
```

Host name Port

cmd = 'GET http://data.pyc.org/promo.dat HTTP/1.0
/n/n'.encode()

depicts 2 enter. if error replace with ctrl+z

mysock.send(cmd)

while True:

```
data = mysock.recv(512)  
if len(data) < 1:  
    break  
print(data.decode())  
mysock.close.
```

in file If no data received

>>> print(ord('H'))

72

ord() helps with
ASCII

• encode() → string to Bytes

Since HTTP is common, we have library that does all socket work.

```
import urllib.request, urllib.parse, urllib.error.  
fhand = urllib.request.urlopen('http://data.pr4e.org/  
romeo.txt')
```

for line in fhand:

print(line.decode().strip())

↳ to remove extra \n as print puts \n for next iteration of loop



NOTE → It gives only data not header.

For headers → headers = dict(fhand.getheaders())

PARSING WEB PAGES / WEB SCRAPPINGS.

→ Intro to Beautiful Soup
Import BeautifulSoup.

url = input()

html = urllib.request.urlopen(url, context=ctx).read()

Soup ← soup = BeautifulSoup(html, 'html.parser').
One string

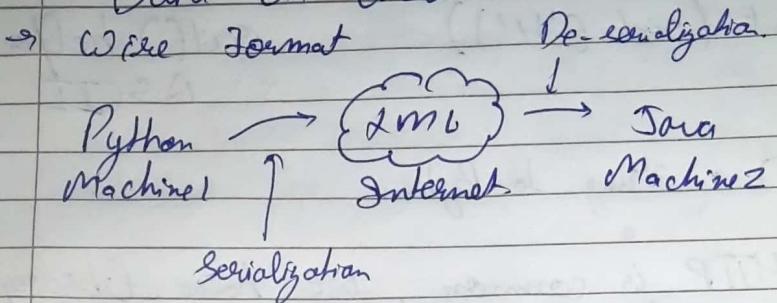
Retrieve all anchor tags

tags = soup('a')

for tag in tags:

print(tag.get('href', None)).

Data on Web.



XML

1. ← </person> → 3
1. Start tag
2. End tag
3. Text content
4. Attribute
5. Self closing tag.
2. ← </person>

Using in python

```
import xmltree as xml.etree.ElementTree as ET
data = '''<person>
    <name> Chuck </name>
    <phone type="int">
        +17343034456
    </phone>
    <email hide="yes"/>
</person>'''
```

- tree = ET.fromstring(data) → forms an XML tree and if fails if bad XML in data
print(tree.find('name').text)

↓ data
tag name

print(tree.find('email').get('hide'))
↓ attribute error
tag name. returns name.
yes ~~error~~

stuff = ET.fromstring (input)

list = stuff.findall('users/user') or (.//user)

list of tags

parent = user
child = user

print (len(list)) → total no. of users.

for item in list:

print (item.find('name').text) → name of all users one by one

XML → complex JSON → simple

very powerful

JSON → Java Script Object Notation.

use in python.

```
import json
data = """
    "name": "Chuck",
    "phone": {
        "type": "int",
        "number": "+1 234"
    },
    "email": {
        "hide": "yes"
    }
"""
print (data)
```

y''' when taking from web make sure to apply .read() first

info = json.loads (data)

print (info ["name"]) / print (info ["email"] ["hide"])

print(`len(info)`)

JSON represents data as nested "lists" and "dictionaries".

API → Application Programming Interface

→ API Security and Rate Limiting.

Google map free limit till → 2500 req. per day

agar json ke output ke case b' logo
ho it means it is logo array and
json is not decoded.

PYTHON & DATABASES

Object Oriented Programming

Classes → a - A template

Method or Message - A defined capability of class

Field or attribute - A bit of data in class

Object or Instance - a particular instance of class

Class PA:

`x=0`

```
def party(self)
    self.x = self.x + 1
    print(self.x)
```

An. PA()

An.party()

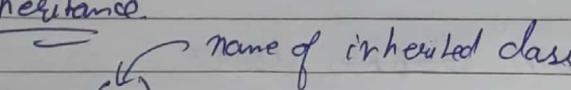
Self.x means an.x
because an is an
object calling party.

Constructor def __init__(self)

Deconstructor def __del__(self)

Constructor is typically used to set up variables
Destructor is seldom used

Inheritance

class (B) 

Relational Database

Database

Relation (or Table)

- Set of tuples with same Attributes

Row (or tuple)

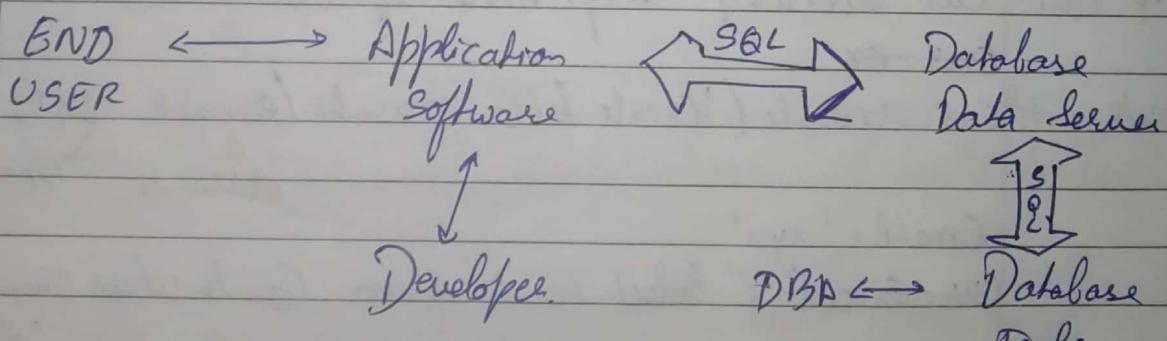
- An object & its info.

Attribute (columns, field)

-

SQL - Structured Query Language

LARGE PROJECT STRUCTURE



Database Model / Schema \Rightarrow Structure & / Format of database

Create Table \rightarrow CREATE TABLE Users {

 name VARCHAR(128),
 email VARCHAR(128)
}

In python `cur.execute('...')` → execute one SQL
`cur.executescript('...')` → multiple SQL
lines separated by ;

`INSERT INTO Users (name, email) VALUES ('Keistin', 'kf@umich.edu')`

Note → `DELETE FROM Users WHERE email = 'ted@umich.edu'`
without from all rows will be deleted.

`UPDATE Users SET name = 'ABL' WHERE email = 'zxy@com'`

`SELECT * FROM users`

`SELECT * FROM users WHERE email = 'zxy@com'`

`Select * From Users order by name`

Program.

`import sqlite3`

(creates a new connection)

`Conn = sqlite3.connect('emaildb.sqlite')`

`Curs = Conn.cursor()`

~~Drop if already exist~~ `Curs.execute('Drop table if exists counts')` table name

~~Create~~

~~Create new table~~ `Curs.execute('Create table counts (email TEXT, count INTEGER)')` attribute type

`Email = 'xyz'`

`Curs.execute('Select count from counts where email = ?', (Email,))`

`Row = Curs.fetchone()`

If Row is None:

else

conn.commit()

sqls = "Select email, count from Counts ORDER BY count DESC LIMIT 10";
for row in cur.execute(sqls):
 print(sqls, row[0], row[1])

cur.close()

rows returned
as in form of
tuples.

Primary Key \rightarrow Main Key
Foreign Key \rightarrow Key referring to Main Key of other table.
Logical Key \rightarrow Might be looking up data using this key.

Create Table Album (

id INTEGER NOT NULL AUTOINCREMENT
UNIQUE PRIMARY KEY

JOIN \rightarrow

table1 field1
table2 field2

Select Album.Title, Artist.name from Album join
Artist on Album.artist-id = Artist.id

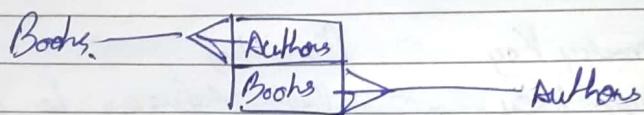
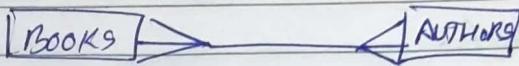
NOTE if we do not use on clause it
basically says all combinations of both tables

INSERT OR IGNORE

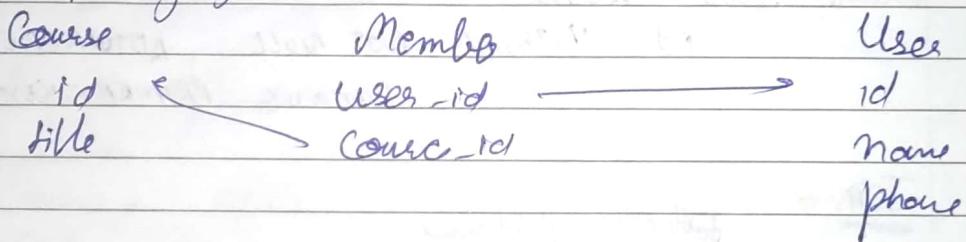
\rightarrow can be used when we have
unique attribute. It will
not add insert same attribute
twice.

MANY TO MANY.

We break the tables having MANY TO MANY using a table called "connection" table with two foreign keys.



There is usually no primary key in such connection table but we can have multiple columns or all columns as primary key.



One user many courses
many courses one course } Many to Many.

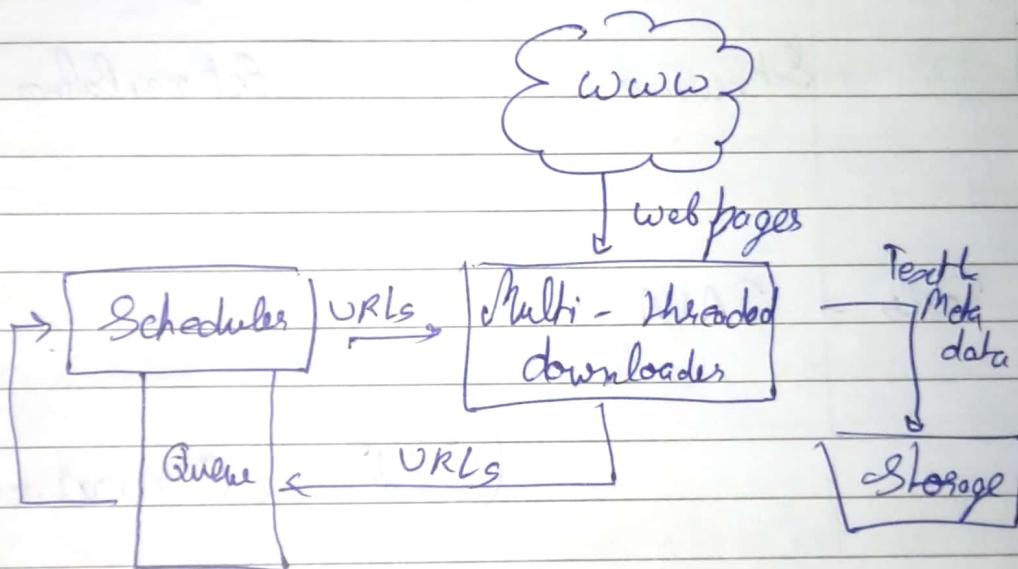
→ Retrieving Processing and Visualizing Data.

Search Engine Architecture → Web Crawling, Index Building, Searching.

Web Crawler → Computer program that browses www in methodical, automated manner. Creates a copy of all visited pages for later

processing by a search engine that will index the downloaded pages to provide fast searches.

- Retrieve a page
- look through the page for links
- Add the links to a list of "to be retrieved" sites
- Refresh



INDEX BUILDING

- Small project to crawl web store in database, page rank, visualizations.