**Extract:**
glob

## Composite functions

```python
import glob

list_csv=glob.glob('*.csv')

list_csv:['source1.csv', 'source2.csv', 'source3.csv']

list_json=glob.glob('*.json')

list_json:['source1.json', 'source3.json', 'source2.json']
```
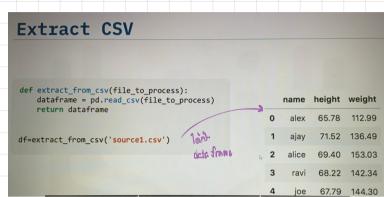
(annotations: import เก็บหลายๆไฟล์ ; output เป็น list ของ file)

Files panel:
- {} source1.json
- {} source3.json
- {} source2.json
- source1.csv
- source3.csv
- source2.csv

## Extract CSV

```python
def extract_from_csv(file_to_process):
    dataframe = pd.read_csv(file_to_process)
    return dataframe

df=extract_from_csv('source1.csv')
```

(annotation: โหลด data frame)

|   | name | height | weight |
|---|------|--------|--------|
| 0 | alex | 65.78 | 112.99 |
| 1 | ajay | 71.52 | 136.49 |
| 2 | alice | 69.40 | 153.03 |
| 3 | ravi | 68.22 | 142.34 |
| 4 | joe | 67.79 | 144.30 |

## Function Extract CSV

```python
def extract():

    # create an empty data frame to hold extracted data
    extracted_data = pd.DataFrame(columns=['name','height','weight'])

    #process all csv files
    for csvfile in glob.glob("*.csv"):
        extracted_data = extracted_data.append(extract_from_csv(csvfile),
        ignore_index=True)

    #process all json files
    for jsonfile in glob.glob("*.json"):
        extracted_data = extracted_data.append(extract_from_json(jsonfile),
        ignore_index=True)

    return extracted_data
```

(annotations: ตั้ง line ให้เป็น Header ใน dataframe ; load csv file ทีละ ; append → dataframe ; เป็น True)

|   | name | height | weight |
|---|------|--------|--------|
| 0 | roger | 65.78 | 110.99 |
| 1 | bob | 63.20 | 136.49 |
| 2 | tod | 69.40 | 190.03 |
| 3 | kate | 78.22 | 262.34 |
| 4 | moe | 66.79 | 194.30 |
| 0 | alex | 65.78 | 112.99 |
| 1 | ajay | 71.52 | 136.49 |
| 2 | alice | 69.40 | 153.03 |
| 3 | ravi | 68.22 | 142.34 |
| 4 | joe | 67.79 | 144.30 |

(annotations: ทีละ ; csv .file ; ดึงข้อมูลออกมาโดยไม่มี Header)

## Ignore index result

(annotation: ถ้าเราไม่ set ignore_index= True)

|   | height | name | weight |
|---|--------|------|--------|
| 0 | 65.78 | alex | 112.99 |
| 1 | 71.52 | ajay | 136.49 |
| 2 | 69.40 | alice | 153.03 |
| 3 | 68.22 | ravi | 142.34 |
| 4 | 67.79 | joe | 144.30 |
| 0 | 65.78 | alex | 112.99 |
| 1 | 71.52 | ajay | 136.49 |
| 2 | 69.40 | alice | 153.03 |
| 3 | 68.22 | ravi | 142.34 |

|   | index | height | name | weight |
|---|-------|--------|------|--------|
| 0 | 0 | 65.78 | alex | 112.99 |
| 1 | 1 | 71.52 | ajay | 136.49 |
| 2 | 2 | 69.40 | alice | 153.03 |
| 3 | 3 | 68.22 | ravi | 142.34 |
| 4 | 4 | 67.79 | joe | 144.30 |
| 5 | 0 | 65.78 | alex | 112.99 |
| 6 | 1 | 71.52 | ajay | 136.49 |
| 7 | 2 | 69.40 | alice | 153.03 |
| 8 | 3 | 68.22 | ravi | 142.34 |
| 9 | 4 | 67.79 | joe | 144.30 |

(annotations: ถ้าเป็น False index จะ index ตาม Original file ทำให้ซ้ำกันได้ ; index= True ก็จะให้ index จริงไปเลย เรียงตามลำดับ)

**Transform:**

## Conversion function

```python
def transform(data):
    #Convert height which is in inches to millimeter
    #Convert inches to meters and round off to two decimals(one inch is 0.0254 meters)
    data['height'] = round(data.height * 0.0254,2)

    #Convert pounds to kilograms and round off to two decimals(one pound is 0.45359237
    kilograms)
    data['weight'] = round(data.weight * 0.45359237,2)
    return data
```

→ ตั้งให้ output มี ทศนิยม 2 ตำแหน่ง

**Load:**
- save Dataframe เป็น .csv

## Load

```python
def load(targetfile,data_to_load):
    data_to_load.to_csv(targetfile)


targetfile = "transformed_data.csv"

load(targetfile,transformed_data)
```

| height | name | weight |
|---|---|---|
| 0.04 | alex | 23.25 |
| 0.05 | ajay | 28.08 |
| 0.04 | alice | 31.48 |
| 0.04 | ravi | 29.28 |
| 0.04 | joe | 23.09 |
| 0.04 | alex | 23.25 |
| 0.05 | ajay | 28.08 |
| 0.04 | alice | 31.48 |
| 0.04 | ravi | 29.28 |
| 0.04 | joe | 23.09 |
| 0.04 | roger | 22.63 |
| 0.04 | bob | 28.08 |
| 0.04 | ted | 30.1 |
| 0.05 | kate | 32.69 |
| 0.04 | max | 35.98 |
| 0.04 | jack | 25.37 |
| 0.04 | tom | 25.11 |
| 0.05 | tracy | 28.09 |
| 0.04 | john | 23.12 |
| 0.04 | jack | 25.37 |
| 0.04 | tom | 25.11 |
| 0.05 | tracy | 28.08 |
| 0.04 | john | 23.12 |
| 0.04 | jack | 25.37 |
| 0.04 | tom | 25.11 |
| 0.05 | tracy | 28.08 |
| 0.04 | john | 23.12 |

**Logging:**
- ETL เก็บ log เพื่อ มีการเกิด ตอนไหน

## Logging Entries

```python
from datetime import datetime

def log(message):
    timestamp_format = '%Y-%h-%d-%H:%M:%S'
    now = datetime.now()
    timestamp = now.strftime(timestamp_format)
    with open( "logfile.txt" , "a" ) as f:
        f.write ( timestamp + ',' + message + '\n' )
```

เป็น ตัวบอกว่า format เวลา แต่ละที่

→ จับเวลา ตอนนี้

เปิด file เพื่อ เขียน → ถ้าเคยมีก่อนจะเป็น append เข้า file

**เรียกใช้ ETL:**

## Final Call

```python
log( "ETL Job Started" )

log( "Extract phase Started" )
extracted_data  = extract()
log( "Extract phase Ended" )

log( "Transform Job Started" )
transformed_data  = transform( extracted_data )
log( "Transform Job Ended" )

log( "Load Job Started" )
load( targetfile, transformed_data )
log( "Load Job Ended" )

log( "ETL Job Ended" )
```