Department of Artificial Intelligence and Machine Learning

III B.Tech I Sem

DATA WRANGLING AND PRE- PROCESSING UNIT-1

Introduction to Data Wrangling

By

Mr.N.Siva,

Assistant Professor,

Department of AI&ML.

Data Wrangling and Preprocessing UNIT-1: Introduction to Data Wrangling what is Data wrangling, Importance of DW, Tasks of Data wrangling, Data wrangling Tools Introduction to python for Data wrangling, Python basics for Data wrangling, Handling Stouctured data: CSV, JSON, and XML formats, Data Meant to Be Read by Machines.

What is Data Wrangling? Data virangling (also known as pala Mugging) is the process of cleaning, structuring, enriching, validating and transforming raw data into a described format for better decession making and analysis. and analysis. 1. Cleaning: Handling missing & Duplicate data. cleaning is the process of handling missing, incorrect, inconsistent or duplicate data to improve data quality. import pandas as pd data = 2 ta 2 = extrag strograf Name: [John, Mary , Sam, Mary Age: [25, None, 30, None], 'Salary': ['\$5000', '\$6000', '7000', \$6000']

df=pd. DataFrame (data) = print ("original Data")

step df=df.drop_duplicate() #Remove duplicate vous. #clean step df ['Age'] = df ['Age']. fillna (df ['Age']. mean()) # fill missing ages with average of ['salary'] = of ['salary']. Str. replace ('\$', '').

Original Data cleaned Data today Salary Name Age Salary Name Age 0 John 25.0 \$5000 5000.0 0 John 25.0 1 Mary 27.5 1 Mary NaN \$6000 6000.0 2 Sam 30.0 7000 7000.0 2 Sam 30.0 3 Mary Nan \$ 6000 a structuring: Split and organize Data. Structuring means converting unorganized or semi structured data ento a well-organized format such as sows and columns. import pandas as pd # Pau unshured data data = { info: [John, 25, 5000] Mary 30, 6000) Sam, 28, 4500] 20000 stoolgub suome & of didney for the dece df = pd. Data Foame (data) print ("original Data") print (df)

split the Info column into Strutured df [['Name', 'Age, 'salary']] = df ['Info']. Str. Split (1,1, expand=me) df. drop ('Info', azis=1, inplace = Toue) print ("In structured Data") print (df). structured data original Data salary Name Age
5 John 25 John, 25, 5000 5000 Mary, 30,6000 1 mary 30 6000 4500 Sam, 28, 4500 Lamendo ash 28 -) we want to split the Info column into three Seperate Columns. Name Age Salary df['Info']. str. split (', ', expand = True) -) This splits each string in the Info column by -) expand = toue means, create seperate df. drop (Info', axis=1, inplace=true) -) This removes original Info column -) axis=1 means we are dropping a column (not a -) inplace = True means the charge is applied directly to off.

New Columns. 3. Enriching: Add adding or creating new Enriching involves existing data to provide data fields from more insights. sur , otil) doctor st import pandas as pd (36) della data = { Name's ['John', 'Mary', Sam') Age 1: [25, 30, 28] Salary': [5000, 6000, 4500] 3 des Ept. potentracon (data) # Add New Columns for Park of [Tax'] = of [Salary'] + 0.20 # Add New Columns for Age Category of ['category'] = of ['Age'] . apply (lamda X: 'young' if xc30 els 1 Adult) print (" Encoched data") print (df). Enriched data John 25 5000 1000.0 Your Sam 28 4500 900.0 Your

4 Validating: Check for Errors and Priconsistencies Validating ensures that data is accurate, consistant and meets required formats or rules. import pandas as pd data = 2 I Name: [John', mary, sam'], 'Age': [25, -30, 28], #Invalved regarive Age 'Salary': [5000, -6000, 4500]
invalid Sakon regative Salary df = pd. Data Frame (data) Data for Validation print (Data for Validation") Age Salary Mame 0 John 25 5000 point (df) Mary -30 -6000 # check for Negative solary 2 sam 28 4500 print ("In Negative salaries") Negative salaries point (df [df ['salary 1] <0]) Mary -30 -600 # Check for Negative Age point ("In Negative Age") Negative Age print (df [df ['Age'] (o)) Name Age Salary 1 mary -30 -6000 # Drop rows where Age or salary is negative de = de [cdf ['age']>=0) se (df ['salary']) print ("In cleaned data") print (df) cleaned data Mame Age Solary John 25 5000 Sam 28 4500 Age: [25, -30, 28] Salasy : [5000 - 6000, 4500 I involved Salow regularion Up= pd. Data Foome (data) grint ("cola for validation") (26) 4 misq # check for Megative solasy be surpered all tracks 100 (of 1 44 1 140) 10099 al Check for Negative 456

5- Transforming: change format or unit Transforming Phyolves converting data analysis or formats or units to meet reporting requirements. 9mpoor pandas as pd 2. Importance of Data wrangling: data = { 'Name': ['John', 'maoy') 'sam'] 2000001 Join Date': ['01/06/2023', 15/07/2022', 10/08/2023] !salary_usp': [5000, 6000, 4500]. ds= pd. Datafoame (data) # convert joinDate from string to date time fromat df [150inDate] = pd. to_datetime (df ['soinDate], format = 1-1-d -1-m -1.Y) # convert salary from USD to INR (1 USD = 83.2 INP) conversion_rate = 83.2 df['salary_INR'] = df['salary_USD'] * Conversion_rate. print ("Froms formed Date:") print (df).

Mansform	red Data:		
Name	John Date	salary_USD	Salary_INR
John	2023-06-01	5000	416000.0
Mary	2022-07-1		499200.0
0-20	202 08-	1.500	374400.0

2. Importance of Data wrangling:

1. Improves Data Phality: Raw data often contain inconsistencies, duplicates or missing values. wrangling fixes thes issues.

En: Sales dataset

customer Name	purchase Date	Amount
John John	12/05/2024	500
John	12-05-2024	500
Smith	NOTE	6500

issues:

-) Duplicate	swor .	for	John	Salas	froumas to
-) Different	date	form	nats (12/05/	2024 VS
	purcha			and L	2-05-2024)

Data wrangling fixes:

- -> Removes duplicates
- standardizes date format (4444-mm-DD)
- -) Fill or flags missing data

2 - Save Time:

Data scientists often spend To-80% of their time for proposing data (process of their time for proposing data (process of transforming saw, messy, incomplete data into a transforming saw, messy, incomplete data into a clean and structured format) for analysis, visualization or machine learning.

Instead of manully cleaning 10,000 rows in Excel
A python script using pandas can clean the
entire file in seconds.

3. Enhance Decession-Making: cleane and mun more accurate structured data leads to more accurate analysis and insights.

4. Supports Machine Learning: Machine learning algorithms require clean, formatted, and normalized data to perform well.

5. Integrates Diverse Sources! combines data from different formats, systems, or databases for a unified view.

Tasks of Data Wrangling

Data wrangling Phuolies multiple tasks convert raw, unstructured or messy that help data Porto clean and structured format suitable for analysis or machine learning.

main tasks of tata wrangling: Demorna Chromes C

- 1. Data collection
- 2. Data cleaning
- 3. Data Transformation
- 4. Data Integration
- 5 Data Reduction
- 6. Data validation
- 7 Data Enrichment

all Encode Acade .

1. Data collection:

Gathering data from various sources like csv files, databases, APIs or web scrapping.

Emport pands as pd [] df = pd. read_csv ("employees.csv")

Load data from

df ['Agail - Filma (df

changing the format.

2. Data cleaning: prisoner stoll zo extent

Fixing or semoving incorrect, corrupted, missing or duplicate data.

Common cleaning Tasks:

-> Handling missing values

-) fixing wrong data types

-) Removing duplicates

-) Correcting Types

df. drop_duplicates (inplace = True) df ['Age'] - filna (df ['Age'] - mean().

Enplace = True

3. Data Transformation:

changing the format, structure or values of data to make it more usable

df ['salary_INR'] = df ['salary_usD'] * 83.2 # convert currenty df['Gender'] = df['Gender'].map(q'male:1,

'female':05) # Encode gender.

4. Data Integration)

combining data from multiple sources into

a single dataset

exi df1= pd. read_csv ('employee. esv')

df2= pd. read_csv ('department.csv')

dfcombined = pd. merge (df1, df2; on = Employee

5 Data Reduction

Reducing the Volume of data by removing Proclevant or redundant information.

df = df [['EmpId', 'Name', 'Salarry']]

Keep only Proportant

columns

6 Data Validation:

Ensures that the data meets Specific quality rules (eg. no negative ages or salaries)

df = df [(df ['Age'] >0) & (df ['salary'] >)]

7) Data Enrichment: Enriching data by adding new relevant Information (from external or derived Sources) df ['experience'] = 2025 - df ['Joining Year'] # calculation of Experience

primarios po estato do melos sat grandas

[C'enfeet, 'mount itempes] 3 hough

Ethan short stone white the some

[(CE16000)] 49) 4 (04 (1904)] 37) 19 49

(some is tops overgen as . 13) block

agetherent street to

Introduction to python for Data wrangling.

Python is a high-level, interpreted and general purpose programming language known for its simplicity, readability and strong community support. It is one of the most used languages in data science, machine learning and data wrongling

why use python for Data wrangling:

- -> Easy to learn and use
- -> Rich set of libraries (pandas, Numpy etc)
- > Supports multiple data formats (csv, Excel, JSON, XML)
 - -) Scalable and works with large datasets.
 - -) Integrates well with visualization tools and databases.

Key Libraries for Data Wrangling in Python. purpose Library Data manipulation & Analysis Pandas Numarical operations & Numpy Lapournal Read write Encel Files openpy XI, XIrd Brush Work with josh data Ison Read/write CSV files esv Regular Expressions (for text

in start nother deren the something

For skides of

Basic woonkflow of Data wrangling in 1. import data 2. inspect and understand data 3. clean data (remove missing values, duplicates 4. Transform data (rename columns, change types 5 Export cleaned data import pandas as pd # Stepl: Load data from csv file df = pd. read_csv ('employees.csv') # Step 2: view first 5 xows (print (df. head ()) # step3: Check for missing values point (df. isnull () · sum()) # step 4: Doop soms with missing values df-clean = df. dropma() # Steps: Rename columns df_clean = rename (columns = { EmpNorme! : Employee Norme) inplace=True) # Step 6: save cleaned data df-clean. to_csv ('cleaned_employees.csv', index = false)

Data wrangling Tasks: Common Python Function/Method Fask Load csv pd. read_csv() Drop missing values df. dropaial) Fill missing values df. fill na (value) df. drop_duplicates() Remove duplicates df [df ['Age'] > 30] filter rows df. rename (columns={ old! Rename columns New 3) df['Age!] = df['Age']. astype
(int) change data types String operations df ['Name']. Str. Upper() 10/7/25: new \$3-03 presenting 441m smae daeg : H days # # Sleps: Rename columns in clean a seriams (columns = & Emphanie : Complete ma # Shap beneat some dates. L'va a septema harred de vas est employee

```
Import pandas as pd to 2018 ad worther # step1: Read the original data
 df = pd. read_csv ('c:/user/ ... /employees.csv')
# Stepa. clean - Remove sows with missing salary
  df clean = df [df [salasy'] . not null ()]. copy ()
# step 3: clean - Rename column
 df-clean = rename (columns = { EmpName :
Employee_Name &
                            inplace = 9 our
# step4: brop duplicates.
   df_clean.drop_duplicate (inplace=True)
 # steps: some cleaned data
               Star Stains (test) Tate
df-cleared.
 df-clean.to-csv ('c:/ user/-- /cleaned_emp.csv',
                                   index = False)
 Print ("cleaned data sould successfully").
```

Python basics por Data Wrangling: 1009mil # stept: Pead The

1. Variables and Data types:

A variable its a name that refers to a value Stored in memory.

Data type define the type of value a variable holds, such as numbers, text, or logical values.

Common python data types:

int- Integer numbers (eg:10)

floot - Decimal numbers (eg: 10.5) +2 #

Str- String (text) data (eg: "Data") bool-Boolean (Toue or False)

Example:

hame = Stra" #str age = 25 #int height=5.7 # float Hold #

is student = True # bool

- GURLAR COCKE

2. Lists:

A list is a mutable, ordered collection

(states) the control

of values.

List allow you to store multiple values in a Single vorviable. You can change, add or remove stems. Dust bis northebros will a appell willed

fouits=["apple", banana", "mango"] print (fruits [1]) # output: banana fourts. append ("grape") # Add a new item

3. Dictionaries

A dictionary is an unordered collection of key-value pairs. Each key maps to a Specific value. It is useful for Stouctured data Student = { (like JSON)

"Name": "Rajy",

"age": 25

& "branch": "AIML" output: point (student [branch"]) # AIML

4. Control How (it, for while) Control flow statements manage the execution path of a program. if: executes code based on condition for: repeats over îtems while: loops while condition is true. 书作 Score = 80 the topost if score > 50° print ("Pass") Dictionocoies # for loop for i in range (3): # Hello o print ("Hello", i) Hello 1 Hello 2 # while loop Count =0 Justin E while count <3; print ("count,", count) # count o count 1 Land Count+=1 Count 2

A functions:

A function is a reasoble block of code that performs a Specific task. we can define own own functions with "def", and call them using their name.

Ex: def greet (name):

return "Hello," + name

print (greet ("Raj")) # output: Hello, Raj

6. file Handling:

File handling refers to reading from and writing to Files. Files can be opened using writing to Files. Files can be opened using open () with modes like r (read), w (write) and a (append).

with open ("demo-txt","w") as f: f. write ("welcome to python")

8) reading csr using pointed

Read from a file
with open ("demo.tat", "12") of f:
print (foread)

7 Imposting Wibraries: Importing allows you to use external python abrances. Ose "Propost" to boing in built in or third party modules like pandas, numpy, etc. Example: import pandas as pd Impost numpy as no impost json. bright (duest (, 201,)) it outlook 8) Roading csv using pandos:

Brathmoth states

Annum CSV (Comma Seperated Values) Ps a file format used to store tabular data. Ex: Emport pandas as pd df = pd. grad_csv ("Employeeg.csv") print (df-head ())

still to more from the

- 9- 10 ("12", " - 1 x 3 - 0 mo)

9) working with JSON Data: JSON (Java Smift Object Motation) is a format for stooing and exchanging data. Use the ison module to load or write Json data Emport json with open ("data. json") as f: data = Sson. load (f) for emp in data ["employees"]: print (emp ["name"], emp ["Jepartment"]) 10) Regular Expression (re modele) Regular expressions are used to match the

text patterns. Python's re module allows pattern marching and text replacement

Ex: import re text = "phone: 9876543210" match = re- search (r"|d {104", text) if match: print ("found numberi", match. group())

· sta host dees U-Numpy Bosics: Numpy is a python library for numarical and array operations. import numpy as Py an = np. amay [1, 2, 3, 4, 5] print (arr. mean ()) # output: 3.0 Etrontripl") quis (Comon) quis) +ming (delpore ex) roicenges roluged (or set determ of been are industry willings ext patterns. Pythions are module allows pasterns to harding and text replacement please: design = despession Anot (For 3 to 1) of diese one of down (Colours of total and purel purels) total Scanned with OKEN Scanner

Handling structured data: csv, ison, xmi CSV is a plain text format where each line represents a data record. Each field is seperated by a comma. * why use csv? -> Easy to read and write -> Common format for tabular data. -> Supported by Excel and databases. Nome: Emphane ditipe: abject employees. CSV EmpID, EmpName, Depostment, Salary 101, Alice, HR, 50000 Jean Head to 23 Most 102, Bob, IT, 60000 103, charlie, Finance, 55000 import pandas as pd # Load CSV df=pd. read_csv ('employees.csv') # Display the data frame point (df) # Access specific columns Print ("In Employee Names:") Print (df [EmpName]).

paretored structured data: csv output: Department Salary Emphlome EmpID HR 50000 Alice 0 101 IT 60000 Bob 1 102 finance 55000 charlie 2 103 Employee Names: o Alice 1 Bob was a series of the seri 2 chartie Name: EmpName, dtype: object 2. JSON (Java Script Object Notation) JSON is a lightweight data-interchange format It is structured like a python dictionary and supports nesting. Sample Lata. json ? "employees":[& "name": "Alice", "department": "HP"3, ¿"name": "Bob", deportment ": "IT"},] { "name": "charlie", "department": "Finance"

Import json #Load Jeon with open ('data-json') as f: data = json-load (f) # Access employee data point (" Employees Details:") for emp in data ['employees']; print (emp['name'], emp['department']) (name) Alice Chame (templostre> Employee Detalls: Alice HR Bob IT (name > Bob (mane) charlie Finance (depositions) from a Lldepositions)

3. XML (eXtensible Markup Language) XML is a markup language that stores data in a structured, hierarchical format using user. defined tags. Sample XML: employees. XMI (company) (employee) ams F'sman'] gms) thing {name} Alile (Iname) (department) HP (| depertment) (lemployue) (employee) (name) Bob (I name) (department) IT (Idepartment) (employee> (employee) < name> Chaolie < Iname> (department) Finance (Idepartment) < lemployee> (company)

import Xml. etree. Element Tree as ET # parse XML tree = ET. parse ('employee.s-xml') root = tree. getroot () # Access employee data print ("Employee Data from XML") for emp in soot. findall ('employee'): name = emp. find ('name'). text dept = emp. find ('department'). text print (name, dept) output: Employee Data from XML: Alice HR Bob IT Charlie Finance

Department of Artificial Intelligence and Machine Learning

III B.Tech I Sem

DATA WRANGLING AND PRE- PROCESSING <u>UNIT-2</u>

Working with Excel Files, PDFs, and Databases

By

Mr.N.Siva,

Assistant Professor,

Department of AI&ML.

working with Excel Files, PDFs and Databases Installing Python Packages for Data Wrangling, Parsing Excel Files, 2000 of remove Programatic Approaches to PDF parsing Converting PDF to Text (pdfminer) Acquiring and Storing Data Introduction to Data bases for Data woungling Relational Databases: MysQL and PostgreSQL. Non Relational Databases: NaSQL and Alternative

Data Storage.

plp install paniles openpyal which parameters

press Enter

Installing Python Packages for Data wrong hing Install enternal libraries that help handle Excel, PDF, and database files.

Common Packages: 2017 200

* pandas - For structured data (Excel, csv)

* openpyxl, xlord - For Encel format

+ pafminer. Six, PYPDF2 - for PDF parsing

* Sqlalchemy, Sqlite3 - For SQL databases

* pymongo - for Mongo DB (No SQ)

Installation steps:

Stepl: Open Terminal or Command prompt press windows + R and type and and press Enter

Stepa: Run installation Command

pro install pandas openpyxl wird pafminer-six Sqlalchemy pymongo

Steps: wait for installation to Complete then

successfully installed pandas openal aird.....

step4: Venify installation
open a python shell IDLE and press control + N
open a python shell IDLE and press control + N
type below packages.

import pandas as pd
import openpayal
import alod
import Sqlalchemy
import sqlalchemy
import pymongo
import pymongo
from pdfminer. high_level import extract_text
print ("All packages imported Successfully")

output:

All packages emported Successfully.

Note: if you are using Python 3.13, Pip is not recognized, try: Py -3.13 -m pip install openpyXI

why use python for Excel:

2 Parsing Encel Files: Excel is a widely used Spread sheet farmal to store structured data in rows and columny it comes in two main file types. 1. . x/8x: Excel Workbook (New format) -) Introduced in Excel 2007 -) XML based (open XML format) -) support over a million rows. -) Commonly used and recommended -> can be read using openpyx 2. 2/s: Excel 97-2003 workbook (old format) -> Binary file format -) Limited to 65,536 rows -) can be read using alrd why use python for Excel: for automation. Python makes it easier to

-) Read Excel data in bulk

-) Analise Filter and Soft

Scanned with OKEN Scanner

-) Add new Columny -) Save results back to new Excel files -) Automate repetitive Excel processing tosks Required Python Libraries: 10% brigage Durpose Library 1 gandas Read and manipulate data openpyxl Read | won'te . 2/5% Excel xlrd Read -xlod .xls Excel How to install: in Command prompt Py -3.13 -m pip install pandas opengul openpyxl aird Read . xls2 File Using openpyxl import pandas as pd # Road new format Excel files (. xlsx) df = pd. read_excel ('sample. 2/52', engine = print (df)

* Pd. read - excel(): reads pandas Datatoame + engine = 'openpyxl': tells pandas to use the .xlsx files. north q horner openpyn ubary for Read .xls Files using xlrd import pandas as pd # Read old format Excel file (.xls) dj=pd. read_excel ('sample. als', engine=x print (df) og bromme ni : Motoni ot evol Common operations: 1) Filter Rows - Get Students with Score 78 high_score=df[df['score'] 780] ~~~ print (high_score) 2) Sort Dota by Score (Desending) Sosted_df = df. sost_values (by = 'score', ascending = talk Print (sorted_df)

paragrammatic Approaches tomming comments of df ['Grade'] = df ['score']. apply (lambda s: A if S>=90 'B' if S>=80 else 'c') print (df) 4. Save the filtered Data to a New Excel File; high scores. to excel (high score - x/sx', index: pateminer six - Best der to Data . Exist = 16989 10 579989 (-Name Score of Essamulgable Alice 88 1. Install Regulated Libraries Bab 75 pip forstall parfuntner - six Charolie David 67.000967 6080 2709 600009 5 * Extract At Text from PDF from paternines six high level impost Aland of text from PDI text = extract text ("surple pot)

Parogrammatic Approaches to PDF parising in Porting parsing PDFs programatically allows automates extraction of data like: -) Text Content (3th) troop -) Tables evic the followed take to a mela data Python offers powerful libraries for PDF parsing -> pafminerosix -> Best for text extraction -> PYPDF2 or pypdf -> Basic text + metadok -) pafplumber -> for table and structured text data 1- install Required Libraries: PYPDF2 pdfplumber pip install pdfminer. six 2. parsing PDFs using pafminer. six (Best for plain H) * Extract All Text from PDF from pafminer-sixhigh-level impost extract-test #Read all text from PDF text = extract_text ("sample. paf") print (text)

Save Extracted Text to tile

with open ('sample tat', 'w', encoding='utf-8')

as f:

f.write(text)

3 Parsing PDFs Using PyPDF2 (or Pypdf)

Biasix Text Extraction 27199 100901

PyPDF is a python library that allows you to read, extract, split and merge PDF files. It's useful for basic operations such as extracting the text content from pages or merge multiple PDFs.

Use cases: " sold lotor") thing

- -> Extracting text from PDFs
- -> Reading pages Individually
- -) Merging or splitting PDF documents
- -> Getting PDF metadata

Installing PyPDF2 pip install PyPDF2 pip Enstall PyPDF2 Spython exe on pip install outproposed from PDF: Let's say you have a PDF file named "sampling with a pages. Emport PyPDF2 and that the # Open PDF File 1199 3 7999 with open ('sample pdf', 'rb') as file: read = PyPDF2. PdfReader(file) # Print total pages point ("Total Pages:", len (reader. pages)) # Extract text foom each page for i, page in enumerate (reader. pages): print (f" |n--- Page { i+1}---) point (page. extract_text())

Sample Output If your sample pdf has the following text: · page 1: "Welcome to Data Wrangling" - page 2: "This PDF demonstrates PYPDF2" output is most test toosted appeal Total pages: 2 insent papplamber --- page 1 ---. welcome to Data wrangling poll-polled / noviz force from step 169 169 ---page2--This PDF demonstrates PDB PYPDF2 4 (Pofflumber): Parsing PDFs Using pofflumber:

4 (Pdf plumber): Parsing PDFs Using pdf plumber:

pdf plumber & a python library designed

for extracting text, tables, and metadata

from PDF files. It works especially well

with PDFs that contain complex layouts or

tables, making it a popular choice when

the boilt in methods (such as those from PyPDF2)

may not extract information correctly.

* Installing pafflumber: Before ounning Code, you need to install pdfplumber. Open your command prompt or [-3.13 -m -pip install pdfplumber]

pip install pdfplumber Basic Usage: Entract Tent from a PDF: import polyplumber # specify the path to your PDF file pdf_path = r"c: | user (sivan | Desktop) sampl. pd # open the PDF Using polyplumber with poffplumber. open (pdf-path) as pdf: # print to tal number of pages print ("Total pages:", len (pdf. pages)) # 2 terate through each page and extract & for i, page in enumerati (pdf. page text=page.extract_text() print (f" |n -- page {i+igpoint (text) (stages most sent ap does) charlom willied

output: sample . pdf page 1 content "welcome to Pdf.plumber" pag 2 content
"Welcome to Data Wrangling class" output: solder of Adomet status. Total pages: 2 -- page 1 ---"welcome to postplumber. -- page 2 welcome to Data Wrangling class 23/73 19 Mar is was able to de Total Red ITA LO algority del 24 contract 189 1012 2 , 90 physics with some swing of and

Acquiring and storing Data:

Data acquistion and storage is the foundation of any data wangling workflow. It refers to the process of "gathering data toom various sources and storing it in Surtable formats for further processing and

I sources of Data Acquistion:

you can aquire data from the following

Sources:

Tools methods Sources Files CSV, Excel, JSON, XML, Tes Databases Mysel, Postgresel, sality Web API REST APIS using requests, ison Nob Scraping HTML data using Beautifulson Selenium PDFs Polymines, 174PDF2, polyfumber dow storage Google Drive, AWS, 53, Firebo

2 stooning formats: After acquiring data, It can be stored in: * structured formats: CSV, Excel, Databases * semi_ stouctured fromats: JSON, XML * Unstructured formats: Text, [mages, PDF 1. Reading from Pa CSV File and Saving to Excel: emport pandas as poder to # Road CSV file df = pd. read_ssv ("data. csv") # Display Data tob Jegmub most I tom print ("csv Data: In", df. head ()) # store to Excel of . to_excel ("output data.xlsx", index= false print (" bata stored in Excel successfully")

& Fetching Data from a REST API and Storing to JSON ! Emport requests import json # API Endpoint (Example) arl = "https://isonplaceholder. typicode. / posts" # send GET request response = requests. get (url) data = response. json() # print first 2 records print (json. dumps (data [:2], index # save data to a JSON tile with open ("api_data. sson", 'w') as f: ison · dump (data, f, indent = 4) Point ("API data Stored to JSON) sicressfully

3. Reading from an Excel tile and Saving to CSV : import pandas as pd # Read Excel file df = pd - read_excel(" "nput_data . x/sx") # Display first few rows print (df. head()) # save to csv df. to_csv ("converted_data.csv", index= Falx) point ("Excel data converted to CSV successfully") + csv to Json import pandas as pd df=pd. read_csv ("data.csv") # convert to Json and save df. to-ison ("converted_data.ison", orient = "records", indent=4) point ("csv converted to Json successfully") 5) Excel to CSV Proport pandas as pd # Road Excel df = pd. read_ excel (" data - x/s x") # convert to csv df. to_csv ("converted_data.csv", index=for point ("Excel converted to CSV Successfully) 125 of 2450 10 6) JSON to CSV: (STONE) VERSON JE import pandas as pd #Read JSON df=pd.read_ison("data.ison") # convert to csv df. to_csv("converted_from_ison.csv", index = Falsy point (" Json converted to CSV successfully

Mars 2011/2 (1036 37 potrovinos vas") freis 3

7) JSON to Excel impost- pandas as po df = pd. read_ ison ("data. ison") df. to excel ("converted_from ison. xIsx", index = false) print (" JSON converted to Excel Successfully") with Xunstanna Xunx and Xuns gord spring of god to see ad about : 3777003 60004007 Quering with BOLL

Introduction to Databases for Data Warangeong: -) what is Data wrangling; Data wrangeing is the process of: * collecting Lata from various sources * collecting and formatting it

* cleaning, transforming and formatting it

* preparing it for analysis or machine

Leaning they use Databases in Data wrangling: Databases helps manage large and structure data sets in a reliable, Scalable and efficient Benifits: why it Helps in DW Features Tables with yours 4 columns Structured Data make wrangling east we can filter, join, group Quering with sol and sort data easily

millions of your Efficient Storage efficiently control who can view or Security & Access edit data connect with Excel, APIS, Integration Scripts and more. Types of Databases: Example Tools Description Type Data in tables Myser, Relational (SQL) Postgre SOL, with fixed Schema Solite Florible, JSON-like MongoDB, Non-Relational (No SQL) fireball Python: Tools Commonly Used with purpose Tool Library Connect to solite databases Sqlite3 Connect to Mysel my sal-connector - python Connect to Postgreson psycopg2 unified DRM fordatabases Sal al chemy

work all C. JONE CHOLD pymongo pandas read_sq1() - Read query results ag DataFramos Using Socite for Data Wrangling: Python Code to Creete, Insert and Read Table producted the in import sqlite3 import pandas as pd # step1: connect to database Correcte file if not conn = squite3. connect ('student.db') cursor = conn. cursor() # Step 2: Create table cursor. execute ("1 CREATE TABLE IF NOT EXISTS Students id INTEGER PRIMAR KEY, name TEXT, grade INTEGER restatored "1) s leasterne

carsor. execute many ("INSEKT INTO Students (Ed, name, grade) VALUES (2,2,2)", [(1, 'AUG', 85), (2, Bob', 90), (3, 'char lie', 78)]) someretel mantrel borness Worken? conn. commit() # Step 4: Road data Using pandas If = pd. read_sq1 ("SELECT * FROM Students", conn) point ("student Data"). point (df) is side connectose() student Data name grade desamoles side Alice 3 charlie 78

Mysal Version:

Desetup Mysol on your Machine:

-) Download from: https://dev.mysql.com/download

-) Install Myral server & workbeach Coptional &

-) During Setup:

* set user name : root

* password: * ***

@ Install Required Python Libraries:

pip install mysql-connector-python pandas

openpyxl

-) mysql-connector-python: Mysal Connector

-> pandas: Data manipulation

-) openpyxl: Export to Excel

3) Create Database and Table in mysal!

use mysqL workbench or terminal:

CREATE DATABASE Customerds;

USE customerab;

CREATE TABLE customers (
id INT AUTOLINCREMENT PRIMARY KEY,
name VARCHAR (100),
email VARCHAR (100),
status VARCHAR (20),
Purchase amount FLOAT
);

Insert Sample Data (with duplicates)

INSERT INTO customers (name, email, status, Poourose

Alice, 'alice Rexample.com', lactive, 250.50),

('Alice', 'alice@example.com', lactive', 250.50),

('Bob', Ibob@ermple.com', "nactive', 0),

('Alice', lalice@example.com', lactive', 250.50),

('Eve', 'Eve@example.com', 'active', 120.75),

('charlie', 'charlie @example.com', 'active', 300.00);

if = pt. med_ent (quesy, com)

5. python script (mysal) - wrangling mysal.p. import mysal, connector 1200 200 Import pandas as pd 3 am sav #step1: connect to pat Mysal conn = mysql. connector. connect (host = "local host", user = "root", and significant password = "your-password" database = " customerdb" # step2: clean data using sQL query = 1111 1 SELECT DISTINCT name, email, Status, purchase_amount FROM customers WHERE constonners stoatus = 'active'; 11 11 11

df = pd. read_sal (query, conn)

Step3: Export to Excel df. to_excel ("cleaned_ customers_mysql.xlsx" index = false) paint ("cleaned customer data exported to 'cleaned_customers_mysql-xlsx'.") all pandos openpaxi conh. close()

Salite Version

this code:

1) No Installation Needed

Solite comes with Python (squite's module is built-in)

Dinstall Required Python Libraries:
pip install pandas openpyxl

3 create Database and Table in SQLite
create a file called setup-squite. Py and run

import squites

conn = squites. connect ("customerdb. squite")

cur = conn. cursor()

Create table

cur-execute ("""

CREATE TABLE IS

CREATE TABLE IF NOT EXISTS CUSTOMERS

ALL INTEGER PRIMARY KEY AUTOINCREMENT

Name TEXT,

eamail TEXT,

Status TEXT,

purchase amount FLOAT

```
# Insert data
   cur executemany (" "
INSERT INTO customers (name, email, Status, Purchase
                         amout) NATREZ (5'5'5'S)
 come a septite a comment ("commonate spire"
 ('Alice', 'alice Dexample. com', 'active', 250.50)
 (1806, 1806@example.com', 'inactive', 0),
 ('Alice', lalice @ example.com', 'active', 250.50)
 ('Eve', 'eve cerample com', 'autive', 120.75),
 ( charlie, charlie example com, active, 300.00)
conn. commit ()
conn. close()
        est being a stability business " ) frees
```

1 Python script (salite) - wrangle_sqlite.py Propost pandas as pd

#step1: connect to salite

conn = squite3. connect ("customordb. squite")

Stepa: clean data using sol query = "11" mos sigmons surg! dod!

SELECT DISTINCT name, email, status, porchase_amount

From customers

WHERE Status = 'active';

df = pd. read_sql (query, conn)

Step3: Export to Eacel

df. to_excel ("cleaned_customers_squite.x/sx",

Index = false print (" cleaned customerdata exported to

cleaned customers south x15x.

Conn. close()

() 920/ 5 . pho)

No SQL Databases: MongoDB & Fixabase in

Datawrangling: Non-reational (Nosel) data bases store data in formats like documents (Ison) or Key-Value pairs, offering flexible schema idea for modern apps. Popular NoSal systems include -> Mongo DB - Document_based data base storing JSON-like structures

-> Firebase (Firestore) - A cloud-based Nosal database by Google, optimized for real time data Why use Mosal for Data Wrongling:

510 AN +A 10 HE

["zonotous"] db = no" pollo

of the cheest ["Customeral 15]

^{*} Schema flexibility

⁴ Easy integration with apps

It Greate for Semi Structured data (eg. user profiles , Sensor logs etc)

Setup and Installation!

MongoDB: (Local or Altascloud)

1. Install MongoDB:

https://www.mongodb.com/try/download/commu

. Install Required libraries

pip install pymongo pandas openpyx1

3. Start mongoDB locally or create a free chuse, https://www.mongodb.com/cloud/atlas.

Python script: Insert Unstructured Data into Mongolls

"insert_unstructured_mongo.py"

from pymongo import Mongollient

import pandas as pd

Step1: connect to Mongolls

elient= Mongollient ("mongolb://localhost:2")

or Atlas OE

collection=db["customerdb"]

collection=db["customerdb"]

Hinsest unstructured customer data collection. Prisert_many ([{"name": "Alice", "email": "alice @ enample.com", "status": "achive", "purchased amount": 250.503 S'name : Bob", "status": "inactive", "notes": "missing email and purchasely" # missing email: " Status: achiveenail: 4 Status: " -> tags": ["loyal", "frequent"], "purchen amount": 120.79 > "address": 5"airy": "Ny", 21p": "10001"33,) "purchase history":[100,200, 195]? f Stept: filter 20072 # love only delive with a volle purchase print (" Unsmutured date "Inserted Into mongood") SVHOPS" == ["Zutate"]] == beroting To (1) Horston ! "- Houses grato reg!" (46)

Connect, filter and Export to in python from pymongo 9mport Mongo Client Proport pandas as pd from pandas import ison_normalize # stepl: connect to MongoDB client = Mongo Client ("mongodb: 1) localhost: 2701) db = client [" customerdb"] collection = db["customers"] #step2: Retrive all documents docs = list (collection . find ()) # step3: Normalize (flatten) the documents df = Ison _normalize (docs) # Step4: Filter rows # Keep only active users with a valid purchase and df_filtered = df[df["status"] == "active") & (df["porchase_amount"]. notnoll())

clear up cotumns # Steps: Export to Excel 25_ filtered. to excel ("filtered_customery-21sa", Prodex = false) print (" Filtored data exported to 'filtery customers. wish") output Tribell Pother Ubeanes: email Status purchase am tage alice@ex.com active 250.50 N/A Alice Evelex.com active 120.75 ('loyal' frequent import frebase dans from Anchor, admin impost coedentials, forestor orador 2 soutours : 1900 of cred = credentles cestificate ("polis to gow) firstos Ley Jan 1 Peplace (loo) 79- 250 Hair o mindos 20 darif Chaill . Noteria = db

Firebase (cloud store) 1. Goto Firebase Console 2. Create a new project -) Add firestore data 3- Download your service account JSON: * project settings >> Service Accounts > Generate new private ke La Enstall python Libraries: firebases-admin pandas Leges svides mas 2390 de Openpyx Python Connection & Export Example & dw-firebase_to_excel-py import firebase_Admin from firebase_admin impost credentials, firestill Import pandas as pd # stepl: Initialize firebase cred = credentials. Certificate ("path to your firebox

firebase_admin_ initialide_app (cred)

db = firestore · Client ()

key ison") # peplat

Step 2: Read data from tirestore docs = db. collection ("customers") . stream () Vaw_data = [] for doc in docs: data = doc. to_dict() You data append (data) dif = pd. Data Frame (raw_data) # step3: filter and dedupticate df = df [df ['status'] = = 'active'] df = df. drop duplicates (subset = ["name" "email"]) # step4: Export to Excel df. to_excel ("cleaned_customers_firebose.xlsx", index = false) Sirebake cleaned data exported to "cleaned customers firebase x15x") steplistetup Firebase project

1. Go to https:// console.firebase.google.com/

2. Click "Add project" -> follow the prompts

3. Enable firestore Database:

x In the left menu: Build > Firestore

se choose "start in test mode"

* click create database.

Step 2: Generate Firebase Admin SDK Key

1- Goto project Settings -> Service Accounts

2. Account clicks "Generate New Private Key"

3 This downloads a . ison file >

· Save Pt (eg. firebase-ky.

4. This file will be used in your python six

to connect.

D

Create and Instat Sample Data sample data - Py import firebase admin from firebase_admin import credentials, firestore # step1: Initialize firebase cred = credentials. Certificate ("firebase-key.json") # Replace with your path firebase_admin.initialize_app (cred) db = firestore - client() Stepa: Create Unstructured customer data customers = status: "active" { "name": "Alice, 'email': pyrchase_amount: 250 3

step 3: Insert into firestore put prop Dec. 1 for customer in customers. db. collection ('customers'). add (uistomer) print ("Sample customer data added to firebase firestore") (" at out manded ") startitude (" factor to) = 1500 If Replace of the soul of the frebuse admin - initialize app (coul) () tools noteson chente) stops create unstructural customer state - Banchana Sydno" sytota " Lighas" " Della " " "