

Experiment-1

Date: 30/08/24

AIM

Understanding Data, what is data, where to find data, data wrangling, data clean up basics - formatting, outliers, duplicates, normalizing and standardizing data.

PROCEDURE

Step-1: Study about data and its importance

Step-2: Understanding about data wrangling procedure

Step-3: Understanding the importance of normalizing and standardizing the data

SOURCE CODE

```
import pandas as pd
import numpy as np
# generate data
data = [
    {"id": [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 3],
     "date": ["2023-01-01", "2023-02-15", "2023-01-01", "2023-03-10",
              "2023-04-25", "2023-05-03", "2025-01-01", "2025-07-20",
              "2023-09-01", "15/08/2023", "2023-07-01"]},
```

```

"ag": [25, 30, 25, 22, 40, 33, 150, 29, 27, 35, 28],
"income": [50000, 60000, 50000, 45000, 12000, 50000, 60000,
            58000, 72000, 50000, 40000],
"score": [55, 90, 85, 75, 95, 86, 85, 88, 93, 85, 45]
}

```

Convert to data frame

```
df = pd.DataFrame(data)
```

Save to excel file

```
file_path = "Sample_dataSet.xlsx"
```

```
df.to_excel(file_path, index=False)
```

2) # Load Sample data set

```
file_path = "Sample_dataSet.xlsx"
```

```
df = pd.read_excel(file_path)
```

display the original data set.

```
print("original data set:")
```

```
print(df)
```

3) # i) formatting

Convert data column to date time format

```
df[["date"]] = pd.to_datetime(df[["date"]], errors='coerce')
```

handle non convertible dates

replace with not or a specific date

```
df[["date"]].fillna(pd.Timestamp("2023-08-15"), inplace=True)
```

```
print("formatted data (column)")
```

```
print(df[["date"]])
```

④) # 2. handle outliers

detect outliers using Z Score for age column

$$\text{df}["\text{age} - \text{Z Score}"] = (\text{df}["\text{age}"] - \text{df}["\text{age}"].\text{mean}()), \text{df}["\text{age}"].\text{std}()$$

outliers = df[np.abs(df["age - Z Score"]) >= 2]

print("In outliers detected")

print(outliers)

⑤) # remove outliers

~~df = df[np.abs(df["age - Z Score"]) <= 2]~~

drop the particular helper column

~~df.drop(columns = ["age - Z Score"], inplace = True)~~

⑥) # duplicate

duplicates = df[df.duplicated()]

print("In duplicates detected")

print(duplicates)

⑦) df["norm-income"] =

~~(df["income"] - (df["income"].min() - df["income"].max()) / df["income"].max())~~

min())

⑧) df["Std-Score"] =

~~(df["Score"] - df["Score"].mean()) / df["Score"].std()~~

Print("In cleaned dataset: ")

OUTPUT

print(df)

9) cleaned_file_path = "Cleaned_Sample_dataset.xlsx"
 df.to_excel(cleaned_file_path, index=False)

OUTPUT:

	Id	data	age	Income	Score
0	1	2023-01-01	25	50000	85
1	2	2023-02-15	30	60000	90
2	3	2023-01-01	25	50000	85
3	4	2023-03-10	22	45000	78
4	5	2023-04-25	40	120000	85
5	6	2023-05-30	33	50000	80
6	7	2023-03-01	150	60000	85
7	8	2023-09-20	29	58000	88
8	9	2023-09-03	27	72000	91
9	10	15/08/2023	35	50000	85
10	3	2023-01-01	25	40000	75

Cleaned dataset

age	Income	Score	Name_Income	Std_Score
-----	--------	-------	-------------	-----------

Cleaned dataset		Date	Age	Income	Score	Name	Income	Std_Score
1	2023-01-01	25	50000	85	0.1250		0.219270	
2	2023-02-15	30	60000	90	0.2500		0.572932	
3	2023-01-01	25	50000	85	0.1250		0.219270	
4	2023-03-10	22	45000	75	0.0625		-0.498053	
5	2023-04-25	40	120000	95	1.0000		0.926593	
6	2023-05-30	33	50000	80	0.1250		-0.134391	
7	=	-	-	-	-		-	-
8	2023-07-20	29	58000	88	0.2250		0.431467	
9	2023-09-01	27	72000	92	0.4000		0.643664	
10	2023-08-15	35	50000	85	0.1250		0.219270	
11	2023-01-01	25	40000	45	0.0000		-2.610022	

Experiment-2

Date: 10/02/24

AIM

Develop the python script to parse the pdf files using pdfminer.

PROCEDURE

Step-1: Set up PDFMiner using !pip install pdfminer.six.

Step-2: Use extract_text method found in pdfminer.high_level to extract text from the PDF file

Step-3: Tokenize the text file using NLTK.tokenize RegexpTokenizer

Step-4: Perform operations such as getting frequency distributions of the words, getting words more than some length etc.

Step-5: Use method such as collocations or collocation_list to get most frequently sequence of words occurring in the text

SOURCE CODE

```
from nltk.tokenize import RegexpTokenizer  
from pdfminer.high_level import extract_text  
# Extract the text from pdf file  
text = extract_text("2010.00462.pdf")
```

Data Wrangling and Visualization - Python/R Programming/Power BI

```
# Create an instance of tokenizer & using NLTK Reg  
tokenizer = RegexpTokenizer(r'\w+')
```

Tokenizes the text read from PDF

```
tokens = tokenizer.tokenize(text)
```

find frequency distribution

```
fd = freqDist(tokens)
```

find words whose lengths is greater than 5 and frequency greater than 20.

```
long_frequent_words = [words for words in tokens  
if len(words) > 5 and fd(words) > 20]
```

- long-frequent-words

Output:

C:\PREPRINT,
'different',
'insurance',
'language',
'information',
'analysis',
'mining',
'different',
'models',
'Sentence -',
len(long_frequent_words)

712
freq Dist
→ PDF
→ PDF
from PDF
analyse
if use
data
process
→ PyPDF
split
format
→ PDF
data
need
accuracy
→ PDF
test
high
How
com
A

LTKR

Date
2023-09-20

Q10. What is long - frequent - words). plot().

long dist stands for Portable Document Format

→ PDF stands for portable document format
→ pdfminer is a Python library to extract text

from PDF documents. It provides away to posse.

analyze & work with content of PDF files, making
it useful for various applications, including
extracting, text analysis and document +
data processing.

processing.

→ PyPDF2: It is a Python library used for merging
splitting, extracting text, rotating pages and
formatting PDF files. It does all the manipulation.

→ PDF parser: It is used to extract all types of
data from PDF documents.

need for passing it to avoid errors and get
accurate data.

→ PDF Query: It is used to extract images and
text data from PDF documents. It has a
high performance in terms of text extraction.
How to install Pdf parser?

Command: pip install pdf miner -six.

✓ ✓ ✓

```

import pandas as pd.

# Load the two Excel files.
file1 = 'xl5 - child - marriage - database May - 2024,
file2 = 'xl5 - child - marriage - database May - 2024 - XL5.

df1 = pd.read_excel(file1, sheet_name='child labour')
df2 = pd.read_excel(file2, sheet_name='child marriage')

# Merge the two datasets on a common column.
# Replace 'Common - Column' with actual column names present in both data sets.
merged_df = pd.merge(df1, df2, how='outer', on='Country')

# Check for duplicates and remove them
merged_df.drop_duplicates(inplace=True)

# Check for the missing data.
missing_data_summary = merged_df.isnull().sum()

# Eliminate mismatches (removing rows with missing values)
merged_df = merged_df.dropna(inplace=True)

# Clean line breaks, spaces, and special characters
# Assuming text columns need this cleaning.
def clean_text(text):
    if isinstance(text, str):
        # Remove line breaks and extra spaces.
        text = ' '.join(text.split())
        # Remove special characters (keeping only alphanumeric and spaces)
        text = ' '.join([c for c in text if c.isalnum() or
                        c.isspace()])
    return text

```


16
 married by 18
 unnamed : 4-9
 Reference years
 observation footnote
 Data Source
 married by 18.
 unnamed : 9
 Reference years.
 observation footnote .
 Data Source :
 dtype: int. 64
 summary of cleaned data:
 class: pandas. Cose - DataFrame - Data Frame's
 Index : 4 entities, 18 to 135.
 Data columns (total: 20 columns):
 # - column - non-null count - Dtype
 0 - countries and areas) 4 non-NULL object.
 1 child_labour(%.)+ 4 non-NULL float 64
 2 (2015 - 2023)* 4
 3 unnamed : 2-X 4 non-NULL object
 4 unnamed : 3 4 non-NULL object
 5 unnamed : 4 non-NULL float 64
 6 unnamed : 6 4 non-NULL object
 7 unnamed : 7 4 non-NULL object .

Data Wrangling and Visualization - Python/R Programming/Power BI

OUTPUT

8	married by 15	↳ non-NULL	float 64
9	unnamed : 2-9	↳ non-NULL	object
10	married by 18	↳ non-NULL	float 64
11	unnamed : 6-9	↳ non-NULL	object
12	Reference year	↳ non-NULL	object
13	observation footnote	↳ non-NULL	object
14	Data Source	↳ non-NULL	object
15	married by 18.1 - other	↳ non-NULL	float 64
16	unnamed 19	↳ non-NULL	object
17	Reference year 2	↳ non-NULL	object
18	Source 1	↳ non-NULL	object
19	Data Source 1	↳ non-NULL	object

Computerized child labour(%)

unnamed unnamed unnamed unnamed unnamed

CSA 2018 unicef-d ILO Calculations
CSA 2018 unicef-d ILO Calculations
CSA 2018 unicef-d ILO Calculations
3 4 5 6 7 8.9 X X 19.4 Y Y 20.4 Y Y 33.0 Y Y 10.1 Y Y 12.5 Y Y 11.4 Y Y 31.5 Y Y 19.9 Y Y 18.8 14.1 13.0 135

~~Access~~ ~~Education~~ ~~18.8~~ ~~14.1~~ ~~13.0~~ ~~135~~

+1n (2015-2003) COOP Unnamed

Experiment-4

Date: 24/08/2024

AIM

Draw the chart between perceived corruption scores compared to the child labour percentages using matplotlib.

PROCEDURE

Step-1: Install the Matplotlib package

Step-2: Read the required data using read.csv method of pandas library

Step-3: Define the x-axis and corresponding y-axis values as lists.

Step-4: Plot them on canvas using .plot() function.

Step-5: Give a name to x-axis and y-axis using xlabel() and ylabel() functions.

Step-6: Give a title to your plot using .title() function.

Step-7: To view your plot, use .show() function.

SOURCE CODE

```
import pandas as pd  
# attempt to load the datasets  
corruption - df = pd.read_excel('Corruption_Perception_index.xlsx')  
unicef - df = pd.read_excel('UNICEF_Oct_2014 - XLSX')
```

```
# display the first few rows of each dataset &
# understand their structure.
corruption_df_head = corruption_df.head()
unicef_df_head = unicef_df.head()
```

```
unicef_df_head
unicef_df_head
```

```
# extract and clean the relevant columns from the
# corruption data.
```

```
corruption_cleaned = corruption . df . iloc [2:(1,)]
corruption_cleaned = corruption . iloc [2:(1,)]
```

```
corruption_cleaned = corruption - corruption . dropna ()
```

```
unicef_cleaned = unicef . dropna (dropop = True)
unicef_cleaned = unicef . dropna (dropop = True)
```

```
unicef_cleaned = unicef . dropna (dropop = True)
unicef_cleaned = unicef . dropna (dropop = True)
```

```
unicef_cleaned = unicef . dropna (dropop = True)
unicef_cleaned = unicef . dropna (dropop = True)
```

```
# merge the datasets
```

```
messy_df = pd . merge (corruption_cleaned,
                      unicef_cleaned,
                      left_on = 'Country',
                      right_on = 'Country')
```

to appropriate types

convert the columns to appropriate types
[corruption score] = pd.to_numeric (merged -
merged - df['Corruption Score'], errors = 'Coerce')
[Child Labour Percentage] - pd.to_numeric
[Child Labour Percentage] - pd.to_numeric
merge - df [child labour percentage] axis=0 (Coerce)
merged . df [child labour percentage] (Coerce)
drop rows with any non values.
drop rows with any non values.
merged . df = merged . df . dropna (C) . reset_index (drop=True)

Country	Corruption Score	Child Labour Percentage	Corruption Score
Uruguay	6	2.9	6
Chile	9	6.6	9
Saint Lucia	3	3.9	3
Botswana	7	0.0	7
Bhutan	4	2.9	4

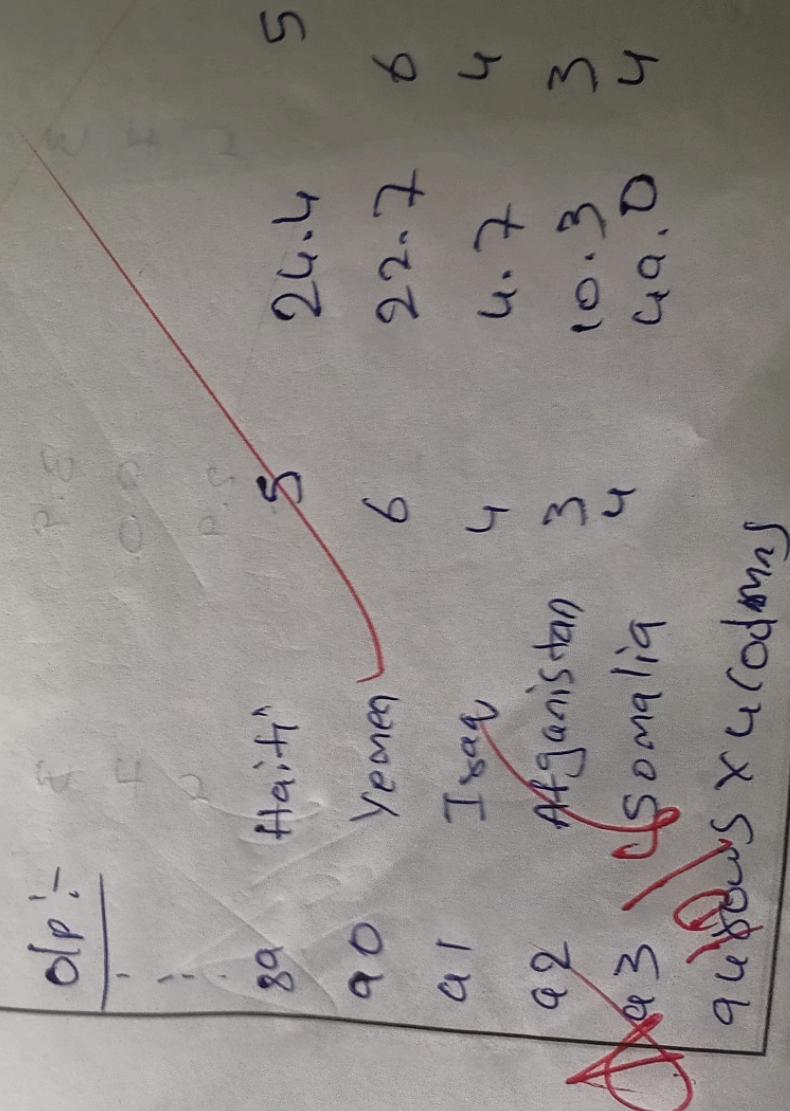
(1,7)
(1,7)
(1,7)
(1,7)
Percentage

Rounds
Ricel.

```

import matplotlib.pyplot as plt
import matplotlib
matplotlib.style.use('ggplot')
plt.scatter(merged['Corruption Score'],
            merged['Child Labour Percentage'],
            color='red',
            s=100)
plt.xlabel('Corruption Perception Index (CPI)')
plt.ylabel('Child Labour Percentage')
plt.title('Comparison of Corruption Perception Index & Child Labour Percentage')
plt.grid(True)
# Show the graph
plt.show()

```



Experiment-5

Date: 10/02/29

AIM

Write a python program to download & display content of robot.txt for en.wikipedia.org.

PROCEDURE

Step-1: Use Requests (HTTP for Humans) Library for Web Scraping

Step-2: Scrape the robots.txt file

Step-3: print the data

SOURCE CODE

Web Crawler:- A web crawler or web spider, is a computer program that is used to search and automatically index website content and other information over the internet.

Crawling:- It is the discovery process in which Search engine sends out a team of robots (spiders) to find new and updated content.

Indexing: Website indexing is the process by which a search engine adds web content to its index. This is done by "crawling" web pages for keywords, meta-data and placed signals that tell search engines it and where to rank content. Indexed website should have navigable, findable and clearly understood content strategy.

Robots.txt file:

A robot.txt file is a set of instruction used by website to tell search engine which pages should and should not be crawled.

web Scapping:

→ web Scapping refers to collection and structuring the data from web source in a more convenient format. It involves no processing or review of data.

web Scapping can be used to build the data set that are to be used in data mining.

Data mining:

Data mining refers to analyzing large data sets to reveal useful information and patterns. It doesn't require data processing extraction.

Data mining refers or is a the process of large data to uncover trend and valuable insights.

Request module:

Python request a library os making HTTP request it provide an easy - to - use interface that make working with HTTP very simple, which mean it simplifies the process of studying and receiving the data from websites.

CODE:

~~import requests as rq.~~

~~response = rq.get("https://en.wikipedia.org/w/index.php?title=Robot&oldid=5800000")~~

data = response.text

print(data)

print("==")

print(data)

~~Output:~~
robots. text for https://en.wikipedia.org/

stq

sets

OUTPUT

```
# dobots.txt for https://en.wikipedia.org/w/index.php?title=Output&oldid=92050749#friends.
```

```
# Please feel free to use this file. There are a lot of pages on this site,
# and I don't want to break anything.
# Some misbehaving spiders out there that go
# always fast. If you're
# responsible, your access to the site may be
# blocked.
```

```
# Observing spamming (large amount of http://)
# wikipedia.org/?curid=NNNNNN
# and ignoring certain responses, claims to
# respect so bots:
# http://microsoftbot.com/
# http://mini12bot.com/
# user-agent: microsoftbot
Disallow: /
```

~~# advertising-spotted bots:~~

~~# user-agent: mediapartners-google~~

~~Disallow: /~~

~~wikipedia.org work bots:~~

~~# user-agent: IsraBot~~

~~Disallow: /~~