

Practical 1:

Creating Data Model using Cassandra.

Cassandra Data Model

Step-1:

- Open a folder Datascience\apache-cassandra-3.11.4-bin\apache-cassandra-3.11.4\bin\cassandra.bat
- now open IDLE (PYTHON GUI)
- go to file -> open -> select (Datascience\apache-cassandra-3.11.4-bin\apache-cassandra-3.11.4\bin\select-cqlsh.py
- inside sqlsh.py -> run -> run module

Step-2: command to Create keyspace:

create keyspace DATASCI WITH replication={ 'class': 'SimpleStrategy', 'replication_factor': 3};

```
cqlsh> create keyspace DATASCI WITH replication={ 'class': 'SimpleStrategy', 'replication_factor': 3};
```

step – 3: command to use keyspace run this command

cqlsh> use datasci;

```
cqlsh> use datasci;
```

step – 4: command to create a new table

```
cqlsh:datasci> create table student(  
    student_id int PRIMARY KEY,  
    student_name text,  
    student_city text,  
    student_fees varint,  
    student_phone varint  
);
```

```
cqlsh> use datasci;  
cqlsh:datasci> create table student(  
    student_id int PRIMARY KEY,  
    student_name text,  
    student_city text,  
    student_fees varint,  
    student_phone varint  
);  
...  
...  
...  
...  
...  
... cqlsh
```

Step – 5: command to display created keyspace list

Desc keyspace;

```
Python 2.7.15 Shell  
File Edit Shell Debug Options Window Help  
:datasci> desc keyspace;  
  
CREATE KEYSPACE datasci WITH replication = {'class': 'SimpleStrategy', 'replication_factor': '3'  
'}) AND durable_writes = true;  
  
CREATE TABLE datasci.student (  
    student_id int PRIMARY KEY,  
    student_city text,  
    student_fees varint,  
    student_name text,  
    student_phone varint  
) WITH bloom_filter_fp_chance = 0.01  
    AND caching = {'keys': 'ALL', 'rows_per_partition': 'NONE'}  
    AND comment = ''  
    AND compaction = {'class': 'org.apache.cassandra.db.compaction.SizeTieredCompactionStrategy'  
, 'max_threshold': '32', 'min_threshold': '4'}  
    AND compression = {'chunk_length_in_kb': '64', 'class': 'org.apache.cassandra.io.compress.L  
Z4Compressor'}  
    AND crc_check_chance = 1.0  
    AND dclocal_read_repair_chance = 0.1  
    AND default_time_to_live = 0  
    AND gc_grace_seconds = 864000  
    AND max_index_interval = 2048  
    AND memtable_flush_period_in_ms = 0  
  
    AND min_index_interval = 128  
    AND read_repair_chance = 0.0  
    AND speculative_retry = '99PERCENTILE';
```

Step – 6: command to alter keyspace

```
alter keyspace datasci with replication={'class':'SimpleStrategy','replication_factor':2};
```

```
cqlsh:datasci> alter keyspace datasci with replication={'class':'SimpleStrategy','replication_factor':2};
```

step – 7 : command to display all the tables of the keyspaces

```
cqlsh:datasci> desc tables;
```

```
cqlsh:datasci> desc tables;

student
```

step- 8: command to alter table

```
cqlsh:datasci>
```

```
alter table student
```

```
add student_gender text;
```

```
cqlsh:datasci>
alter table student
add student_gender text;
```

step- 9: command to insert data into table

```
insert into student(student_id,student_city,student_fees,student_name,student_phone)
```

```
values(1,'Bhy',5000,'pooja',0939293939);
```

(you can only add one value at a time)

```
cqlsh:datasci> insert into student(student_id,student_city,student_fees,student_name,student_phone)
values(1,'Bhy',5000,'pooja',0939293939);
... cqlsh:datasci> insert into student(student_id,student_city,student_fees,student_name,student_phone) values(2,'Bhy',1000,'hima',0923233939);
cqlsh:datasci> insert into student(student_id,student_city,student_fees,student_name,student_phone) values(3,'Bhy',4500,'mira',1234293939);
cqlsh:datasci> insert into student(student_id,student_city,student_fees,student_name,student_phone) values(4,'Bhy',7800,'kavita',1234293656);
cqlsh:datasci> insert into student(student_id,student_city,student_fees,student_name,student_phone) values(5,'Bhy',9200,'soniya',1234298989);
cqlsh:datasci> select * from student;
```

Step – 10: command to show the table

```
cqlsh:datasci> select * from student;
```

```
cqlsh:datasci> select * from student;

student_id | student_city | student_fees | student_gender | student_name | student_phone
-----+-----+-----+-----+-----+-----
5 | Bhy | 9200 | null | soniya | 1234298989
1 | Bhy | 5000 | null | pooja | 939293939
2 | Bhy | 1000 | null | hima | 923233939
4 | Bhy | 7800 | null | kavita | 1234293656
3 | Bhy | 4500 | null | mira | 1234293939

(5 rows)
```

step – 11: command to update table

```
cqlsh:datasci> update student set student_fees=200000,student_name='hima' where student_id=2;
```

```
cqlsh:datasci> update student set student_fees=200000,student_name='hima' where student_id=2;
cqlsh:datasci> select * from student;
```

```
student_id | student_city | student_fees | student_gender | student_name | student_phone
-----+-----+-----+-----+-----+-----
5 | Bhy | 9200 | null | soniya | 1234298989
1 | Bhy | 5000 | null | pooja | 939293939
2 | Bhy | 200000 | null | hima | 923233939
4 | Bhy | 7800 | null | kavita | 1234293656
3 | Bhy | 4500 | null | mira | 1234293939

(5 rows)
```

step – 12: command to refresh the table

cqlsh:datasci> truncate student;

```
cqlsh:datasci> truncate student;
cqlsh:datasci> select * from student;

\ student_id | student_city | student_fees | student_gender | student_name | student_phone
-----+-----+-----+-----+-----+-----
(0 rows)
```

step – 13: command to delete the specific column data from the table

cqlsh:datasci> delete student_city from student where student_id=2;

```
cqlsh:datasci> delete student_city from student where student_id=2;
cqlsh:datasci> select * from student;

 student_id | student_city | student_fees | student_gender | student_name | student_phone
-----+-----+-----+-----+-----+-----
          5 |          Bhy |         10000 |             null |          prima | 1234293939
          1 |          Bhy |          5000 |             null |          seema | 1234293939
          2 |         null |          5000 |             null |          zoya | 1234293939
          4 |          Bhy |          7000 |             null |          kajal | 1234293939
          3 |          Bhy |          6000 |             null |          rohini | 1234293939
(5 rows)
```

PRACTICAL NO2

A.Text delimited CSVto HORUS format

```
import pandas as pd
sInputFileName='C:/VKHCG/05-DS/9999-Data/Country_Code.csv'
InputData=pd.read_csv(sInputFileName,encoding="latin-1")
ProcessData=InputData
ProcessData.drop(['ISO-3-Code', 'ISO-2-CODE'], axis=1,inplace=True)
ProcessData.rename(columns={'ISO-M49': 'CountryNumber', 'Country': 'CountryName'}, inplace=True)
ProcessData.set_index('CountryNumber', inplace=True)
ProcessData.sort_values('CountryName', axis=0, ascending=True, inplace=True)
print(ProcessData.head(10))
```

```
===== RESTART: C:\gaurav\pr:
CountryName
CountryNumber
4              Afghanistan
248            Aland Islands
8              Albania
12             Algeria
16            American Samoa
20              Andorra
24              Angola
660            Anguilla
10             Antarctica
28            Antigua and Barbuda
```

B>XML to HORUS Format

```
# Utility Start XML to HORUS =====
# Standard Tools
import pandas as pd
import xml.etree.ElementTree as ET
def df2xml(data):
    header = data.columns
    root = ET.Element('root')
    for row in range(data.shape[0]):
        entry = ET.SubElement(root,'entry')
        for index in range(data.shape[1]):
            schild=str(header[index])
            child = ET.SubElement(entry, schild)
            if str(data[schild][row]) != 'nan':
                child.text = str(data[schild][row])
            else:
                child.text = 'n/a'
            entry.append(child)
    result = ET.tostring(root)
    return result
def xml2df(xml_data):
    root = ET.XML(xml_data)
    all_records = []
    for i, child in enumerate(root):
        record = { }
        for subchild in child:
            record[subchild.tag] = subchild.text
        all_records.append(record)
    return pd.DataFrame(all_records)
sInputFileName='C:/VKHCG/05-DS/9999-Data/Country_Code.xml'
InputData = open(sInputFileName).read()
print('Input Data Values =====')
ProcessDataXML=InputData
ProcessData=xml2df(ProcessDataXML)
ProcessData.drop('ISO-2-CODE', axis=1,inplace=True)
ProcessData.drop('ISO-3-Code', axis=1,inplace=True)
ProcessData.rename(columns={'Country': 'CountryName'}, inplace=True)
ProcessData.rename(columns={'ISO-M49': 'CountryNumber'}, inplace=True)
ProcessData.set_index('CountryNumber', inplace=True)
ProcessData.sort_values('CountryName', axis=0, ascending=False, inplace=True)
print('Process Data Values =====')
print(ProcessData.head(5))
print('=====')
OutputData=ProcessData
sOutputFileName='C:/VKHCG/05-DS/9999-Data/HORUS-XML-Country.csv'
OutputData.to_csv(sOutputFileName, index = False, encoding="latin-1")
print('XML to HORUS - Done')
```

```
===== RESTART: C:/gaurav/practicals/XML_TO_HRS_COPY.py =====
Process Data Values =====
      CountryName
CountryNumber
716      Zimbabwe
716      Zimbabwe
716      Zimbabwe
716      Zimbabwe
716      Zimbabwe
=====
XML to HORUS - Done
>>> |
```

C>JSON to HORUS Format

```

import pandas as pd
sInputFileName='C:/VKHCG/05-DS/9999-Data/Country_Code.json'
InputData=pd.read_json(sInputFileName, orient='index', encoding="latin-1")
ProcessData=InputData
ProcessData.drop(['ISO-3-Code', 'ISO-2-CODE'], axis=1,inplace=True)
ProcessData.rename(columns={'ISO-M49': 'CountryNumber', 'Country': 'CountryName'}, inplace=True)
ProcessData.set_index('CountryNumber', inplace=True)
ProcessData.sort_values('CountryName', axis=0, ascending=False, inplace=True)
print(ProcessData.head(5))
OutputData=ProcessData
sOutputFileName='c:/VKHCG/05-DS/9999-Data/HORUS-JSON-Country.csv'
OutputData.to_csv(sOutputFileName, index = False, encoding="latin-1")
print('JSON to HORUS - Done')

```

```

===== RESTART: C:\gaurav\practicals\json_to_hrs.py =====
CountryName
CountryNumber
716 Zimbabwe
894 Zambia
887 Yemen
732 Western Sahara
876 Wallis and Futuna Islands
JSON to HORUS - Done
>>> |

```

D>MySql Database to HORUS Format

```

import pandas as pd
import sqlite3 as sq
conn = sq.connect('C:/VKHCG/05-DS/9999-Data/utility.db')
sSQL='select * FROM ' + 'Country_Code' + ';'
InputData=pd.read_sql_query(sSQL, conn)
ProcessData=InputData
ProcessData.drop(['ISO-3-Code', 'ISO-2-CODE'], axis=1,inplace=True)
ProcessData.rename(columns={'ISO-M49': 'CountryNumber', 'Country': 'CountryName'}, inplace=True)
ProcessData.set_index('CountryNumber', inplace=True)
ProcessData.sort_values('CountryName', axis=0, ascending=False, inplace=True)
print('Process Data Values =====')
print(ProcessData.head(5))
OutputData=ProcessData
sOutputFileName='C:/VKHCG/05-DS/9999-Data/HORUS-CSV-Country.csv'
OutputData.to_csv(sOutputFileName, index = False, encoding="latin-1")
print('Database to HORUS - Done')

```

```

===== RESTART: C:/gaurav/practicals/db_to_hs.py =====
Process Data Values =====
index CountryName
CountryNumber
716 246 Zimbabwe
894 245 Zambia
887 244 Yemen
732 243 Western Sahara
876 242 Wallis and Futuna Islands
Database to HORUS - Done
>>> |

```

E>Picture (JPEG) to HORUS Format

```

import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import imageio

```

```

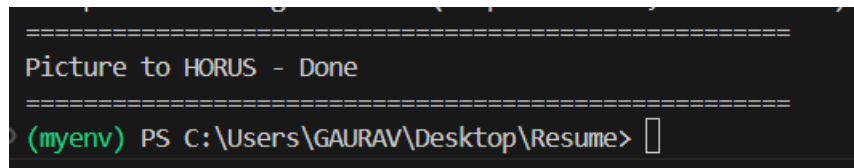
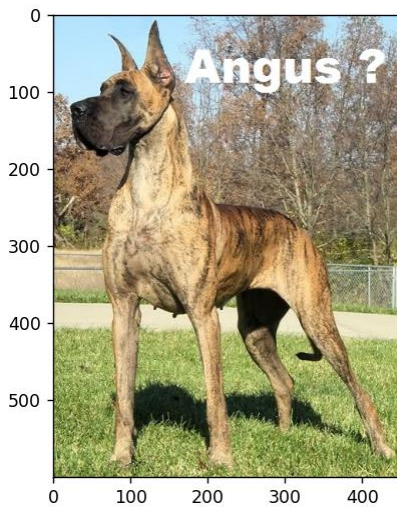
sInputFileName='C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/Angus.jpg'
InputData = imageio.imread(sInputFileName, mode='RGBA')

```

```

ProcessRawData=InputData.flatten()
y=InputData.shape[2] + 2
x=int(ProcessRawData.shape[0]/y)
ProcessData=pd.DataFrame(np.reshape(ProcessRawData, (x, y)))
sColumns= ['XAxis','YAxis','Red', 'Green', 'Blue','Alpha']
ProcessData.columns=sColumns
ProcessData.index.names =['ID']
plt.imshow(InputData)
plt.show()
OutputData=ProcessData
sOutputFileName='C:/VKHCG/05-DS/9999-Data/HORUS-Picture.csv'
OutputData.to_csv(sOutputFileName, index = False)
print('=====')
print('Picture to HORUS - Done')
print('=====')

```



F>Video to HORUS Format

```

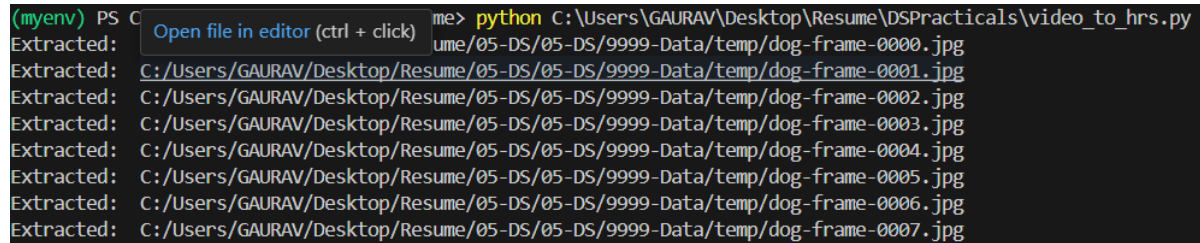
1st =====
import os
import shutil
import cv2
sInputFileName = 'C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/dog.mp4'
sDataBaseDir = 'C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/temp'
if os.path.exists(sDataBaseDir):
    shutil.rmtree(sDataBaseDir)
if not os.path.exists(sDataBaseDir):
    os.makedirs(sDataBaseDir)
vidcap = cv2.VideoCapture(sInputFileName)
if not vidcap.isOpened():
    print('Error: Could not open video file')
    exit()
count = 0
while True:
    success, image = vidcap.read()
    if not success:
        break
    sFrame = sDataBaseDir + '/dog-frame-' + str(format(count, '04d')) + '.jpg'
    print('Extracted: ', sFrame)
    cv2.imwrite(sFrame, image)
    if os.path.getsize(sFrame) == 0:
        os.remove(sFrame)
        print('Removed: ', sFrame)
        continue

```

```

count += 1
if cv2.waitKey(10) == 27:
    break
print('Generated: ', count, ' Frames')
print('=====')
print('Movie to Frames HORUS - Done')
print('=====')

```



```

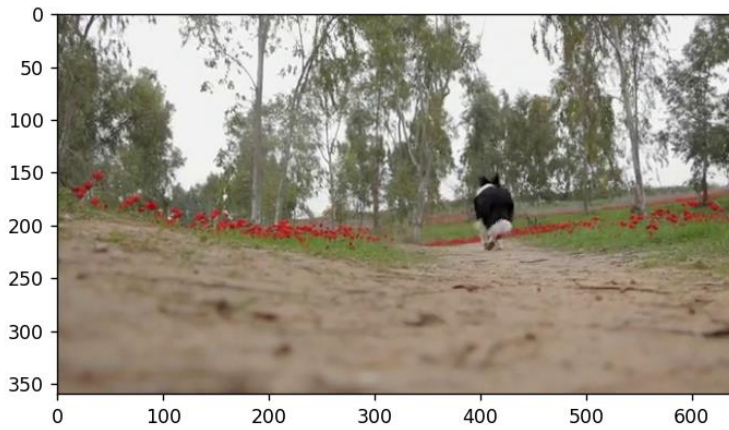
(myenv) PS C:\Users\GAURAV\Desktop\Resume\DSPracticals\video_to_hrs.py
Extracted: C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/temp/dog-frame-0000.jpg
Extracted: C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/temp/dog-frame-0001.jpg
Extracted: C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/temp/dog-frame-0002.jpg
Extracted: C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/temp/dog-frame-0003.jpg
Extracted: C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/temp/dog-frame-0004.jpg
Extracted: C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/temp/dog-frame-0005.jpg
Extracted: C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/temp/dog-frame-0006.jpg
Extracted: C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/temp/dog-frame-0007.jpg

```

```

2nd part =====
import imageio
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import os
sDataBaseDir='C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/temp'
f=0
for file in os.listdir(sDataBaseDir):
    if file.endswith(".jpg"):
        f += 1
sInputFileName=os.path.join(sDataBaseDir, file)
InputData = imageio.imread(sInputFileName, mode='RGBA')
ProcessRawData=InputData.flatten()
y=InputData.shape[2] + 2
x=int(ProcessRawData.shape[0]/y)
ProcessFrameData=pd.DataFrame(np.reshape(ProcessRawData, (x, y)))
ProcessFrameData['Frame']=file
plt.imshow(InputData)
plt.show()
ProcessData = []
if f == 1:
    ProcessData=ProcessFrameData
else:
    ProcessData=ProcessData.append(ProcessFrameData)
if f > 0:
    # ProcessData = pd.DataFrame(ProcessFrameData)
    print(ProcessData)
    sColumns= ['XAxis','YAxis','Red', 'Green', 'Blue','Alpha','FrameName']
    ProcessData.columns=sColumns
    ProcessFrameData.index.names =['ID']
    print('Rows: ',ProcessData.shape[0])
    print('Columns :',ProcessData.shape[1])
ProcessData.to_csv('C:/VKHCG/05-DS/9999-Data/HORUS-Movie-Frame.csv' , index = False)
print('Processed ; ', f, ' frames')

```

	0	1	2	3	4	5	Frame
0	94	95	89	255	183	184	dog-frame-0048.jpg
1	178	255	207	208	202	255	dog-frame-0048.jpg
2	114	115	109	255	98	99	dog-frame-0048.jpg
3	93	255	163	164	158	255	dog-frame-0048.jpg
4	155	156	150	255	175	176	dog-frame-0048.jpg
...
153595	108	255	151	129	106	255	dog-frame-0048.jpg
153596	150	128	107	255	149	127	dog-frame-0048.jpg
153597	106	255	149	127	106	255	dog-frame-0048.jpg
153598	150	128	107	255	151	129	dog-frame-0048.jpg
153599	108	255	151	129	108	255	dog-frame-0048.jpg

G. Audio to HORUS Format

```

from scipy.io import wavfile
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
def show_info(aname, a,r):
    print(' ')
    print("Audio:", aname)
    print(' ')
    print("Rate:", r)
    print(' ')
    print("shape:", a.shape)
    print("dtype:", a.dtype)
    print("min, max:", a.min(), a.max())
    print(' ')
    plot_info(aname, a,r)
def plot_info(aname, a,r):
    sTitle= 'Signal Wave - '+ aname + ' at ' + str(r) + 'hz'
    plt.title(sTitle)
    sLegend=[]
    for c in range(a.shape[1]):
        sLabel = 'Ch' + str(c+1)
        sLegend=sLegend+[str(c+1)]
        plt.plot(a[:,c], label=sLabel)
    plt.legend(sLegend)
    plt.show()
#=====
sInputFileName='C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/2ch-sound.wav'
InputRate, InputData = wavfile.read(sInputFileName)
show_info("2 channel", InputData,InputRate)
ProcessData=pd.DataFrame(InputData)
sColumns= ['Ch1','Ch2']
ProcessData.columns=sColumns
OutputData=ProcessData
sOutputFileName='C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/HORUS-Audio-2ch.csv'
OutputData.to_csv(sOutputFileName, index = False)
sInputFileName='C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/4ch-sound.wav'
InputRate, InputData = wavfile.read(sInputFileName)
show_info("4 channel", InputData,InputRate)
ProcessData=pd.DataFrame(InputData)
sColumns= ['Ch1','Ch2','Ch3', 'Ch4']
ProcessData.columns=sColumns
OutputData=ProcessData
sOutputFileName='C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/HORUS-Audio-4ch.csv'
OutputData.to_csv(sOutputFileName, index = False)
sInputFileName='C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/6ch-sound.wav'
print('=====')
print('Processing : ', sInputFileName)

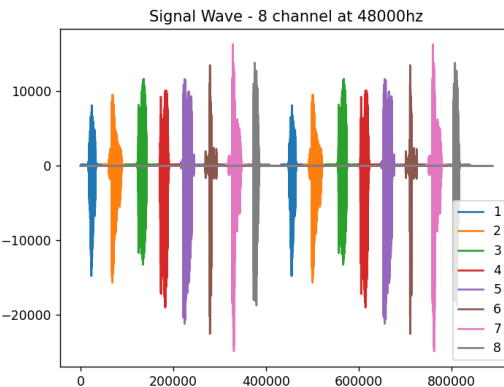
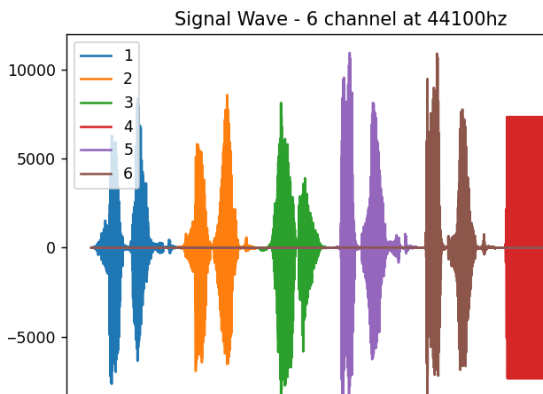
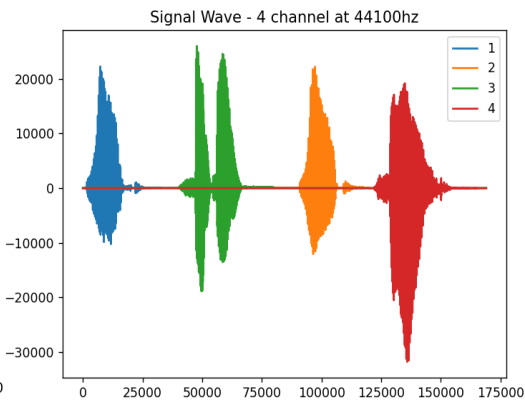
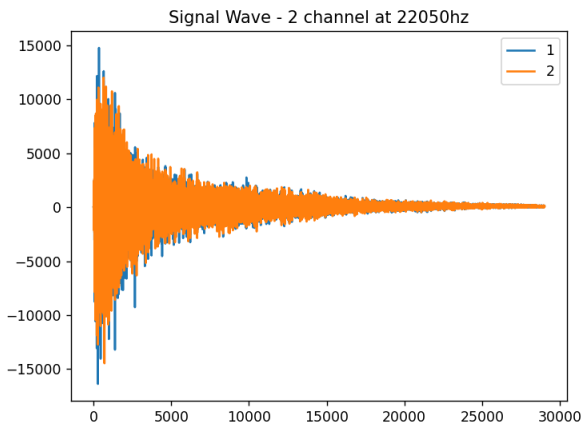
```



```

print('=====')
InputRate, InputData = wavfile.read(sInputFileName)
show_info("6 channel", InputData, InputRate)
ProcessData=pd.DataFrame(InputData)
sColumns= ['Ch1','Ch2','Ch3', 'Ch4', 'Ch5','Ch6']
ProcessData.columns=sColumns
OutputData=ProcessData
sOutputFileName='C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/HORUS-Audio-6ch.csv'
OutputData.to_csv(sOutputFileName, index = False)
#=====
sInputFileName='C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/8ch-sound.wav'
print('=====')
print('Processing : ', sInputFileName)
print('=====')
InputRate, InputData = wavfile.read(sInputFileName)
show_info("8 channel", InputData, InputRate)
ProcessData=pd.DataFrame(InputData)
sColumns= ['Ch1','Ch2','Ch3', 'Ch4', 'Ch5','Ch6','Ch7','Ch8']
ProcessData.columns=sColumns
OutputData=ProcessData
sOutputFileName='C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/HORUS-Audio-8ch.csv'
OutputData.to_csv(sOutputFileName, index = False)
print('=====')
print('Audio to HORUS - Done')

```



```

=====
Processing : C:/Users/GAURAV/Desktop/Resume/05-DS/05-DS/9999-Data/8ch-sound.wav
=====

```

Audio: 8 channel

Rate: 48000

shape: (888000, 8)

dtype: int16

min, max: -24859 16303

Practical 3: Utilities and Auditing

Basic Utility Design

A. Fixers Utilities:

Fixers enable your solution to take your existing data and fix a specific quality issue.

```
import string
```

```
import datetime as dt
```

1 Removing leading or lagging spaces from a data entry

```
print('#1 Removing leading or lagging spaces from a data entry');
```

```
baddata = " Data Science with too many spaces is bad!!! "
```

```
print('>',baddata,<')
```

```
cleandata=baddata.strip()
```

```
print('>',cleandata,<')
```

2 Removing nonprintable characters from a data entry

```
print('#2 Removing nonprintable characters from a data entry')
```

```
printable = set(string.printable)
```

```
baddata = "Data\x00Science with\x02 funny characters is \x10bad!!!"
```

```
cleandata="".join(filter(lambda x: x in string.printable,baddata))
```

```
print('Bad Data : ',baddata);
```

```
print('Clean Data : ',cleandata)
```

3 Reformatting data entry to match specific formatting criteria.

```
# Convert YYYY/MM/DD to DD Month YYYY
```

```
print('# 3 Reformatting data entry to match specific formatting criteria.')
```

```
baddate = dt.date(2019, 10, 31)
```

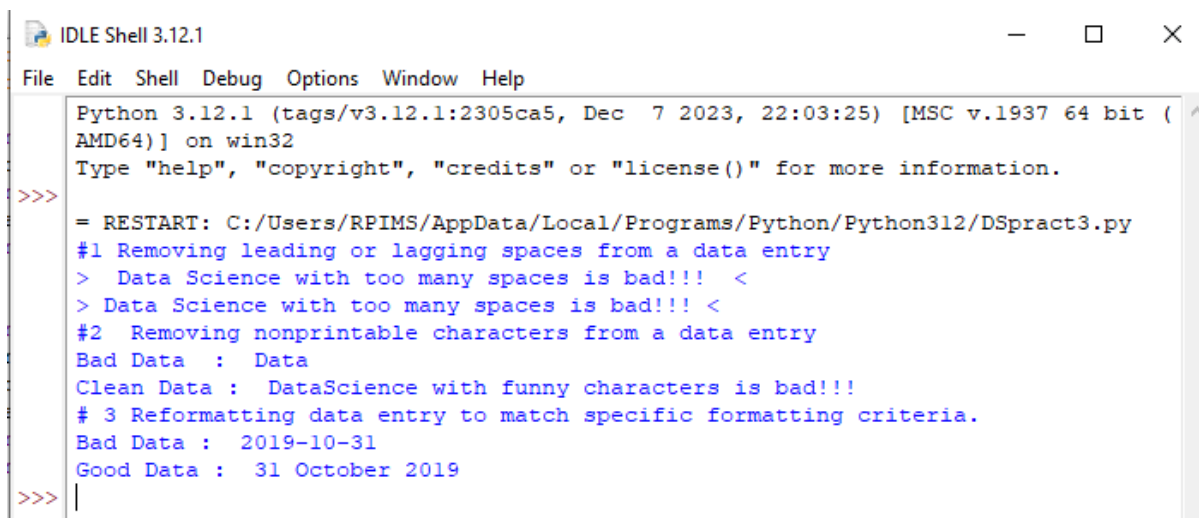
```
baddata=format(baddate,'%Y-%m-%d')
```

```
gooddate = dt.datetime.strptime(baddata,'%Y-%m-%d')
```

```
gooddata=format(gooddate,'%d %B %Y')
```

```
print('Bad Data : ',baddata)
```

```
print('Good Data : ',gooddata)
```



```
IDLE Shell 3.12.1
Python 3.12.1 (tags/v3.12.1:2305ca5, Dec 7 2023, 22:03:25) [MSC v.1937 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license()" for more information.
>>>
= RESTART: C:/Users/RPIMS/AppData/Local/Programs/Python/Python312/DSpract3.py
#1 Removing leading or lagging spaces from a data entry
> Data Science with too many spaces is bad!!! <
> Data Science with too many spaces is bad!!! <
#2 Removing nonprintable characters from a data entry
Bad Data : Data
Clean Data : DataScience with funny characters is bad!!!
# 3 Reformatting data entry to match specific formatting criteria.
Bad Data : 2019-10-31
Good Data : 31 October 2019
>>> |
```

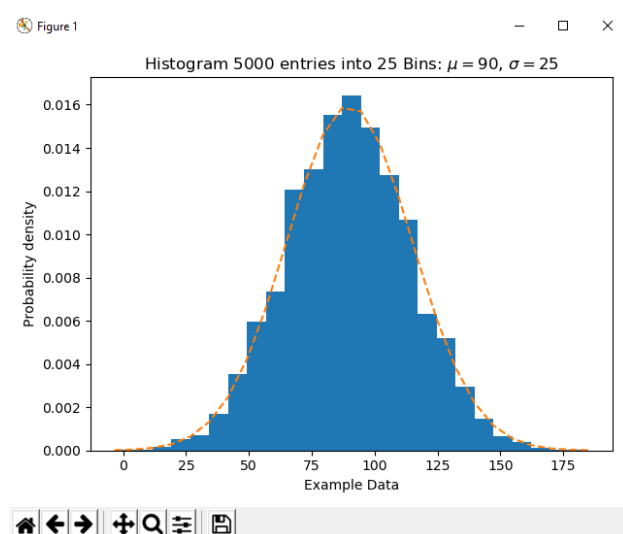
B. Data Binning or Bucketing

```

import numpy as np
import matplotlib.mlab as mlab
import matplotlib.pyplot as plt
np.random.seed(0)
# example data
mu = 90 # mean of distribution
sigma = 25 # standard deviation of distribution
x = mu + sigma * np.random.randn(5000)
num_bins = 25
fig, ax = plt.subplots()
# the histogram of the data
n, bins, patches = ax.hist(x, num_bins, normed=1)
# add a 'best fit' line
y = mlab.normpdf(bins, mu, sigma)
ax.plot(bins, y, '--')
ax.set_xlabel('Example Data')
ax.set_ylabel('Probability density')
sTitle=r'Histogram ' + str(len(x)) + ' entries into ' + str(num_bins) + ' Bins:  $\mu=$ ' + str(mu) + ',  $\sigma=$ ' + str(sigma) + '$'
ax.set_title(sTitle)
fig.tight_layout()
sPathFig='C:/VKHCG/05-DS/4000-UL/0200-DU/DU-Histogram.png'
fig.savefig(sPathFig)
plt.show()

```

Output:



C. Averaging of Data

Input:

```

import pandas as pd
InputFileName='IP_DATA_CORE.csv'
OutputFileName='Retrieve_Router_Location.csv'
Base='C:/VKHCG'
print('Working Base :',Base, ' using ')
sFileName=Base + '/01-Vermeulen/00-RawData/' + InputFileName
print('Loading :',sFileName)
IP_DATA_ALL=pd.read_csv(sFileName,header=0,low_memory=False, usecols=['Country','Place
Name','Latitude','Longitude'], encoding="latin-1")
IP_DATA_ALL.rename(columns={'Place Name': 'Place_Name'}, inplace=True)

```

```
AllData=IP_DATA_ALL[['Country', 'Place_Name','Latitude']]
print(AllData)
MeanData=AllData.groupby(['Country', 'Place_Name'])['Latitude'].mean()
print(MeanData)
```

Output:

```
===== RESTART: C:/Users/RPIMS/AppData/Local/
Working Base : C:/VKHCG using
Loading : C:/VKHCG/01-Vermeulen/00-RawData/IP_DATA_CORE.csv
   Country Place_Name Latitude
0        US   New York  40.7528
1        US   New York  40.7528
2        US   New York  40.7528
3        US   New York  40.7528
4        US   New York  40.7528
...      ...      ...
3557     DE    Munich  48.0915
3558     DE    Munich  48.1833
3559     DE    Munich  48.1000
3560     DE    Munich  48.1480
3561     DE    Munich  48.1480

[3562 rows x 3 columns]
Country Place_Name
DE      Munich      48.143223
GB      London      51.509406
US      New York      40.747044
Name: Latitude, dtype: float64
```

D. Outlier Detection

```
import pandas as pd
InputFileName='IP_DATA_CORE.csv'
OutputFileName='Retrieve_Router_Location.csv'
Base='C:/VKHCG'
print('Working Base :',Base)
sFileName=Base + '/01-Vermeulen/00-RawData/' + InputFileName
print('Loading :',sFileName)
IP_DATA_ALL=pd.read_csv(sFileName,header=0,low_memory=False, usecols=['Country','Place
Name','Latitude','Longitude'], encoding="latin-1")
IP_DATA_ALL.rename(columns={'Place Name': 'Place_Name'}, inplace=True)
LondonData=IP_DATA_ALL.loc[IP_DATA_ALL['Place_Name']=='London']
AllData=LondonData[['Country', 'Place_Name','Latitude']]
print('All Data')
print(AllData)
MeanData=AllData.groupby(['Country', 'Place_Name'])['Latitude'].mean()
StdData=AllData.groupby(['Country', 'Place_Name'])['Latitude'].std()
print('Outliers')
UpperBound=float(MeanData+StdData)
print('Higher than ', UpperBound)
OutliersHigher=AllData[AllData.Latitude>UpperBound]
print(OutliersHigher)
LowerBound=float(MeanData-StdData)
print('Lower than ', LowerBound)
OutliersLower=AllData[AllData.Latitude<LowerBound]
print(OutliersLower)
print('Not Outliers')
OutliersNot=AllData[(AllData.Latitude>=LowerBound) & (AllData.Latitude<=UpperBound)]
print(OutliersNot)
```

Output:

```

= RESTART: C:/Users/RPIMS/AppData/Local/Programs/Python/Python38-64/Python.exe
Working Base : C:/VKHCG
Loading : C:/VKHCG/01-Vermeulen/00-RawData/IP_DATA_CORE.csv
All Data
Country Place_Name Latitude
1910 GB London 51.5130
1911 GB London 51.5508
1912 GB London 51.5649
1913 GB London 51.5895
1914 GB London 51.5232
... ... ...
3434 GB London 51.5092
3435 GB London 51.5092
3436 GB London 51.5163
3437 GB London 51.5085
3438 GB London 51.5136

[1502 rows x 3 columns]
Outliers
Warning (from warnings module):
  File "C:/Users/RPIMS/AppData/Local/Programs/Python/Python38-64/Python.exe", line 17
    UpperBound=float(MeanData+StdData)
FutureWarning: Calling float on a single element Series is deprecated and will raise a TypeError in the future. Use float(ser.iloc[0]) instead.
Higher than 51.512635507867415
Country Place_Name Latitude
1910 GB London 51.5130
1911 GB London 51.5508
1912 GB London 51.5649
1913 GB London 51.5895
1914 GB London 51.5232
1916 GB London 51.5491
1919 GB London 51.5161
1920 GB London 51.5198

[1485 rows x 3 columns]
Warning (from warnings module):
  File "C:/Users/RPIMS/AppData/Local/Programs/Python/Python38-64/Python.exe", line 21
    LowerBound=float(MeanData-StdData)
FutureWarning: Calling float on a single element Series is deprecated and will raise a TypeError in the future. Use float(ser.iloc[0]) instead.
Lower than 51.506176875621264
Country Place_Name Latitude
1915 GB London 51.4739
Not Outliers
Country Place_Name Latitude
1917 GB London 51.5085
1918 GB London 51.5085
1922 GB London 51.5085
1928 GB London 51.5085
1929 GB London 51.5085
... ... ...
3432 GB London 51.5092
3433 GB London 51.5092
3434 GB London 51.5092
3435 GB London 51.5092
3437 GB London 51.5085

[1485 rows x 3 columns]

```

```
import sys
import os
import logging
import uuid
import shutil
import time

if sys.platform == 'linux':
    Base=os.path.expanduser('~') + '/VKHCG'
else:
    Base='C:/VKHCG'

sCompanies=['01-Vermeulen','02-Krennwallner','03-Hillman','04-Clark']
sLayers=['01-Retrieve','02-Assess','03-Process','04-Transform','05-Organise','06-Report']
sLevels=['debug','info','warning','error']

for sCompany in sCompanies:
    sFileDir=Base + '/' + sCompany
    if not os.path.exists(sFileDir):
        os.makedirs(sFileDir)
    for sLayer in sLayers:
        log = logging.getLogger()
        for hdlr in log.handlers[:]:
            log.removeHandler(hdlr)
        sFileDir=Base + '/' + sCompany + '/' + sLayer + '/Logging'
        if os.path.exists(sFileDir):
            shutil.rmtree(sFileDir)
        time.sleep(2)
```

```

if not os.path.exists(sFileDir):
    os.makedirs(sFileDir)
skey=str(uuid.uuid4())
sLogFile=Base + '/' + sCompany + '/' + sLayer + '/Logging/Logging_'+skey+'.log'
print('Set up:',sLogFile)
logging.basicConfig(level=logging.DEBUG,
                    format='%(asctime)s %(name)-12s %(levelname)-8s %(message)s',
                    datefmt='%m-%d %H:%M',
                    filename=sLogFile,
                    filemode='w')

console = logging.StreamHandler()
console.setLevel(logging.INFO)
formatter = logging.Formatter('%(name)-12s: %(levelname)-8s %(message)s')
console.setFormatter(formatter)
logging.getLogger('').addHandler(console)
logging.info('Practical Data Science is fun!.')
for sLevel in sLevels:
    sApp='Appllication-' + sCompany + '-' + sLayer + '-' + sLevel
    logger = logging.getLogger(sApp)
    if sLevel == 'debug':
        logger.debug('Practical Data Science logged a debugging message.')
    if sLevel == 'info':
        logger.info('Practical Data Science logged information message.')
    if sLevel == 'warning':
        logger.warning('Practical Data Science logged a warning message.')
    if sLevel == 'error':
        logger.error('Practical Data Science logged an error message.')

```

Output:

```

RESTART: D:\yukta\Datascience\Datascience\VKHCG\practical-data-science\VKHCG\77
-Yoke\Yoke_Logging.py
('Set up:', 'C:/VKHCG/01-Vermeulen/01-Retrieve/Logging/Logging_e48e608b-23f3-43e
d-885d-5eff531f83ad.log')
root      : INFO      Practical Data Science is fun!.
Appllication-01-Vermeulen-01-Retrieve-info: INFO      Practical Data Science logge
d information message.
Appllication-01-Vermeulen-01-Retrieve-warning: WARNING  Practical Data Science lo
gged a warning message.
Appllication-01-Vermeulen-01-Retrieve-error: ERROR     Practical Data Science logg
ed an error message.

```

Practical 7 - Transform-Gunnarsson_is_Born.py

```

import sys import os
from datetime import datetime from pytz
import timezone import pandas as pd
import sqlite3 as sq import uuid
pd.options.mode.chained_assignment = None if
sys.platform == 'linux':
    Base=os.path.expanduser('~') + '/VKHCG'
else:
    Base='C:/VKHCG'
print('Working Base :',Base, ' using ', sys.platform)
Company='01-Vermeulen' InputDir='00-
RawData' InputFileName='VehicleData.csv'
sDataBaseDir=Base + '/' + Company + '/04-Transform/SQLite' if not
os.path.exists(sDataBaseDir):
    os.makedirs(sDataBaseDir)
sDatabaseName=sDataBaseDir + '/Vermeulen.db' conn1 =

```

```

sq.connect(sDatabaseName)
sDataVaultDir=Base + '/88-DV'
if not os.path.exists(sDataVaultDir):
    os.makedirs(sDataVaultDir)
sDatabaseName=sDataVaultDir + '/datavault.db' conn2 =
sq.connect(sDatabaseName)
sDataWarehouseDir=Base + '/99-DW'
if not os.path.exists(sDataWarehouseDir):
    os.makedirs(sDataWarehouseDir)
sDatabaseName=sDataWarehouseDir + '/datawarehouse.db' conn3
= sq.connect(sDatabaseName)
print('Time Category') print('UTC Time')
BirthDateUTC = datetime(1960,12,20,10,15,0)
BirthDateZoneUTC=BirthDateUTC.replace(tzinfo=timezone('UTC'))
BirthDateZoneStr=BirthDateZoneUTC.strftime("%Y-%m-%d %H:%M:%S")
BirthDateZoneUTCStr=BirthDateZoneUTC.strftime("%Y-%m-%d %H:%M:%S (%Z) (%z)")
print(BirthDateZoneUTCStr)
print('Birth Date in Reykjavik :')
BirthZone = 'Atlantic/Reykjavik'
BirthDate = BirthDateZoneUTC.astimezone(timezone(BirthZone))
BirthDateStr=BirthDate.strftime("%Y-%m-%d %H:%M:%S (%Z) (%z)")
BirthDateLocal=BirthDate.strftime("%Y-%m-%d %H:%M:%S")
print(BirthDateStr)
IDZoneNumber=str(uuid.uuid4())
sDateTimeKey=BirthDateZoneStr.replace(' ','-').replace(':','-')
TimeLine=[('ZoneBaseKey', ['UTC']),
          ('IDNumber', [IDZoneNumber]),

          ('DateTimeKey', [sDateTimeKey]),
          ('UTCDateTimeValue', [BirthDateZoneUTC]), ('Zone',
          [BirthZone]),
          ('DateTimeValue', [BirthDateStr])]
TimeFrame = pd.DataFrame.from_dict(dict(TimeLine))
TimeHub=TimeFrame[['IDNumber','ZoneBaseKey','DateTimeKey','DateTimeValue']]
TimeHubIndex=TimeHub.set_index(['IDNumber'],inplace=False)
sTable = 'Hub-Time-Gunnarsson'
print('Storing :',sDatabaseName,'\n Table:',sTable)
TimeHubIndex.to_sql(sTable, conn2, if_exists="replace") sTable =
'Dim-Time-Gunnarsson' TimeHubIndex.to_sql(sTable, conn3,
if_exists="replace")
TimeSatellite=TimeFrame[['IDNumber','DateTimeKey','Zone','DateTimeValue']]
TimeSatelliteIndex=TimeSatellite.set_index(['IDNumber'],inplace=False)
BirthZoneFix=BirthZone.replace(' ','-').replace('/','-')
sTable = 'Satellite-Time-' + BirthZoneFix + '-Gunnarsson'
print('Storing :',sDatabaseName,'\n Table:',sTable)
TimeSatelliteIndex.to_sql(sTable, conn2, if_exists="replace") sTable
= 'Dim-Time-' + BirthZoneFix + '-Gunnarsson'
TimeSatelliteIndex.to_sql(sTable, conn3, if_exists="replace")
print('Person Category') FirstName =
'Guðmundur' LastName = 'Gunnarsson'
print('Name:',FirstName,LastName) print('Birth
Date:',BirthDateLocal) print('Birth
Zone:',BirthZone)
print('UTC Birth Date:',BirthDateZoneStr)
IDPersonNumber=str(uuid.uuid4())
PersonLine=[('IDNumber', [IDPersonNumber]),

```



```

        ('FirstName', [FirstName]),
        ('LastName', [LastName]),
        ('Zone', ['UTC']),
        ('DateTimeValue', [BirthDateZoneStr])) PersonFrame =
pd.DataFrame.from_dict(dict(PersonLine))
TimeHub=PersonFrame
TimeHubIndex=TimeHub.set_index(['IDNumber'],inplace=False)

sTable = 'Hub-Person-Gunnarsson'
print('Storing :',sDatabaseName,'\n Table:',sTable)
TimeHubIndex.to_sql(sTable, conn2, if_exists="replace") sTable =
'Dim-Person-Gunnarsson' TimeHubIndex.to_sql(sTable, conn3,
if_exists="replace")

```

OUTPUT

```

#####
Working Base : C:/VKHCG using win32

#####
Time Category
UTC Time
1960-12-20 10:15:00 (UTC) (+0000)
#####
Birth Date in Reykjavik :
1960-12-20 10:15:00 (GMT) (+0000)

#####
Storing : C:/VKHCG/99-DW/datawarehouse.db
Table: Hub-Time-Gunnarsson

#####
Storing : C:/VKHCG/99-DW/datawarehouse.db
Table: Satellite-Time-Atlantic-Reykjavik-Gunnarss

#####
Person Category
Name: Guðmundur Gunnarsson
Birth Date: 1960-12-20 10:15:00
Birth Zone: Atlantic/Reykjavik
UTC Birth Date: 1960-12-20 10:15:00
#####

#####
Storing : C:/VKHCG/99-DW/datawarehouse.db
Table: Hub-Person-Gunnarsson

```

Practical 7 - Transform-Gunnarsson-Sun-Model.py

```

import sys import os
from datetime import datetime from pytz
import timezone import pandas as pd
import sqlite3 as sq import uuid
pd.options.mode.chained_assignment = None
if sys.platform == 'linux': Base=os.path.expanduser('~') + '/VKHCG'
else:
    Base='C:/VKHCG'
print('Working Base :',Base, ' using ', sys.platform)
Company='01-Vermeulen'
sDataBaseDir=Base + '/' + Company + '/04-Transform/SQLite' if not
os.path.exists(sDataBaseDir):
    os.makedirs(sDataBaseDir)
sDatabaseName=sDataBaseDir + '/Vermeulen.db' conn1 =
sq.connect(sDatabaseName)
sDataWarehousetDir=Base + '/99-DW'
if not os.path.exists(sDataWarehousetDir):
    os.makedirs(sDataWarehousetDir)

```

```

sDatabaseName=sDataWarehousetDir + '/datawarehouse.db' conn2 =
sq.connect(sDatabaseName)
print('Time Dimension') BirthZone =
'Atlantic/Reykjavik'
BirthDateUTC = datetime(1960,12,20,10,15,0)
BirthDateZoneUTC=BirthDateUTC.replace(tzinfo=timezone('UTC'))
BirthDateZoneStr=BirthDateZoneUTC.strftime("%Y-%m-%d %H:%M:%S")
BirthDateZoneUTCStr=BirthDateZoneUTC.strftime("%Y-%m-%d %H:%M:%S (%Z) (%z)") BirthDate =
BirthDateZoneUTC.astimezone(timezone(BirthZone)) BirthDateStr=BirthDate.strftime("%Y-%m-%d
%H:%M:%S (%Z) (%z)") BirthDateLocal=BirthDate.strftime("%Y-%m-%d %H:%M:%S")
IDTimeNumber=str(uuid.uuid4()) TimeLine=[('TimeID',
[IDTimeNumber]),
('UTCDate', [BirthDateZoneStr]), ('LocalTime',
[BirthDateLocal]), ('TimeZone', [BirthZone])]
TimeFrame = pd.DataFrame.from_dict(dict(TimeLine))
DimTime=TimeFrame DimTimeIndex=DimTime.set_index(['TimeID'],inplace=False)
sTable = 'Dim-Time'
print('Storing :',sDatabaseName,'\n Table:',sTable)
DimTimeIndex.to_sql(sTable, conn1, if_exists="replace")
DimTimeIndex.to_sql(sTable, conn2, if_exists="replace")
print('Dimension Person') FirstName =
'Guðmundur' LastName = 'Gunnarsson'
IDPersonNumber=str(uuid.uuid4()) PersonLine=[('PersonID',
[IDPersonNumber]),
('FirstName', [FirstName]),
('LastName', [LastName]),
('Zone', ['UTC']),
('DateTimeValue', [BirthDateZoneStr])] PersonFrame =
pd.DataFrame.from_dict(dict(PersonLine))
DimPerson=PersonFrame DimPersonIndex=DimPerson.set_index(['PersonID'],inplace=False)
sTable = 'Dim-Person'
print('Storing :',sDatabaseName,'\n Table:',sTable)
DimPersonIndex.to_sql(sTable, conn1, if_exists="replace") DimPersonIndex.to_sql(sTable, conn2,
if_exists="replace")
print('Fact - Person - time')
IDFactNumber=str(uuid.uuid4()) PersonTimeLine=[('IDNumber',
[IDFactNumber]),
('IDPersonNumber', [IDPersonNumber]),
('IDTimeNumber', [IDTimeNumber])]
PersonTimeFrame = pd.DataFrame.from_dict(dict(PersonTimeLine))
FctPersonTime=PersonTimeFrame
FctPersonTimeIndex=FctPersonTime.set_index(['IDNumber'],inplace=False)
sTable = 'Fact-Person-Time'
print('Storing :',sDatabaseName,'\n Table:',sTable)
FctPersonTimeIndex.to_sql(sTable, conn1, if_exists="replace")
FctPersonTimeIndex.to_sql(sTable, conn2, if_exists="replace")

```

OUTPUT

```
#####
Working Base : C:/VKHCG  using  win32

#####
Time Dimension

#####
Storing : C:/VKHCG/99-DW/datawarehouse.db
Table: Dim-Time

#####
Dimension Person

#####
Storing : C:/VKHCG/99-DW/datawarehouse.db
Table: Dim-Person

#####
Fact - Person - time

#####
Storing : C:/VKHCG/99-DW/datawarehouse.db
Table: Fact-Person-Time
```

Practical 8 : Organize-Horizontal

```
import sys import os
import pandas as pd import sqlite3 as sq
if sys.platform == 'linux': Base=os.path.expanduser('~') + '/VKHCG'
else:
    Base='C:/VKHCG'
print('Working Base :',Base, ' using ', sys.platform)
Company='01-Vermeulen' sDataWarehouseDir=Base + '/99-DW'
if not os.path.exists(sDataWarehouseDir):
    os.makedirs(sDataWarehouseDir)
sDatabaseName=sDataWarehouseDir + '/datawarehouse.db' conn1 =
sq.connect(sDatabaseName)
sDatabaseName=sDataWarehouseDir + '/datamart.db' conn2 =
sq.connect(sDatabaseName)
sTable = 'Dim-BMI'
print('Loading :',sDatabaseName,' Table:',sTable) sSQL="SELECT *
FROM [Dim-BMI];"
PersonFrame0=pd.read_sql_query(sSQL, conn1)
sTable = 'Dim-BMI'
print('Loading :',sDatabaseName,' Table:',sTable) sSQL="SELECT
PersonID,\
    Height,\ Weight,\ bmi,\
    Indicator\
FROM [Dim-BMI]\ WHERE \
Height > 1.5 \ and Indicator = 1\
ORDER BY \
    Height,\ Weight;"
PersonFrame1=pd.read_sql_query(sSQL, conn1)
DimPerson=PersonFrame1
DimPersonIndex=DimPerson.set_index(['PersonID'],inplace=False)
```

```

sTable = 'Dim-BMI-Horizontal'
print('Storing :',sDatabaseName,'\n Table:',sTable)
DimPersonIndex.to_sql(sTable, conn2, if_exists="replace")
sTable = 'Dim-BMI-Horizontal'
print('Loading :',sDatabaseName,' Table:',sTable) sSQL="SELECT *
FROM [Dim-BMI];"
PersonFrame2=pd.read_sql_query(sSQL, conn2)
print('Full Data Set (Rows):', PersonFrame0.shape[0]) print('Full Data Set (Columns):',
PersonFrame0.shape[1]) print('Horizontal Data Set (Rows):', PersonFrame2.shape[0]) print('Horizontal Data
Set (Columns):', PersonFrame2.shape[1])

```

OUTPUT

```

#####
Working Base : C:/VKHCG using win32
#####
Loading : C:/VKHCG/99-DW/datamart.db Table: Dim-BMI
#####
Loading : C:/VKHCG/99-DW/datamart.db Table: Dim-BMI

#####
Storing : C:/VKHCG/99-DW/datamart.db
Table: Dim-BMI-Horizontal
#####
Loading : C:/VKHCG/99-DW/datamart.db Table: Dim-BMI-Horizontal
#####
Full Data Set (Rows): 1080
Full Data Set (Columns): 5
Horizontal Data Set (Rows): 194
Horizontal Data Set (Columns): 5

```

Practical 8 : Organize-Vertical

```

import sys import os
import pandas as pd import sqlite3 as sq
if sys.platform == 'linux': Base=os.path.expanduser('~') + '/VKHCG'
else:
    Base='C:/VKHCG'
print('Working Base :',Base, ' using ', sys.platform)
Company='01-Vermeulen' sDataWarehouseDir=Base + '/99-DW'
if not os.path.exists(sDataWarehouseDir):
    os.makedirs(sDataWarehouseDir)
sDatabaseName=sDataWarehouseDir + '/datawarehouse.db' conn1 =
sq.connect(sDatabaseName)
sDatabaseName=sDataWarehouseDir + '/datamart.db' conn2 =
sq.connect(sDatabaseName)
sTable = 'Dim-BMI'
print('Loading :',sDatabaseName,' Table:',sTable) sSQL="SELECT *
FROM [Dim-BMI];"
PersonFrame0=pd.read_sql_query(sSQL, conn1)
sTable = 'Dim-BMI'
print('Loading :',sDatabaseName,' Table:',sTable)
sSQL="SELECT \
    Height,\ Weight,\ Indicator\
FROM [Dim-BMI];"
PersonFrame1=pd.read_sql_query(sSQL, conn1)
DimPerson=PersonFrame1

```

```

DimPersonIndex=DimPerson.set_index(['Indicator'],inplace=False)
sTable = 'Dim-BMI-Vertical')
print('Storing :',sDatabaseName,'\n Table:',sTable)
DimPersonIndex.to_sql(sTable, conn2, if_exists="replace")
sTable = 'Dim-BMI-Vertical'
print('Loading :',sDatabaseName,' Table:',sTable) sSQL="SELECT *
FROM [Dim-BMI-Vertical];"
PersonFrame2=pd.read_sql_query(sSQL, conn2)
print('Full Data Set (Rows):', PersonFrame0.shape[0]) print('Full Data Set
(Columns):', PersonFrame0.shape[1])
print('Horizontal Data Set (Rows):', PersonFrame2.shape[0]) print('Horizontal Data
Set (Columns):', PersonFrame2.shape[1])

```

OUTPUT

```

#####
Working Base : C:/VKHCG using win32
#####
Loading : C:/VKHCG/99-DW/datamart.db Table: Dim-BMI
Loading : C:/VKHCG/99-DW/datamart.db Table: Dim-BMI
#####

#####
Storing : C:/VKHCG/99-DW/datamart.db
Table: Dim-BMI-Vertical
#####
Loading : C:/VKHCG/99-DW/datamart.db Table: Dim-BMI-Vertical
#####
Full Data Set (Rows): 1080
Full Data Set (Columns): 5
Horizontal Data Set (Rows): 1080
Horizontal Data Set (Columns): 3

```

Practical 8 : Organize-island

```

import sys import os
import pandas as pd import sqlite3 as sq
if sys.platform == 'linux': Base=os.path.expanduser('~') + '/VKHCG'
else:
    Base='C:/VKHCG'
print('Working Base :',Base, ' using ', sys.platform)
Company='01-Vermeulen' sDataWarehouseDir=Base + '/99-DW'
if not os.path.exists(sDataWarehouseDir):
    os.makedirs(sDataWarehouseDir)
sDatabaseName=sDataWarehouseDir + '/datawarehouse.db' conn1 =
sq.connect(sDatabaseName)
sDatabaseName=sDataWarehouseDir + '/datamart.db' conn2 =
sq.connect(sDatabaseName)
sTable = 'Dim-BMI'
print('Loading :',sDatabaseName,' Table:',sTable) sSQL="SELECT *
FROM [Dim-BMI];"
PersonFrame0=pd.read_sql_query(sSQL, conn1)
sTable = 'Dim-BMI'
print('Loading :',sDatabaseName,' Table:',sTable)
sSQL="SELECT \
    Height,\ Weight,\ Indicator\
FROM [Dim-BMI]\

```

```

WHERE Indicator > 2\ ORDER BY \
    Height,\ Weight;"
PersonFrame1=pd.read_sql_query(sSQL, conn1)
DimPerson=PersonFrame1
DimPersonIndex=DimPerson.set_index(['Indicator'],inplace=False)
sTable = 'Dim-BMI-Vertical'
print('Storing :',sDatabaseName,'\n Table:',sTable)
DimPersonIndex.to_sql(sTable, conn2, if_exists="replace")
sTable = 'Dim-BMI-Vertical'
print('Loading :',sDatabaseName,' Table:',sTable)
sSQL="SELECT * FROM [Dim-BMI-Vertical];"
PersonFrame2=pd.read_sql_query(sSQL, conn2)
print('Full Data Set (Rows):', PersonFrame0.shape[0]) print('Full Data Set
(Columns):', PersonFrame0.shape[1])
print('Horizontal Data Set (Rows):', PersonFrame2.shape[0]) print('Horizontal Data
Set (Columns):', PersonFrame2.shape[1])

```

OUTPUT

```

#####
Working Base : C:/VKHCG using win32
#####
Loading : C:/VKHCG/99-DW/datamart.db Table: Dim-BMI
#####
Loading : C:/VKHCG/99-DW/datamart.db Table: Dim-BMI

#####
Storing : C:/VKHCG/99-DW/datamart.db
Table: Dim-BMI-Vertical
#####
Loading : C:/VKHCG/99-DW/datamart.db Table: Dim-BMI-Vertical
#####
Full Data Set (Rows): 1080
Full Data Set (Columns): 5
Horizontal Data Set (Rows): 771
Horizontal Data Set (Columns): 3

```

Practical 8 : Organize-secure-vault

```

import sys import os
import pandas as pd import sqlite3 as sq
if sys.platform == 'linux': Base=os.path.expanduser('~') + '/VKHCG'
else:
    Base='C:/VKHCG'
print('Working Base :',Base, ' using ', sys.platform)
Company='01-Vermeulen' sDataWarehouseDir=Base + '/99-DW'
if not os.path.exists(sDataWarehouseDir):
    os.makedirs(sDataWarehouseDir)

sDatabaseName=sDataWarehouseDir + '/datawarehouse.db' conn1 =
sq.connect(sDatabaseName)
sDatabaseName=sDataWarehouseDir + '/datamart.db' conn2 =
sq.connect(sDatabaseName)
sTable = 'Dim-BMI'
print('Loading :',sDatabaseName,' Table:',sTable) sSQL="SELECT *
FROM [Dim-BMI];"
PersonFrame0=pd.read_sql_query(sSQL, conn1)
sTable = 'Dim-BMI'
print('Loading :',sDatabaseName,' Table:',sTable)
sSQL="SELECT \
    Height,\ Weight,\ Indicator,\
    CASE Indicator\

```

```

    WHEN 1 THEN 'Pip'\
    WHEN 2 THEN 'Norman'\ WHEN 3 THEN
    'Grant'\ ELSE 'Sam'\
    END AS Name\ FROM [Dim-BMI]\
WHERE Indicator > 2\ ORDER BY \
    Height,\ Weight;"

```

```

PersonFrame1=pd.read_sql_query(sSQL, conn1)
DimPerson=PersonFrame1
DimPersonIndex=DimPerson.set_index(['Indicator'],inplace=False)
sTable = 'Dim-BMI-Secure'
print('Storing :',sDatabaseName,'\n Table:',sTable) DimPersonIndex.to_sql(sTable,
conn2, if_exists="replace")
sTable = 'Dim-BMI-Secure'
print('Loading :',sDatabaseName,' Table:',sTable)
sSQL="SELECT * FROM [Dim-BMI-Secure] WHERE Name = 'Sam';"
PersonFrame2=pd.read_sql_query(sSQL, conn2)
print('Full Data Set (Rows):', PersonFrame0.shape[0]) print('Full Data Set
(Columns):', PersonFrame0.shape[1])
print('Horizontal Data Set (Rows):', PersonFrame2.shape[0]) print('Horizontal Data
Set (Columns):', PersonFrame2.shape[1]) print('Only Sam Data')
print(PersonFrame2.head())

```

OUTPUT

```

#####
Working Base : C:/VKHCG using win32
#####
Loading : C:/VKHCG/99-DW/datamart.db Table: Dim-BMI
Loading : C:/VKHCG/99-DW/datamart.db Table: Dim-BMI

#####
Storing : C:/VKHCG/99-DW/datamart.db
Table: Dim-BMI-Secure
#####
Loading : C:/VKHCG/99-DW/datamart.db Table: Dim-BMI-Secure
#####
Full Data Set (Rows): 1080
Full Data Set (Columns): 5
Horizontal Data Set (Rows): 692
Horizontal Data Set (Columns): 4
Only Sam Data
   Indicator  Height  Weight Name
0           4      1.0      35  Sam
1           4      1.0      40  Sam
2           4      1.0      45  Sam
3           4      1.0      50  Sam
4           4      1.0      55  Sam

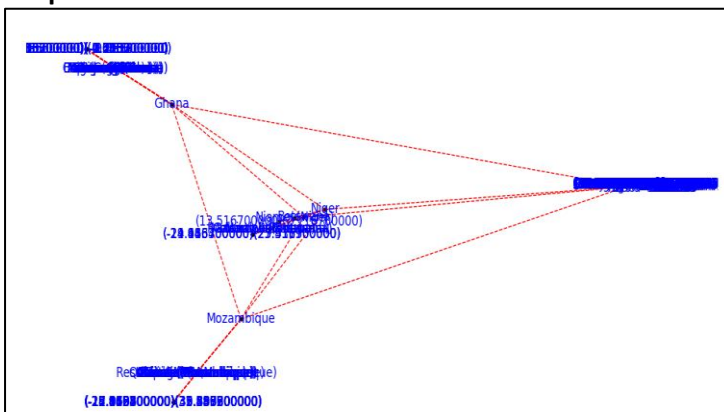
```


Practical No.9

A) Generating Reports

```
import sys
import os
import pandas as pd
import networkx as nx
import matplotlib.pyplot as plt
pd.options.mode.chained_assignment = None
if sys.platform == 'linux':
    Base=os.path.expanduser('~') + 'VKHCG'
else:
    Base='C:/VKHCG'
sInputFileName='02-Assess/01-EDS/02-Python/Assess-Network-Routing-Customer.csv'
sOutputFileName1='06-Report/01-EDS/02-Python/Report-Network-Routing-Customer.gml'
sOutputFileName2='06-Report/01-EDS/02-Python/Report-Network-Routing-Customer.png'
Company='01-Vermeulen'
sFileName=Base + '/' + Company + '/' + sInputFileName
CustomerDataRow=pd.read_csv(sFileName, header=0, low_memory=False, encoding="latin-1")
CustomerData=CustomerDataRow.head(100)
G=nx.Graph()
for i in range(CustomerData.shape[0]):
    for j in range(CustomerData.shape[0]):
        Node0=CustomerData['Customer_Country_Name'][i]
        Node1=CustomerData['Customer_Country_Name'][j]
        if Node0 != Node1:
            G.add_edge(Node0,Node1)
for i in range(CustomerData.shape[0]):
    Node0=CustomerData['Customer_Country_Name'][i]
    Node1=CustomerData['Customer_Place_Name'][i] + '(' + CustomerData['Customer_Country_Name'][i] + ')'
    Node2='(' + "{:.9f}".format(CustomerData['Customer_Latitude'][i]) + ')(' +
"{:.9f}".format(CustomerData['Customer_Longitude'][i]) + ')'
    if Node0 != Node1:
        G.add_edge(Node0,Node1)
    if Node1 != Node2:
        G.add_edge(Node1,Node2)
sFileName=Base + '/' + Company + '/' + sOutputFileName1
nx.write_gml(G, sFileName)
sFileName=Base + '/' + Company + '/' + sOutputFileName2
plt.figure(figsize=(25, 25))
pos=nx.spectral_layout(G, dim=2)
nx.draw_networkx_nodes(G, pos, node_color='k', node_size=10, alpha=0.8)
nx.draw_networkx_edges(G, pos, edge_color='r', arrows=False, style='dashed')
nx.draw_networkx_labels(G, pos, font_size=12, font_family='sans-serif', font_color='b')
plt.axis('off')
plt.savefig(sFileName, dpi=600)
plt.show()
```

Output:



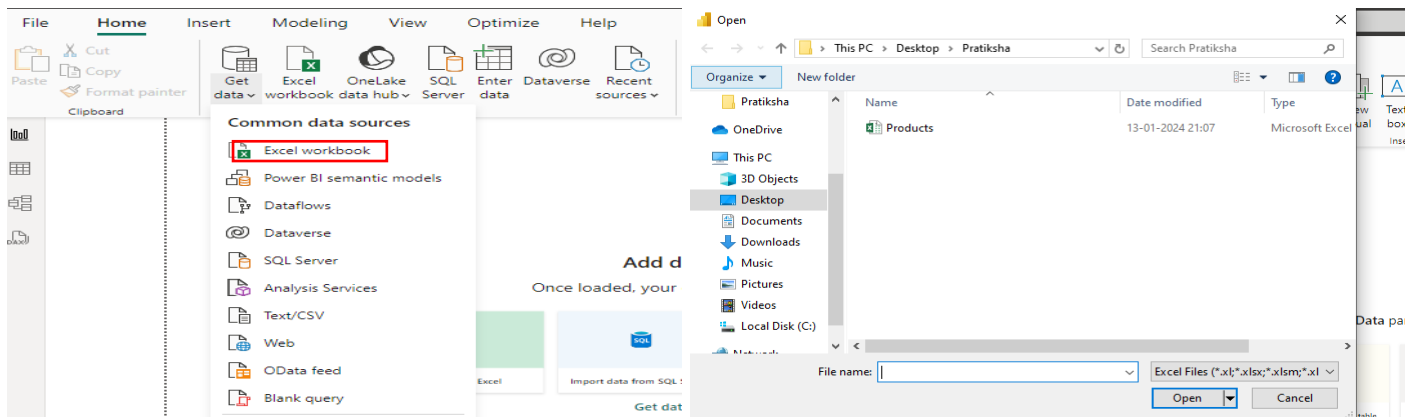
Practical No. 10

Data Visualization with power bi

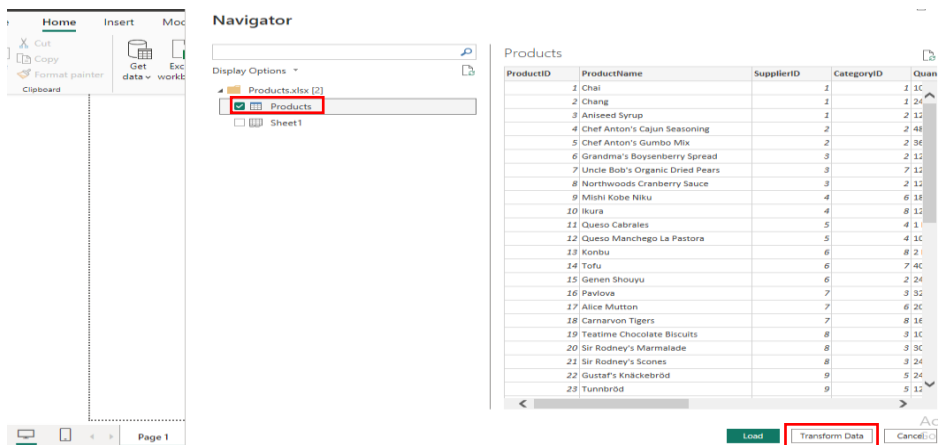
Case Study : Sales Data

Step 1: connect to an Excel Workbook

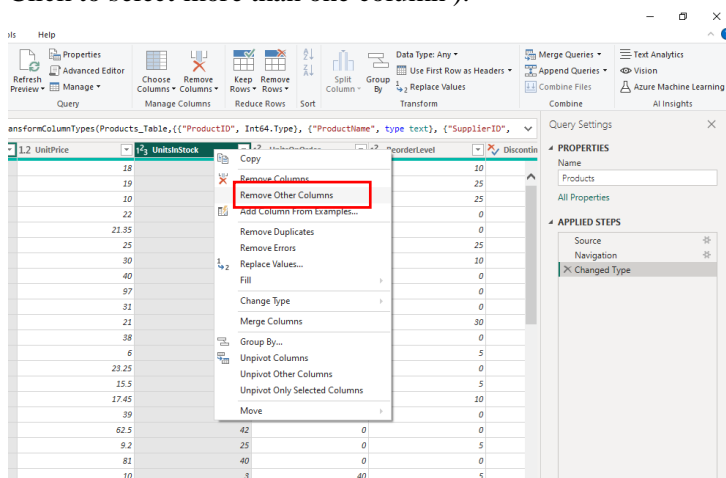
1. Launch power Bi Desktop.
2. From the Home Ribbon, Select Get Data → Select Excel Workbook .
3. In the Open File Dialog Box, Select the Product.xlsx file.



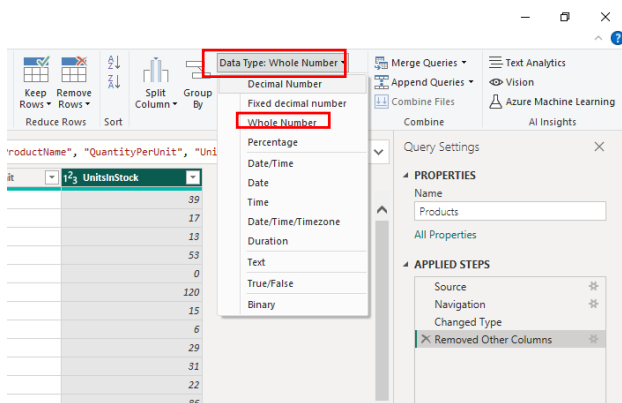
4. In the Open File Dialog Box, Select the Product.xlsx file.



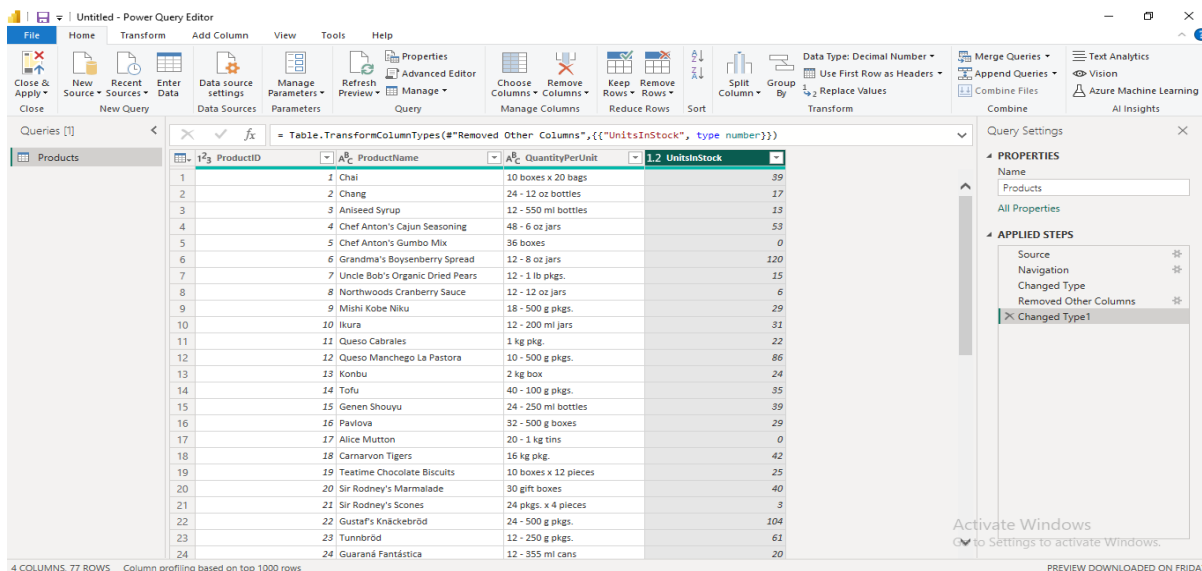
5. Click on Products Check Box. & We Will see the product Table. Select The Transform Data.
6. In Query Editor , Select the ProductID, ProductName, QuantityPerUnit, and unitsInStock Columns. (Use Ctrl + Click to select more than one column).



7. Right Click on Column Header and Click Remove Other Columns.

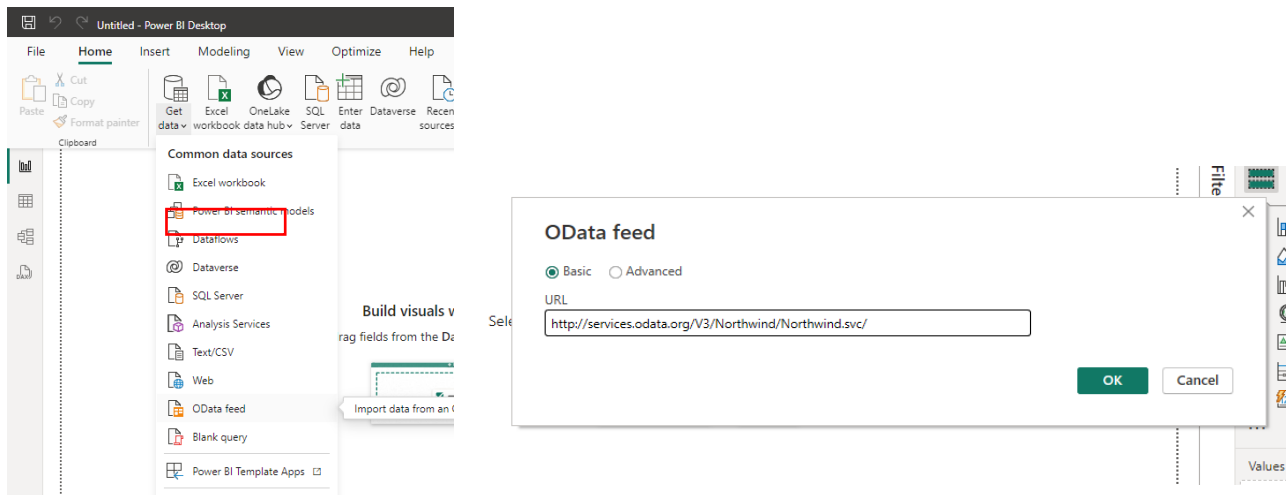


8. Select Close & Apply from Home Ribbon.

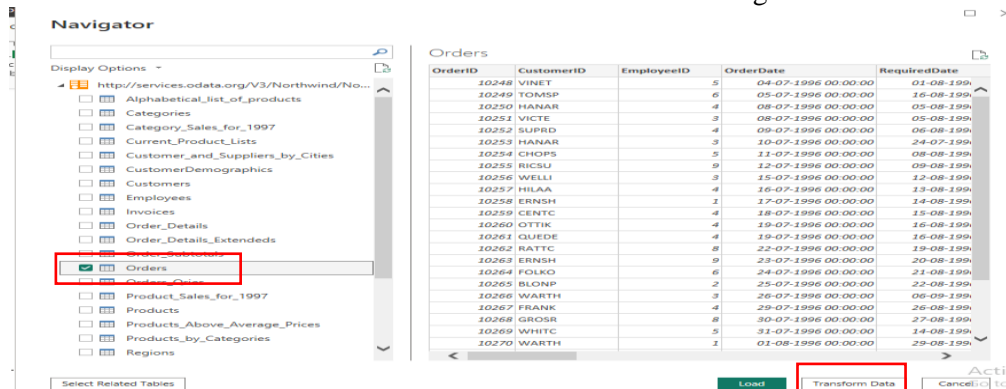


9. Another Window on → Select Get Data & Select the OData feed. And Copy the link given below. & Paste it to OData feed URL Box and Click Ok.

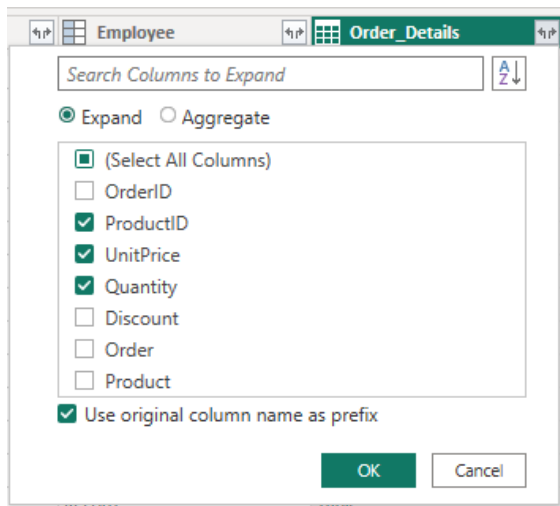
<http://services.odata.org/V3/Northwind/Northwind.svc/>



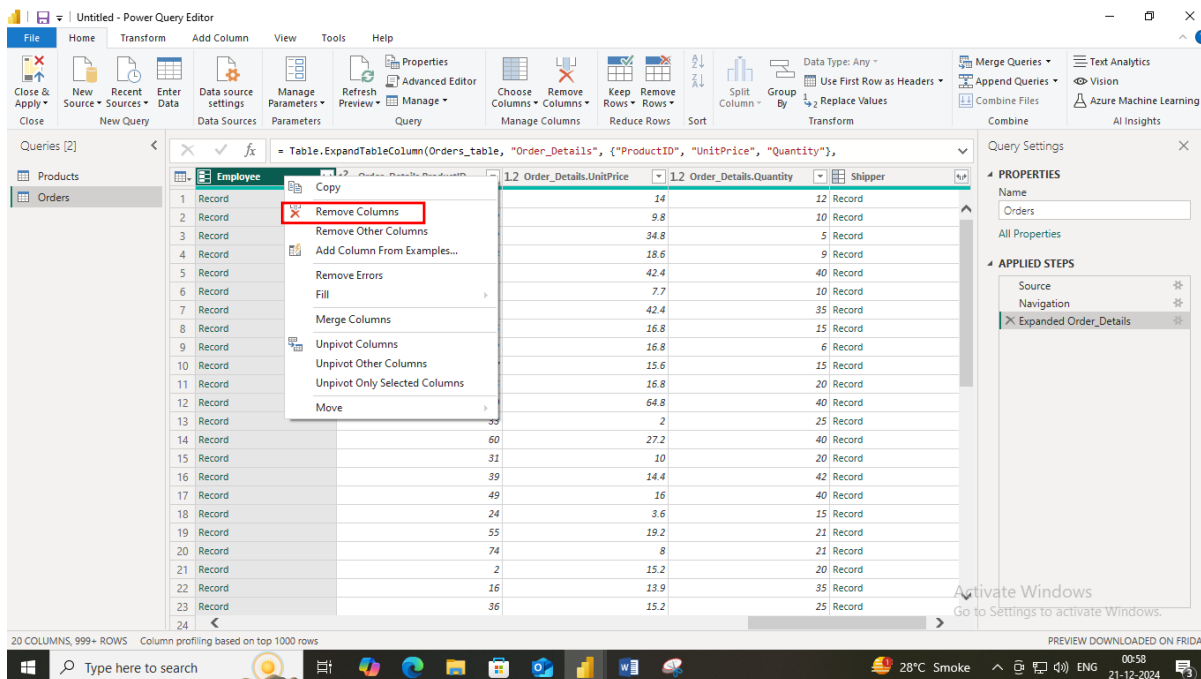
10. We will see the below screen like this then select from navigator orders checkbox. & Click on Transform Data.



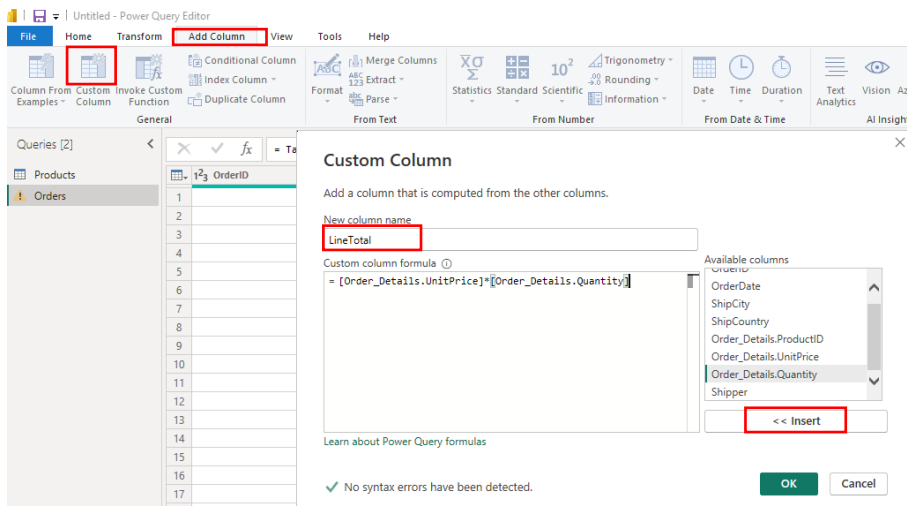
11. Expand the Order_Details column & select the ProductID, UnitPrice, Quantity & Click OK.



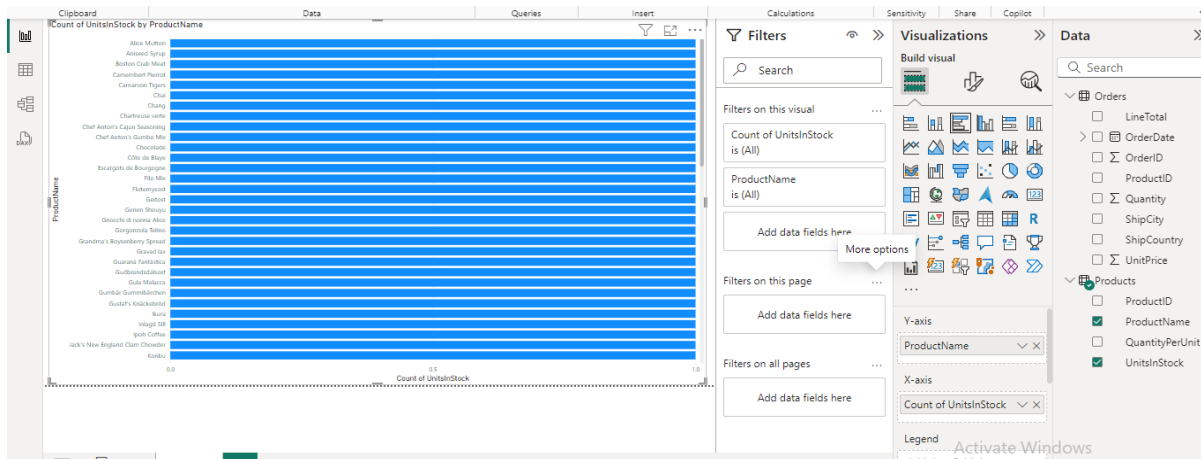
12. Remove other Column to only display column of interest, In this step you remove all Column except OrderID, OrderDate, ShipCity, Order_Details.ProductID, Order_Details.Unitprice, Order_Details.Quantity, Shipper columns. & Remove Columns.



13. From Add Column Ribbon Select Custom Column. Add New Name in new Column name LineTotal. From Available column Select order_Details.Unitprice and Click insert Add “*” and select Order_Details.Quantity and insert → Ok. We Will see the a New Column Name LineTotal Appears.



14. In Query Editor, drag the LineTotal Column to the left , After ShipCountry. → Double Click on Order_Details.ProductID, Order_Details.Unitprice, Order_Details.Quantity change name to Only ProductID, Unitprice, Quantity.
15. From Home Ribbon , Select Close and apply. We Will get new Window of Power Bi. Select From Data Paneel From Products select ProductName And UnitInStock. If output is not seen then Change X-axis and Y-axis from Visualizations.



16. For Orders Select Map from Visualizations. And From data Column Select From Orders Select LineTotal And ShipCity.

