



**Department of Computer Science and Engineering (Data Science)**  
**Subject: Time Series Analysis**

**Jhanvi Parekh**

**60009210033**

**D11**

**Experiment 3**

**(Data Wrangling and Preparation for Time Series Data)**

**Aim:** Implementation of Data Wrangling and Preparation for Time Series Data.

**Theory:**

**Data Wrangling:**

Data wrangling is the art of transformation and mapping raw data into valuable information with the help of preprocessing strategies. It aims to provide rich data that can be utilized to gain maximum insights. It comprises operations such as loading, imputing, applying transformations, treating outliers, cleaning, integrating, dealing with inconsistency, reducing dimensionality, and engineering features. Data wrangling is an integral and main part of machine learning modeling.

Real-world data is undoubtedly messy, and it is not reasonable to use data directly for modeling without performing some wrangling.

**Loading Data into Pandas:**

Pandas is the most notable framework for data wrangling. It includes data manipulation and analysis tools intended to make data analysis rapid and convenient. In the real world, data is generated via various tools and devices; hence, it comes in different formats such as CSV, Excel, and JSON. In addition, sometimes data needs to be read from a URL. All the data comprises several records and variables.

**Loading Data Using CSV:**

This dataset depicts the number of female births over time. Pandas has a built-in function to read CSV files. If files have different separators, use sep = ' ' and fill in the separator, as in (;, |\t, ' '). The following code imports the dataset:

```
import pandas as pd df = pd.read_csv(r'daily-total-female-births-CA.csv') df.head(5)
# sep can be like ; |
```

**Output:**



### Department of Computer Science and Engineering (Data Science)

	date	births
0	1959-01-01	35
1	1959-01-02	32
2	1959-01-03	30
3	1959-01-04	31
4	1959-01-05	44

#### Loading Data Using Excel:

The Istanbul Stock Exchange dataset comprises the returns of the Istanbul Stock Exchange along with seven other international indices (SP, DAX, FTSE, NIKKEI, BOVESPA, MSCE\_EU, MSCI\_EM) from June 5, 2009, to February 22, 2011. The data is organized by workday. This data is available in Excel format, and Pandas has a built-in function to read the Excel file. The following code imports the dataset:

```
import pandas as pd dfExcel = pd.read_excel(r'istambul_stock_exchange.xlsx', sheet_name = 'Data') dfExcel.head(5)
```

#### Output:

	date	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EM
0	2009-01-05	0.035754	0.038376	-0.004679	0.002193	0.003894	0.000000	0.031190	0.012698	0.028524
1	2009-01-06	0.025426	0.031813	0.007787	0.008455	0.012866	0.004162	0.018920	0.011341	0.008773
2	2009-01-07	-0.028862	-0.026353	-0.030469	-0.017833	-0.028735	0.017293	-0.035899	-0.017073	-0.020015
3	2009-01-08	-0.062208	-0.084716	0.003391	-0.011726	-0.000466	-0.040061	0.028283	-0.005561	-0.019424
4	2009-01-09	0.009860	0.009658	-0.021533	-0.019873	-0.012710	-0.004474	-0.009764	-0.010989	-0.007802

#### Loading Data Using JSON:

This example dataset is in JSON format. The following code imports the dataset:

```
dfJson = pd.read_json(r'test.json') dfJson.head(5)
```

#### Output:

	Names	Age
0	John	33
1	Sal	45
2	Tim	22
3	Rod	54

#### Loading Data from a URL:



### Department of Computer Science and Engineering (Data Science)

The Abalone dataset comprises 4,177 observations and 9 variables. It is not in any file structure; instead, we can display it as text at the specific URL. The following code imports the dataset:

```
dfURL = pd.read_csv(r'https://archive.ics.uci.edu/ml/machine-
learningdatabases/abalone/abalone.data', names =['Sex',
                                              'Length', 'Diameter', 'Height','Whole weight',
                                              'Shucked weight','Viscera weight', 'Shell weight', 'Rings']) dfURL.head(5)
```

#### Output:

Sex	Length	Diameter	Height	Whole weight	Shucked weight	Viscera weight	Shell weight	Rings
M	0.455	0.365	0.095	0.5140	0.2245	0.1010	0.150	15
M	0.350	0.265	0.090	0.2255	0.0995	0.0485	0.070	7
F	0.530	0.420	0.135	0.6770	0.2565	0.1415	0.210	9
M	0.440	0.365	0.125	0.5160	0.2155	0.1140	0.155	10
I	0.330	0.255	0.080	0.2050	0.0895	0.0395	0.055	7

#### Exploring Pandasql and Pandas:

Pandasql is a Python framework for running SQL queries on Pandas DataFrames. It has plenty of essential attributes and a similar purpose as the sqldf framework in R. It helps us to query Pandas DataFrames using SQL syntax. It provides a SQL interface to perform data wrangling on a Pandas DataFrame. Pandasql helps analysts and data engineers who are masters of SQL to transition into Python (Pandas) smoothly.

Use pip install pandasql to install the library.

#### Selecting the Top Five Records:

With the help of the pandas.head() function, we can fetch the first N records from the dataset. We can do the same operation with the help of Pandassql. The example illustrates how to do it with Pandas and Pandasql.

Pandas:

```
import pandas as pd dfp =
pd.read_excel(r'Absenteeism_at_work.xls')
dfp.head(5)
```

#### Output:

ID	Reason_for_absence	Month_of_absence	Day_of_the_week	Seasons	Transportation_expense	Distance_from_Residence_to_Work	Service_time	Age
0	11	26	7	3	1	289	36	13 33
1	36	0	7	3	1	118	13	18 50
2	3	23	7	4	1	179	51	18 38
3	7	7	7	5	1	279	5	14 39
4	11	23	7	5	1	289	36	13 33

5 rows × 21 columns

Pandasql:



### Department of Computer Science and Engineering (Data Science)

```
from pandasql import sqldf
dfpsql = pd.read_excel(r'Absenteeism_at_work.xls') Query_string = """
select * from dfpsql limit 5 """ sqldf(Query_string, globals())
```

#### Output:

ID	Reason_for_absence	Month_of_absence	Day_of_the_week	Seasons	Transportation_expense	Distance_from_Residence_to_Work	Service_time	Age
0	11	26	7	3	1	289	36	13 33
1	36	0	7	3	1	118	13	18 50
2	3	23	7	4	1	179	51	18 38
3	7	7	7	5	1	279	5	14 39
4	11	23	7	5	1	289	36	13 33

5 rows × 21 columns

#### Applying a Filter:

Data filtering is a significant part of data preprocessing. With filtering, we can choose a smaller partition of the dataset and use that subset for viewing and munging; we need specific criteria or a rule to filter the data. This is also known as **subsetting data or drill-down data**. The following example illustrates how to apply a filter with Pandas and Pandasql.

#### Pandas:

```
dfp[(dfp['Age'] >=30) & (dfp['Age'] <=45)]
```

#### Output:

ID	Reason_for_absence	Month_of_absence	Day_of_the_week	Seasons	Transportation_expense	Distance_from_Residence_to_Work	Service_time	A
0	11	26	7	3	1	289	36	13
2	3	23	7	4	1	179	51	18
3	7	7	7	5	1	279	5	14
4	11	23	7	5	1	289	36	13
5	3	23	7	6	1	179	51	18
7	20	23	7	6	1	260	50	11
8	14	19	7	2	1	155	12	14
9	1	22	7	2	1	235	11	14
10	20	1	7	2	1	260	50	11
11	20	1	7	3	1	260	50	11

#### Pandasql:

```
Query_string = """ select * from dfpsql where age>=30 and age<=45 """
sqldf(Query_string, globals())
```

#### Output:



### Department of Computer Science and Engineering (Data Science)

ID	Reason_for_absence	Month_of_absence	Day_of_the_week	Seasons	Transportation_expense	Distance_from_Residence_to_Work	Service_time
0	11	26	7	3	1	289	36
1	3	23	7	4	1	179	51
2	7	7	7	5	1	279	5
3	11	23	7	5	1	289	36
4	3	23	7	6	1	179	51
5	20	23	7	6	1	260	50
6	14	19	7	2	1	155	12
7	1	22	7	2	1	235	11
8	20	1	7	2	1	260	50
9	20	1	7	3	1	260	50

#### Distinct (Unique):

Several duplicate records exist in the dataset. If we want to select the number of unique values for the specific variable, then we can use the unique() function of Pandas. The following example illustrates how to do this with Pandas and Pandasql.

Pandas

```
dfp['ID'].unique()
```

#### Output:

```
array([11, 36, 3, 7, 10, 20, 14, 1, 24, 6, 33, 18, 30, 2, 19, 27, 34,
      5, 15, 29, 28, 13, 22, 17, 31, 23, 32, 9, 26, 21, 8, 25, 12, 16,
      4, 35], dtype=int64)
```

Pandasql:

```
Query_string = """ select distinct ID from dfpsql;"""
sqldf(Query_string, globals())
```

#### Output:

ID
0 11
1 36
2 3
3 7
4 10
5 20
6 14
7 1
8 24
9 6
10 33



## Department of Computer Science and Engineering (Data Science)

### IN:

Data filtering is a process of extracting essential data from a dataset by using some condition. There are several methods to do filtering on a dataset. Sometimes we want to investigate whether the data has been associated with a particular DataFrame or Series. In such an event, we can use Pandas' `isin()` function, which checks whether values are present in the sequence. The procedure can also be carried out in Pandasql. The following example illustrates how to do it with Pandas and Pandasql.

Pandas:

```
dfp[dfp.Age.isin([20,30,40])]
```

### Output:

ID	Reason_for_absence	Month_of_absence	Day_of_the_week	Seasons	Transportation_expense	Distance_from_Residence_to_Work	Service_time
47	15	23	9	5	1	291	31 12
49	15	14	9	2	4	291	31 12
65	22	23	10	5	4	179	26 9
71	15	23	10	5	4	291	31 12
75	15	14	10	3	4	291	31 12
83	17	21	11	5	4	179	22 17
84	15	23	11	5	4	291	31 12
87	15	14	11	2	4	291	31 12
91	17	21	11	4	4	179	22 17
97	15	23	11	5	4	291	31 12

Pandasql:

```
Query_string = """ select * from dfpsql where Age in(20,30,40);"""
sqldf(Query_string, globals())
```

### Output:

ID	Reason_for_absence	Month_of_absence	Day_of_the_week	Seasons	Transportation_expense	Distance_from_Residence_to_Work	Service_time
0	15	23	9	5	1	291	31 12
1	15	14	9	2	4	291	31 12
2	22	23	10	5	4	179	26 9
3	15	23	10	5	4	291	31 12
4	15	14	10	3	4	291	31 12
5	17	21	11	5	4	179	22 17
6	15	23	11	5	4	291	31 12
7	15	14	11	2	4	291	31 12
8	17	21	11	4	4	179	22 17
9	15	23	11	5	4	291	31 12

### NOT IN:

The NOT IN operation is used for a similar purpose as explained earlier. If we want to check whether a value is not part of a sequence, we can use the tilde (~) symbol to perform a NOT IN operation.

Pandas:

```
dfp[~dfp.Age.isin([20,30,40])]
```



## Department of Computer Science and Engineering (Data Science)

### Output:

ID	Reason_for_absence	Month_of_absence	Day_of_the_week	Seasons	Transportation_expense	Distance_from_Residence_to_Work	Service_time
0	11	26	7	3	1	289	36
1	36	0	7	3	1	118	13
2	3	23	7	4	1	179	51
3	7	7	7	5	1	279	5
4	11	23	7	5	1	289	36
5	3	23	7	6	1	179	51
6	10	22	7	6	1	361	52
7	20	23	7	6	1	260	50
8	14	19	7	2	1	155	12
9	1	22	7	2	1	235	11

Pandasql:

```
Query_string = """ select * from dfpsql where Age not in(20,30,40);"""
sqldf(Query_string, globals())
```

### Output:

ID	Reason_for_absence	Month_of_absence	Day_of_the_week	Seasons	Transportation_expense	Distance_from_Residence_to_Work	Service_time
0	11	26	7	3	1	289	36
1	36	0	7	3	1	118	13
2	3	23	7	4	1	179	51
3	7	7	7	5	1	279	5
4	11	23	7	5	1	289	36
5	3	23	7	6	1	179	51
6	10	22	7	6	1	361	52
7	20	23	7	6	1	260	50
8	14	19	7	2	1	155	12
9	1	22	7	2	1	235	11

### Ascending Data Order:

ORDER BY sorts the result-set in ascending or descending order based on the value selected. Pandas has a sort\_value() function that can use different sorting algorithms, such as quicksort, mergesort, and heapsort. The default value is ascending order, whose Boolean value is True. Except for that axis-wise, we do perform sorting such as 0 for index and 1 for columns. SQL has an ORDER BY clause to perform a similar sorting operation.

Pandas:

```
dfp.sort_values(by = ['Age','Service_time'], ascending= True)
```

### Output:



## Department of Computer Science and Engineering (Data Science)

ID	Reason_for_absence	Month_of_absence	Day_of_the_week	Seasons	Transportation_expense	Distance_from_Residence_to_Work	Service_time
40	27	23	9	3	1	184	42
118	27	23	1	5	2	184	42
132	27	23	1	5	2	184	42
137	27	23	2	6	2	184	42
149	27	23	2	3	2	184	42
209	27	7	5	4	3	184	42
269	27	6	8	4	1	184	42
6	10	22	7	6	1	361	52
22	10	13	8	2	1	361	52
25	10	25	8	2	1	361	52

Pandasql:

```
Query_string = """ select * from dfpsql order by Age,Service_time;"""
sqldf(Query_string, globals())
```

**Output:**

ID	Reason_for_absence	Month_of_absence	Day_of_the_week	Seasons	Transportation_expense	Distance_from_Residence_to_Work	Service_time
0	27	23	9	3	1	184	42
1	27	23	1	5	2	184	42
2	27	23	1	5	2	184	42
3	27	23	2	6	2	184	42
4	27	23	2	3	2	184	42
5	27	7	5	4	3	184	42
6	27	6	8	4	1	184	42
7	10	22	7	6	1	361	52
8	10	13	8	2	1	361	52
9	10	25	8	2	1	361	52

**Descending Data Order:**

As mentioned, the sort\_value() function can sort results in ascending or descending order. It needs the parameter ascending = False to sort values in descending order. You can also update some other parameters, as explained for ORDER BY ascending.

Pandas:

```
dfp.sort_values(by = ['Age','Service_time'], ascending= False)
```

**Output:**



### Department of Computer Science and Engineering (Data Science)

ID	Reason_for_absence	Month_of_absence	Day_of_the_week	Seasons	Transportation_expense	Distance_from_Residence_to_Work	Service_time	Age
9	18	8	3	1	228		14	16
9	18	5	4	3	228		14	16
9	1	10	4	4	228		14	16
9	25	3	3	2	228		14	16
9	12	3	3	2	228		14	16
9	25	3	4	2	228		14	16
9	6	7	2	1	228		14	16
9	6	7	3	1	228		14	16
35	0	0	6	3	179		45	14
36	0	7	3	1	118		13	18

Pandasql:

```
Query_string = """ select * from dfpsql order by Age Desc,Service_time Desc;"""
sqldf(Query_string, globals())
```

**Output:**

ID	Reason_for_absence	Month_of_absence	Day_of_the_week	Seasons	Transportation_expense	Distance_from_Residence_to_Work	Service_time	Age
9	18	8	3	1	228		14	16 58
9	18	5	4	3	228		14	16 58
9	1	10	4	4	228		14	16 58
9	25	3	3	2	228		14	16 58
9	12	3	3	2	228		14	16 58
9	25	3	4	2	228		14	16 58
9	6	7	2	1	228		14	16 58
9	6	7	3	1	228		14	16 58
35	0	0	6	3	179		45	14 53
36	0	7	3	1	118		13	18 50

**Aggregation:**

Aggregation is the process of mining data so the data can be investigated, collected, and presented in a summarized manner. It allows you to perform a calculation on single or multiple vectors, usually accompanied by the group by() functions in Pandas. It performs numerous operations such as splitting, applying, and combining. This example illustrates the aggregation operation in Pandas and Pandasql. For the following example, we are leveraging .agg in Pandas to achieve the required result.

Pandas:

```
dfp.agg({'Transportation_expense': ['count', 'min', 'max', 'mean']})
```

**Output:**

Transportation_expense	
count	740.000000
min	118.000000
max	388.000000
mean	221.32973



### Department of Computer Science and Engineering (Data Science)

Pandasql:

```
Query_string = """ select count(Transportation_expense) as count,
min(Transportation_expense) as min, max(Transportation_expense) as max,
avg(Transportation_expense) as mean from dfp;"""
sqldf(Query_string,
globals())
```

**Output:**

	count	min	max	mean
0	740	118	388	221.32973

### GROUP BY:

We can use GROUP BY to arrange identical data into groups based on the aggregation function used. In other words, it groups rows that have similar values into summary rows. We perform this operation with the groupby() function in Pandas. SQL has its GROUP BY clause to perform a similar operation. The example illustrates the GROUP BY operation in Pandas and Pandasql.

Pandas:

```
dfp.groupby('ID')[['Service_time']].sum()
```

**Output:**

ID	
1	322
2	72
3	2034
4	13
5	247
6	104
7	84
8	28
9	128
10	72
11	520
12	7
13	180
14	406
15	444
16	48
17	340

Pandasql:

```
Query_string = """ select ID , sum(Service_time) as Sum_Service_time from dfp
group by ID;"""
sqldf(Query_string, globals())
```

**Output:**



### Department of Computer Science and Engineering (Data Science)

ID Sum\_Service\_time

1	322
2	72
3	2034
4	13
5	247
6	104
7	84
8	28
9	128
10	72
11	520

#### GROUP BY with Aggregation:

The groupby() function provides a group of similar data based on a selected feature with that data; we can perform different aggregation operations with the .agg() function on other features in Pandas. In SQL, we use GROUP BY to apply aggregate functions on groups of data returned from a query. FILTER is a modifier used with an aggregate function to bound the values used in an aggregation.

Pandas:

```
dfp.groupby('Reason_for_absence').agg({'Age': ['mean', 'min', 'max']})
```

#### Output:

Reason_for_absence	Age			
	mean	min	max	
0	39.604651	28	53	
1	37.687500	28	58	
2	28.000000	28	28	
3	40.000000	40	40	
4	45.000000	41	49	
5	41.666667	37	50	
6	38.500000	27	58	
7	32.866667	27	46	
8	36.500000	28	40	



### Department of Computer Science and Engineering (Data Science)

Pandasql:

```
Query_string = """ select Reason_for_absence , avg(Age) as mean, min(Age) as min, max(Age) as max from dfp group by Reason_for_absence;"""
sqldf(Query_string, globals())
```

#### Output:

Reason_for_absence	mean	min	max
0	39.604651	28	53
1	37.687500	28	58
2	28.000000	28	28
3	40.000000	40	40
4	45.000000	41	49
5	41.666667	37	50
6	38.500000	27	58
7	32.866667	27	46
8	36.500000	28	40

#### Join (Merge):

The terms join and merge mean the same thing in Pandas, Python, and other languages like SQL and R. In reality, their fundamental operation is equivalent, but their way of executing an operation is different. A join in Pandas utilizes an index to consolidate two data sources, while a merge looks for overlapping columns to merge. In Pandas, the merge() and join() functions are used. In SQL, the MERGE and JOIN operators perform the same operation as in Pandas. This example illustrates the aggregation operation in Pandas and Pandasql. We've created Sample Employee and Department tables to illustrate the join examples.

Pandas:

```
import pandas as pd
data1 = {
    'Empid': [1011, 1012, 1013, 1014, 1015],
    'Name': ['John', 'Rahul', 'Rick', 'Morty', 'Tim'],
    'Designation': ['Manager', 'Research Engineer', ' Research Engineer',
'VP', 'Delivery Manager'],
    'Date_of_joining': ['01-Jan-2000', '23-sep-2006', '11-Jan-2012', '21-
Jan-1991', '12-Jan-1990']}
Emp_df = pd.DataFrame(data1, columns = ['Empid', 'Name',
'Designation','Date_of_joining'])
Emp_df
```

#### Output:



### Department of Computer Science and Engineering (Data Science)

Empid	Name	Designation	Date_of_joining
1011	John	Manager	01-Jan-2000
1012	Rahul	Research Engineer	23-sep-2006
1013	Rick	Research Engineer	11-Jan-2012
1014	Morty	VP	21-Jan-1991
1015	Tim	Delivery Manager	12-Jan-1990

Pandasql:

```
data2 = {
    'Empid': [1011, 1017, 1013, 1019, 1015],
    'Deptartment': ['Management', 'Research', 'Research', 'Management',
'Delivery'],
    'Total_Experience': [18, 10, 10, 28, 22]}
Dept_df = pd.DataFrame(data2, columns = ['Empid', 'Deptartment',
'Total_Experience'])
Dept_df
```

**Output:**

Empid	Deptartment	Total_Experience
1011	Management	18
1017	Research	10
1013	Research	10
1019	Management	28
1015	Delivery	22

**INNER JOIN:**

An inner merge/inner join keeps rows where the merge value contains an “on” value in both the DataFrames. It selects the records that have matching values (joining columns) in one or more tables.

Pandas:

```
pd.merge(Emp_df, Dept_df, left_on='Empid',right_on='Empid', how='inner')
```

**Output:**

Empid	Name	Designation	Date_of_joining	Deptartment	Total_Experience
1011	John	Manager	01-Jan-2000	Management	18
1013	Rick	Research Engineer	11-Jan-2012	Research	10
1015	Tim	Delivery Manager	12-Jan-1990	Delivery	22



### Department of Computer Science and Engineering (Data Science)

Pandasql:

```
Query_string = """ select * from Emp_df a INNER JOIN Dept_df b ON a.Empid = b.Empid;"""
sqldf(Query_string, globals())
```

#### Output:

Empid	Name	Designation	Date_of_joining	Empid	Deptartment	Total_Experience
1011	John	Manager	01-Jan-2000	1011	Management	18
1013	Rick	Research Engineer	11-Jan-2012	1013	Research	10
1015	Tim	Delivery Manager	12-Jan-1990	1015	Delivery	22

#### LEFT JOIN:

A left merge or left join keeps a row on the left DataFrame when the missing value finds the “on” variable in the right DataFrame and adds the empty or NaN value in that result. All the records from the left table and selected matching records based on the joining condition from one or more tables.

Pandas:

```
pd.merge(Emp_df, Dept_df, left_on='Empid', right_on='Empid', how='left')
```

#### Output:

Empid	Name	Designation	Date_of_joining	Deptartment	Total_Experience
1011	John	Manager	01-Jan-2000	Management	18.0
1012	Rahul	Research Engineer	23-sep-2006	NaN	NaN
1013	Rick	Research Engineer	11-Jan-2012	Research	10.0
1014	Morty	VP	21-Jan-1991	NaN	NaN
1015	Tim	Delivery Manager	12-Jan-1990	Delivery	22.0

Pandasql:

```
Query_string = """ select * from Emp_df a LEFT JOIN Dept_df b ON a.Empid = b.Empid;"""
sqldf(Query_string, globals())
```

#### Output:

Empid	Name	Designation	Date_of_joining	Empid	Deptartment	Total_Experience
1011	John	Manager	01-Jan-2000	1011.0	Management	18.0
1012	Rahul	Research Engineer	23-sep-2006	NaN	None	NaN
1013	Rick	Research Engineer	11-Jan-2012	1013.0	Research	10.0
1014	Morty	VP	21-Jan-1991	NaN	None	NaN
1015	Tim	Delivery Manager	12-Jan-1990	1015.0	Delivery	22.0



## Department of Computer Science and Engineering (Data Science)

### RIGHT JOIN:

A right merge or right join keeps every row on the right DataFrame. Where the missing values find the “on” variable in the left columns, you add the empty/ NaN values to the result. All the records from the left table and selected matching records based on the joining condition from one or more tables.

Pandas:

```
pd.merge(Emp_df, Dept_df, left_on='Empid', right_on='Empid', how='right')
```

### Output:

Empid	Name	Designation	Date_of_joining	Deptartment	Total_Experience
1011	John	Manager	01-Jan-2000	Management	18
1013	Rick	Research Engineer	11-Jan-2012	Research	10
1015	Tim	Delivery Manager	12-Jan-1990	Delivery	22
1017	NaN	NaN	NaN	Research	10
1019	NaN	NaN	NaN	Management	28

Pandasql:

```
Query_string = """ select
a.Empid,Name,Designation,Date_of_joining,Deptartment,Total_Experience from
Dept_df a LEFT JOIN Emp_df b ON a.Empid = b.Empid;""" sqldf(Query_string,
globals())
```

### Output:

Empid	Name	Designation	Date_of_joining	Deptartment	Total_Experience
1011	John	Manager	01-Jan-2000	Management	18
1017	None	None	None	Research	10
1013	Rick	Research Engineer	11-Jan-2012	Research	10
1019	None	None	None	Management	28
1015	Tim	Delivery Manager	12-Jan-1990	Delivery	22

### OUTER JOIN:

An outer merge/full outer join returns all the rows from left and right DataFrames only where it matches or returns NaNs or empty values. It returns all the records from both the tables based on the joining condition regardless of whether they match or not. Not matching ones will be nulls.

Pandas:

```
pd.merge(Emp_df, Dept_df, left_on='Empid', right_on='Empid', how='outer')
```



### Department of Computer Science and Engineering (Data Science)

#### Output:

Empid	Name	Designation	Date_of_joining	Deptartment	Total_Experience
1011	John	Manager	01-Jan-2000	Management	18.0
1012	Rahul	Research Engineer	23-sep-2006	NaN	NaN
1013	Rick	Research Engineer	11-Jan-2012	Research	10.0
1014	Morty	VP	21-Jan-1991	NaN	NaN
1015	Tim	Delivery Manager	12-Jan-1990	Delivery	22.0
1017	NaN	NaN	NaN	Research	10.0
1019	NaN	NaN	NaN	Management	28.0

Pandasql:

```
Query_string = """ select * from Emp_df a left OUTER JOIN Dept_df b ON
a.Empid = b.Empid;"""
sqldf(Query_string, globals())
```

#### Output:

Empid	Name	Designation	Date_of_joining	Empid	Deptartment	Total_Experience
1011	John	Manager	01-Jan-2000	1011.0	Management	18.0
1012	Rahul	Research Engineer	23-sep-2006	NaN	None	NaN
1013	Rick	Research Engineer	11-Jan-2012	1013.0	Research	10.0
1014	Morty	VP	21-Jan-1991	NaN	None	NaN
1015	Tim	Delivery Manager	12-Jan-1990	1015.0	Delivery	22.0

#### Summary of the DataFrame:

Descriptive statistics is a method used to depict data in a meaningful way to represent either the entire population or a sample. It has two sections: the measure of **central tendency**, which is used to measure the summary of a sample and population (e.g., mean, median, and mode), and the **measure of variability**, which is used to measure the spread or dispersion in the dataset (e.g., range, interquartile, variance, and standard deviation).

In Pandas, it provides a built-in function called **describe()**. The following example illustrates the detailed result in Pandas. In Pandas, the describe() function returns a number of descriptive statistics values (e.g., mean, standard deviation, min, median, max, first quartile, third quartile).

```
import pandas as pd
dfsumm =
pd.read_csv(r'https://archive.ics.uci.edu/ml/machine-
learningdatabases/abalone/abalone.data', names =['Sex',
'Length','Diameter', 'Height','Whole weight',
'Shucked weight','Viscera weight', 'Shell weight', 'Rings']) dfsumm.head(5)
```

#### Output:



### Department of Computer Science and Engineering (Data Science)

Sex	Length	Diameter	Height	Whole weight	Shucked weight	Viscera weight	Shell weight	Rings
M	0.455	0.365	0.095	0.5140	0.2245	0.1010	0.150	15
M	0.350	0.265	0.090	0.2255	0.0995	0.0485	0.070	7
F	0.530	0.420	0.135	0.6770	0.2565	0.1415	0.210	9
M	0.440	0.365	0.125	0.5160	0.2155	0.1140	0.155	10
I	0.330	0.255	0.080	0.2050	0.0895	0.0395	0.055	7

```
dfsumm.describe()
```

### Output:

	Length	Diameter	Height	Whole weight	Shucked weight	Viscera weight	Shell weight	Rings
count	4177.000000	4177.000000	4177.000000	4177.000000	4177.000000	4177.000000	4177.000000	4177.000000
mean	0.523992	0.407881	0.139516	0.828742	0.359367	0.180594	0.238831	9.933684
std	0.120093	0.099240	0.041827	0.490389	0.221963	0.109614	0.139203	3.224169
min	0.075000	0.055000	0.000000	0.002000	0.001000	0.000500	0.001500	1.000000
25%	0.450000	0.350000	0.115000	0.441500	0.186000	0.093500	0.130000	8.000000
50%	0.545000	0.425000	0.140000	0.799500	0.336000	0.171000	0.234000	9.000000
75%	0.615000	0.480000	0.165000	1.153000	0.502000	0.253000	0.329000	11.000000
max	0.815000	0.650000	1.130000	2.825500	1.488000	0.760000	1.005000	29.000000

### Resampling:

Resampling is a method to convert the time-based observed data into a different interval. In other words, it is used to change the time frequency into another time-frequency format. For instance, say we want to change monthly data into a year-wise format or upsample week data into hours. We would perform either upsample or downsample operations. For that, every data object must have a DateTime-like index(Datetimeindex, PeriodIndex, TimedeltaIndex).

```
import pandas as pd df = pd.read_csv(r'daily-total-female-births-CA.csv',index_col =0, parse_dates=['date']) df.head(5)
```

### Output:

births	
	date
1959-01-01	35
1959-01-02	32
1959-01-03	30
1959-01-04	31
1959-01-05	44

### Resampling by Month:

```
df.births.resample('M').mean()
```

### Output:



### Department of Computer Science and Engineering (Data Science)

```
date
1959-01-31    39.129032
1959-02-28    41.000000
1959-03-31    39.290323
1959-04-30    39.833333
1959-05-31    38.967742
1959-06-30    40.400000
1959-07-31    41.935484
1959-08-31    43.580645
1959-09-30    48.200000
1959-10-31    44.129032
1959-11-30    45.000000
1959-12-31    42.387097
Freq: M, Name: births, dtype: float64
```

#### Resampling by Quarter:

```
df.births.resample('Q').mean()
```

#### Output:

```
date
1959-03-31    39.766667
1959-06-30    39.725275
1959-09-30    44.532609
1959-12-31    43.826087
Freq: Q-DEC, Name: births, dtype: float64
```

#### Resampling by Year:

```
df.births.resample('Y').mean()
```

#### Output:

```
date
1959-12-31    41.980822
Freq: A-DEC, Name: births, dtype: float64
```

#### Resampling by Week:

```
df.births.resample('W').mean()
```

#### Output:



### Department of Computer Science and Engineering (Data Science)

date	
1959-01-04	32.000000
1959-01-11	37.714286
1959-01-18	44.285714
1959-01-25	41.142857
1959-02-01	35.142857
1959-02-08	40.428571
1959-02-15	42.857143
1959-02-22	42.428571
1959-03-01	40.000000
1959-03-08	39.428571
1959-03-15	36.571429
1959-03-22	40.857143
1959-03-29	39.142857
1959-04-05	41.142857
1959-04-12	37.857143
1959-04-19	37.285714
1959-04-26	40.142857

#### Resampling on a Semimonthly basis:

```
df.births.resample('SM').mean()
```

#### Output:

date	
1958-12-31	37.642857
1959-01-15	41.375000
1959-01-31	38.533333
1959-02-15	43.384615
1959-02-28	38.000000
1959-03-15	39.812500
1959-03-31	38.333333
1959-04-15	40.666667
1959-04-30	39.133333
1959-05-15	39.625000
1959-05-31	40.800000
1959-06-15	38.600000
1959-06-30	43.733333
1959-07-15	41.375000
1959-07-31	43.733333
1959-08-15	43.250000

#### Windowing Function:

A windowing function is used to calculate the number of operations, such as rolling count, rolling sum, rolling mean, rolling median, rolling variance, rolling standard deviation, rolling min, rolling max, rolling correlation, rolling covariance, rolling skewness, rolling kurtosis, rolling quantile, etc. In addition, we can perform some other operations such as expanding window and exponential weighted moving window.

```
import pandas as pd dfExcelwin =
pd.read_excel(r'istambul_stock_exchange.xlsx', sheet_name = 'Data'
               ,index_col =0,
parse_dates=[ 'date']) dfExcelwin.head(5)
```

#### Output:



### Department of Computer Science and Engineering (Data Science)

	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EM
date									
2009-01-05	0.035754	0.038376	-0.004679	0.002193	0.003894	0.000000	0.031190	0.012698	0.028524
2009-01-06	0.025426	0.031813	0.007787	0.008455	0.012866	0.004162	0.018920	0.011341	0.008773
2009-01-07	-0.028862	-0.026353	-0.030469	-0.017833	-0.028735	0.017293	-0.035899	-0.017073	-0.020015
2009-01-08	-0.062208	-0.084716	0.003391	-0.011726	-0.000466	-0.040061	0.028283	-0.005561	-0.019424
2009-01-09	0.009860	0.009658	-0.021533	-0.019873	-0.012710	-0.004474	-0.009764	-0.010989	-0.007802

#### Rolling Window:

The rolling window feature supports the following methods: count, sum, mean, median, var, std, min, max, corr, cov, skew, kurt, quantile, sum, and aggregate.

```
dfExcelwin.rolling(window=4).mean().head(10)
```

#### Output:

	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EM
date									
2009-01-05	NaN								
2009-01-06	NaN								
2009-01-07	NaN								
2009-01-08	-0.007473	-0.010220	-0.005993	-0.004728	-0.003110	-0.004651	0.010624	0.000351	-0.000536
2009-01-09	-0.013946	-0.017400	-0.010206	-0.010244	-0.007261	-0.005770	0.000385	-0.005570	-0.009617
2009-01-12	-0.027600	-0.035943	-0.017858	-0.015739	-0.011734	-0.019070	-0.017807	-0.011518	-0.017468
2009-01-13	-0.016523	-0.029423	-0.009802	-0.015700	-0.006086	-0.023393	-0.007940	-0.010305	-0.013671
2009-01-14	-0.011263	-0.017132	-0.019158	-0.024614	-0.018705	-0.012650	-0.025086	-0.020220	-0.010984
2009-01-15	-0.013563	-0.023863	-0.013443	-0.024533	-0.019112	-0.024143	-0.015066	-0.020491	-0.014891

#### Expanding Window:

Expanding window supports the following methods: count, sum, mean, median, var, std, min, max, corr, cov, skew, kurt, quantile, sum, aggregate, and quantile.

```
dfExcelwin.expanding(min_periods=4).mean().head(10)
```

#### Output:

	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EM
date									
2009-01-05	NaN								
2009-01-06	NaN								
2009-01-07	NaN								
2009-01-08	-0.007473	-0.010220	-0.005993	-0.004728	-0.003110	-0.004651	0.010624	0.000351	-0.000536
2009-01-09	-0.004006	-0.006244	-0.009101	-0.007757	-0.005030	-0.004616	0.006546	-0.001917	-0.001989
2009-01-12	-0.008204	-0.012264	-0.011388	-0.008718	-0.005029	-0.012020	-0.003520	-0.003672	-0.005429
2009-01-13	-0.004825	-0.010551	-0.009510	-0.009998	-0.005188	-0.010303	-0.002507	-0.004894	-0.005343
2009-01-14	-0.009368	-0.013676	-0.012575	-0.014671	-0.010908	-0.008651	-0.007231	-0.009934	-0.005760
2009-01-15	-0.008254	-0.014075	-0.011030	-0.015213	-0.011289	-0.013295	-0.003059	-0.010172	-0.007723

#### Exponentially Weighted Moving Window:



### Department of Computer Science and Engineering (Data Science)

Expanding window supports the following methods: mean, std, var, corr, and cov.

```
dfExcelwin.ewm(com=0.5).mean().head(10)
```

#### Output:

	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EM
date									
2009-01-05	0.035754	0.038376	-0.004679	0.002193	0.003894	0.000000	0.031190	0.012698	0.028524
2009-01-06	0.028008	0.033454	0.004670	0.006890	0.010623	0.003122	0.021987	0.011680	0.013711
2009-01-07	-0.011363	-0.007951	-0.019657	-0.010226	-0.016625	0.012933	-0.018088	-0.008226	-0.009638
2009-01-08	-0.045684	-0.059767	-0.004099	-0.011239	-0.005718	-0.022838	0.013213	-0.006427	-0.016243
2009-01-09	-0.008502	-0.013292	-0.015770	-0.017019	-0.010398	-0.010545	-0.002168	-0.009481	-0.010593
2009-01-12	-0.022313	-0.032698	-0.020478	-0.014687	-0.006812	-0.036242	-0.036670	-0.011464	-0.018628
2009-01-13	0.002871	-0.011071	-0.005648	-0.016679	-0.006365	-0.012070	-0.009830	-0.011968	-0.009423
2009-01-14	-0.026493	-0.027394	-0.024574	-0.037152	-0.036090	-0.002080	-0.030147	-0.034140	-0.008925

#### Shifting:

Shifting is also known as “lag” and moves a value back or forward in time. Pandas has a df.shift() function to shift an index by the desired number of periods with an optional time frequency.

```
import pandas as pd dfshift =  
pd.read_excel(r'istambul_stock_exchange.xlsx', sheet_name = 'Data'  
,index_col =0,  
parse_dates=[ 'date' ]) dfshift.head(5)
```

#### Output:

	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EM
date									
2009-01-05	0.035754	0.038376	-0.004679	0.002193	0.003894	0.000000	0.031190	0.012698	0.028524
2009-01-06	0.025426	0.031813	0.007787	0.008455	0.012866	0.004162	0.018920	0.011341	0.008773
2009-01-07	-0.028862	-0.026353	-0.030469	-0.017833	-0.028735	0.017293	-0.035899	-0.017073	-0.020015
2009-01-08	-0.062208	-0.084716	0.003391	-0.011726	-0.000466	-0.040061	0.028283	-0.005561	-0.019424
2009-01-09	0.009860	0.009658	-0.021533	-0.019873	-0.012710	-0.004474	-0.009764	-0.010989	-0.007802

```
dfshift.shift(periods=3).head(7)
```

#### Output:

	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EM
date									
2009-01-05	NaN								
2009-01-06	NaN								
2009-01-07	NaN								
2009-01-08	0.035754	0.038376	-0.004679	0.002193	0.003894	0.000000	0.031190	0.012698	0.028524
2009-01-09	0.025426	0.031813	0.007787	0.008455	0.012866	0.004162	0.018920	0.011341	0.008773
2009-01-12	-0.028862	-0.026353	-0.030469	-0.017833	-0.028735	0.017293	-0.035899	-0.017073	-0.020015
2009-01-13	-0.062208	-0.084716	0.003391	-0.011726	-0.000466	-0.040061	0.028283	-0.005561	-0.019424

```
dfshift.shift(periods=-1).head(7)
```



### Department of Computer Science and Engineering (Data Science)

#### Output:

	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EM
date									
2009-01-05	0.025426	0.031813	0.007787	0.008455	0.012866	0.004162	0.018920	0.011341	0.008773
2009-01-06	-0.028862	-0.026353	-0.030469	-0.017833	-0.028735	0.017293	-0.035899	-0.017073	-0.020015
2009-01-07	-0.062208	-0.084716	0.003391	-0.011726	-0.000466	-0.040061	0.028283	-0.005561	-0.019424
2009-01-08	0.009860	0.009658	-0.021533	-0.019873	-0.012710	-0.004474	-0.009764	-0.010989	-0.007802
2009-01-09	-0.029191	-0.042361	-0.022823	-0.013526	-0.005026	-0.049039	-0.053849	-0.012451	-0.022630
2009-01-12	0.015445	-0.000272	0.001757	-0.017674	-0.006141	0.000000	0.003572	-0.012220	-0.004827
2009-01-13	-0.041168	-0.035552	-0.034032	-0.047383	-0.050945	0.002912	-0.040302	-0.045220	-0.008677

```
dfshift.shift(periods=3, axis =1).head(7)
```

#### Output:

	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EM
date									
2009-01-05	NaN	NaN	NaN	0.035754	0.038376	-0.004679	0.002193	0.003894	0.000000
2009-01-06	NaN	NaN	NaN	0.025426	0.031813	0.007787	0.008455	0.012866	0.004162
2009-01-07	NaN	NaN	NaN	-0.028862	-0.026353	-0.030469	-0.017833	-0.028735	0.017293
2009-01-08	NaN	NaN	NaN	-0.062208	-0.084716	0.003391	-0.011726	-0.000466	-0.040061
2009-01-09	NaN	NaN	NaN	0.009860	0.009658	-0.021533	-0.019873	-0.012710	-0.004474
2009-01-12	NaN	NaN	NaN	-0.029191	-0.042361	-0.022823	-0.013526	-0.005026	-0.049039
2009-01-13	NaN	NaN	NaN	0.015445	-0.000272	0.001757	-0.017674	-0.006141	0.000000

```
dfshift.shift(periods=3,fill_value=0).head(7)
```

#### Output:

	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EM
date									
2009-01-05	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
2009-01-06	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
2009-01-07	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
2009-01-08	0.035754	0.038376	-0.004679	0.002193	0.003894	0.000000	0.031190	0.012698	0.028524
2009-01-09	0.025426	0.031813	0.007787	0.008455	0.012866	0.004162	0.018920	0.011341	0.008773
2009-01-12	-0.028862	-0.026353	-0.030469	-0.017833	-0.028735	0.017293	-0.035899	-0.017073	-0.020015
2009-01-13	-0.062208	-0.084716	0.003391	-0.011726	-0.000466	-0.040061	0.028283	-0.005561	-0.019424

#### Handling Missing Data:

A missing value or missing data occurs when no data value is found at a respective number of places in either a DataFrame or dataset. A missing value presents many problems in the dataset. It reduces the statistical power of data, and the lost data can increase the bias in the dataset. So, it is essential to handle this missing value to maintain the characteristics of the data. Pandas has a number of methods to handle the missing values such as bfill(), ffill(), interpolate(), and fillna().



### Department of Computer Science and Engineering (Data Science)

```
import pandas as pd dfmiss = pd.read_csv(r'daily-total-female-births-CA-
with_nulls.csv',index_col
=0,
parse_dates=[ 'date' ]) dfmiss
```

#### Output:

```
births
      date
1959-01-01    35.0
1959-01-02    32.0
1959-01-03    30.0
1959-01-04    31.0
1959-01-05    44.0
1959-01-06    29.0
1959-01-07    45.0
1959-01-08     NaN
1959-01-09    38.0
1959-01-10    27.0
# Snippet to check for nulls dfmiss.isnull().sum()
```

#### Output:

```
births    16
dtype: int64
```

#### BFILL:

The backward method fills the missing values in the dataset and uses the next valid observation to fill the gap.

```
dfmiss.bfill()
```

#### Output:



### Department of Computer Science and Engineering (Data Science)

births

date

1959-01-01	35.0
1959-01-02	32.0
1959-01-03	30.0
1959-01-04	31.0
1959-01-05	44.0
1959-01-06	29.0
1959-01-07	45.0
1959-01-08	38.0
1959-01-09	38.0
1959-01-10	27.0
1959-01-11	38.0

#### FFILL:

The forward fill method propagates the last valid observation forward to the next valid one.

`dfmiss.fillna()`

#### Output:

births

date

1959-01-01	35.0
1959-01-02	32.0
1959-01-03	30.0
1959-01-04	31.0
1959-01-05	44.0
1959-01-06	29.0
1959-01-07	45.0
1959-01-08	45.0
1959-01-09	38.0
1959-01-10	27.0
1959-01-11	38.0

#### FILLNA:

This replaces NA/NaN values using a constant value.

`dfmiss.fillna(10)`

#### Output:



## Department of Computer Science and Engineering (Data Science)

births

date

1959-01-01	35.0
1959-01-02	32.0
1959-01-03	30.0
1959-01-04	31.0
1959-01-05	44.0
1959-01-06	29.0
1959-01-07	45.0
1959-01-08	10.0
1959-01-09	38.0
1959-01-10	27.0
1959-01-11	38.0

### INTERPOLATE:

This method interpolates values linearly in a forward direction.

```
#interpolate
dfmiss.interpolate(method='linear', limit_direction='forward')
```

### Output:

births

date

1959-01-01	35.0
1959-01-02	32.0
1959-01-03	30.0
1959-01-04	31.0
1959-01-05	44.0
1959-01-06	29.0
1959-01-07	45.0
1959-01-08	41.5
1959-01-09	38.0
1959-01-10	27.0
1959-01-11	38.0

### Lab Assignments to complete:

Perform the following tasks using the datasets mentioned. Download the datasets from the link given:

**Link:** <https://drive.google.com/drive/folders/107ecI29wJUZdnKI--PqvRfkUe8MaQHse?usp=sharing>

### Dataset 1: Smart City Index Headers

1. Implement data wrangling with the help of various pre-processing strategies.

```
In [1]: import pandas as pd
```

```
In [2]: df = pd.read_csv("D:\Lab setup\Smart_City_index_headers.csv")
```

```
In [3]: df.head(5)
```

Out[3]:

	<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Mobility</b>	<b>Smart_Environment</b>	<b>Smart_Government</b>	<b>Smart_Econc</b>
<b>0</b>	1	Oslo	Norway	6480	6512	7516	4
<b>1</b>	2	Bergen	Norway	7097	6876	7350	4
<b>2</b>	3	Amsterdam	Netherlands	7540	5558	8528	8
<b>3</b>	4	Copenhagen	Denmark	7490	7920	8726	5
<b>4</b>	5	Stockholm	Sweden	6122	7692	8354	4

```
In [4]: !pip install pandasql
```

Requirement already satisfied: pandasql in c:\users\djsce.student\appdata\local\anaconda3\lib\site-packages (0.7.3)  
Requirement already satisfied: numpy in c:\users\djsce.student\appdata\local\anaconda3\lib\site-packages (from pandasql) (1.24.3)  
Requirement already satisfied: pandas in c:\users\djsce.student\appdata\local\anaconda3\lib\site-packages (from pandasql) (1.5.3)  
Requirement already satisfied: sqlalchemy in c:\users\djsce.student\appdata\local\anaconda3\lib\site-packages (from pandasql) (1.4.39)  
Requirement already satisfied: python-dateutil>=2.8.1 in c:\users\djsce.student\appdata\roaming\python\python311\site-packages (from pandas->pandasql) (2.8.2)  
Requirement already satisfied: pytz>=2020.1 in c:\users\djsce.student\appdata\local\anaconda3\lib\site-packages (from pandas->pandasql) (2022.7)  
Requirement already satisfied: greenlet!=0.4.17 in c:\users\djsce.student\appdata\local\anaconda3\lib\site-packages (from sqlalchemy->pandasql) (2.0.1)  
Requirement already satisfied: six>=1.5 in c:\users\djsce.student\appdata\roaming\python\python311\site-packages (from python-dateutil>=2.8.1->pandas->pandasql) (1.16.0)

```
In [5]: pip install pandasql
```

Requirement already satisfied: pandasql in c:\users\djsce.student\appdata\local\anaconda3\lib\site-packages (0.7.3)  
Requirement already satisfied: numpy in c:\users\djsce.student\appdata\local\anaconda3\lib\site-packages (from pandasql) (1.24.3)  
Requirement already satisfied: pandas in c:\users\djsce.student\appdata\local\anaconda3\lib\site-packages (from pandasql) (1.5.3)  
Requirement already satisfied: sqlalchemy in c:\users\djsce.student\appdata\local\anaconda3\lib\site-packages (from pandasql) (1.4.39)  
Requirement already satisfied: python-dateutil>=2.8.1 in c:\users\djsce.student\appdata\roaming\python\python311\site-packages (from pandas->pandasql) (2.8.2)  
Requirement already satisfied: pytz>=2020.1 in c:\users\djsce.student\appdata\local\anaconda3\lib\site-packages (from pandas->pandasql) (2022.7)  
Requirement already satisfied: greenlet!=0.4.17 in c:\users\djsce.student\appdata\local\anaconda3\lib\site-packages (from sqlalchemy->pandasql) (2.0.1)  
Requirement already satisfied: six>=1.5 in c:\users\djsce.student\appdata\roaming\python\python311\site-packages (from python-dateutil>=2.8.1->pandas->pandasql) (1.16.0)  
Note: you may need to restart the kernel to use updated packages.

```
In [6]: from pandasql import sqldf
dfpsql= pd.read_csv("D:\Lab setup\Smart_City_index_headers.csv")
Query_string= """ select * from dfpsql limit 5 """
sqldf(Query_string, globals())
```

Out[6]:

	<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Mobility</b>	<b>Smart_Environment</b>	<b>Smart_Government</b>	<b>Smart_Econc</b>
<b>0</b>	1	Oslo	Norway	6480	6512	7516	4
<b>1</b>	2	Bergen	Norway	7097	6876	7350	4
<b>2</b>	3	Amsterdam	Netherlands	7540	5558	8528	8
<b>3</b>	4	Copenhagen	Denmark	7490	7920	8726	5
<b>4</b>	5	Stockholm	Sweden	6122	7692	8354	4

```
In [7]: df[(df['Smart_People'] >=2825) & (df['Smart_People'] <=9695)]
```

Out[7]:

	<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Mobility</b>	<b>Smart_Environment</b>	<b>Smart_Government</b>	<b>Smart_Econc</b>
<b>0</b>	1	Oslo	Norway	6480	6512	7516	
<b>1</b>	2	Bergen	Norway	7097	6876	7350	
<b>2</b>	3	Amsterdam	Netherlands	7540	5558	8528	
<b>3</b>	4	Copenhagen	Denmark	7490	7920	8726	
<b>4</b>	5	Stockholm	Sweden	6122	7692	8354	
...	...	...	...	...	...	...	...
<b>97</b>	98	Riga	Latvia	4152	4584	4616	
<b>98</b>	99	Beijing	China	7610	2998	2806	
<b>99</b>	100	St Petersburg	Russia	4588	2908	3622	
<b>100</b>	101	Calgary	Canada	6675	4052	5946	
<b>101</b>	102	Edmonton	Canada	5801	4499	6396	

102 rows × 11 columns

```
In [12]: Query_string= """ select * from dfpsql where Smart_People>=2825 and Smart_People<=9695
sqldf(Query_string, globals())
```

Out[12]:

	<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Mobility</b>	<b>Smart_Environment</b>	<b>Smart_Government</b>	<b>Smart_Ec</b>
<b>0</b>	1	Oslo	Norway	6480	6512	7516	
<b>1</b>	2	Bergen	Norway	7097	6876	7350	
<b>2</b>	3	Amsterdam	Netherlands	7540	5558	8528	
<b>3</b>	4	Copenhagen	Denmark	7490	7920	8726	
<b>4</b>	5	Stockholm	Sweden	6122	7692	8354	
...	...	...	...	...	...	...	...
<b>97</b>	98	Riga	Latvia	4152	4584	4616	
<b>98</b>	99	Beijing	China	7610	2998	2806	
<b>99</b>	100	St Petersburg	Russia	4588	2908	3622	
<b>100</b>	101	Calgary	Canada	6675	4052	5946	
<b>101</b>	102	Edmonton	Canada	5801	4499	6396	

102 rows × 11 columns

In [13]:

```
Query_string= """ select * from dfpsql where Smart_People>=2825 and Smart_People<=9695
sqldf(Query_string, globals())
```

Out[13]:

	<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Mobility</b>	<b>Smart_Environment</b>	<b>Smart_Government</b>	<b>Smart_Ec</b>
<b>0</b>	1	Oslo	Norway	6480	6512	7516	
<b>1</b>	2	Bergen	Norway	7097	6876	7350	
<b>2</b>	3	Amsterdam	Netherlands	7540	5558	8528	
<b>3</b>	4	Copenhagen	Denmark	7490	7920	8726	
<b>4</b>	5	Stockholm	Sweden	6122	7692	8354	
...	...	...	...	...	...	...	...
<b>97</b>	98	Riga	Latvia	4152	4584	4616	
<b>98</b>	99	Beijing	China	7610	2998	2806	
<b>99</b>	100	St Petersburg	Russia	4588	2908	3622	
<b>100</b>	101	Calgary	Canada	6675	4052	5946	
<b>101</b>	102	Edmonton	Canada	5801	4499	6396	

102 rows × 11 columns

In [14]:

```
df['Id'].unique()
```

```
Out[14]: array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13,
       14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26,
       27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39,
       40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52,
       53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65,
       66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78,
       79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91,
       92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102], dtype=int64)
```

```
In [15]: Query_string= """ select distinct Id from dfpsql;"""
sqldf(Query_string, globals())
```

```
Out[15]:   Id
0    1
1    2
2    3
3    4
4    5
...
97   98
98   99
99   100
100  101
101  102
```

102 rows × 1 columns

```
In [16]: df[df.Smart_People.isin([3745,8050,5390])]
```

	<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Mobility</b>	<b>Smart_Environment</b>	<b>Smart_Government</b>	<b>Smart_Econom</b>
1	2	Bergen	Norway	7097	6876	7350	490
97	98	Riga	Latvia	4152	4584	4616	738
99	100	St Petersburg	Russia	4588	2908	3622	451

```
In [17]: Query_string= """ select * from dfpsql where Smart_People in(3745,8050,5390);"""
sqldf(Query_string, globals())
```

	<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Mobility</b>	<b>Smart_Environment</b>	<b>Smart_Government</b>	<b>Smart_Economy</b>
<b>0</b>	2	Bergen	Norway	7097	6876	7350	4905
<b>1</b>	98	Riga	Latvia	4152	4584	4616	7380
<b>2</b>	100	St Petersburg	Russia	4588	2908	3622	4515

In [18]: `df[df.Smart_People.isin([3745,8050,5390])]`

	<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Mobility</b>	<b>Smart_Environment</b>	<b>Smart_Government</b>	<b>Smart_Econom</b>
<b>1</b>	2	Bergen	Norway	7097	6876	7350	490
<b>97</b>	98	Riga	Latvia	4152	4584	4616	738
<b>99</b>	100	St Petersburg	Russia	4588	2908	3622	451

In [19]: `Query_string= """ select * from dfpsql where Smart_People not in(3745,8050,5390);"""
sqldf(Query_string, globals())`

	<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Mobility</b>	<b>Smart_Environment</b>	<b>Smart_Government</b>	<b>Smart_Eco</b>
<b>0</b>	1	Oslo	Norway	6480	6512	7516	
<b>1</b>	3	Amsterdam	Netherlands	7540	5558	8528	
<b>2</b>	4	Copenhagen	Denmark	7490	7920	8726	
<b>3</b>	5	Stockholm	Sweden	6122	7692	8354	
<b>4</b>	6	Montreal	Canada	7490	4848	6624	
...	...	...	...	...	...	...	...
<b>94</b>	96	Naples	Italy	3175	6802	4008	
<b>95</b>	97	Moscow	Russia	5015	2772	5078	
<b>96</b>	99	Beijing	China	7610	2998	2806	
<b>97</b>	101	Calgary	Canada	6675	4052	5946	
<b>98</b>	102	Edmonton	Canada	5801	4499	6396	

99 rows × 11 columns

In [20]: `df.sort_values(by= ['Smart_People', 'City'], ascending= True)`

Out[20]:

	<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Mobility</b>	<b>Smart_Environment</b>	<b>Smart_Government</b>	<b>Smart_Economy</b>
88	89	Ancona	Italy	4208	7496	5826	7910
76	77	Bayreuth	Germany	3932	7696	5630	6391
93	94	Kuala Lumpur	Malaysia	4347	2076	5890	2720
58	59	Osaka	Japan	7233	4758	5732	2105
59	60	Bordeaux	France	5550	7586	6832	7720
...	...	...	...	...	...	...	...
5	6	Montreal	Canada	7490	4848	6624	6180
6	7	Vienna	Austria	5683	7608	6232	5419
0	1	Oslo	Norway	6480	6512	7516	4561
62	63	Daejeon	South Korea	7018	4192	4304	2250
8	9	Singapore	Singapore	5790	4344	5560	5531

102 rows × 11 columns

In [21]: `Query_string= """ select * from dfpsql order by Smart_People,City;"""
sqldf(Query_string, globals())`

Out[21]:

	<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Mobility</b>	<b>Smart_Environment</b>	<b>Smart_Government</b>	<b>Smart_Econom</b>
0	89	Ancona	Italy	4208	7496	5826	7910
1	77	Bayreuth	Germany	3932	7696	5630	6391
2	94	Kuala Lumpur	Malaysia	4347	2076	5890	2720
3	59	Osaka	Japan	7233	4758	5732	2105
4	60	Bordeaux	France	5550	7586	6832	7720
...	...	...	...	...	...	...	...
97	6	Montreal	Canada	7490	4848	6624	6180
98	7	Vienna	Austria	5683	7608	6232	5419
99	1	Oslo	Norway	6480	6512	7516	4561
100	63	Daejeon	South Korea	7018	4192	4304	2250
101	9	Singapore	Singapore	5790	4344	5560	5531

102 rows × 11 columns

In [22]: `df.sort_values(by= ['Smart_People','City'], ascending= False)`

Out[22]:

	<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Mobility</b>	<b>Smart_Environment</b>	<b>Smart_Government</b>	<b>Smart_Economy</b>
8	9	Singapore	Singapore	5790	4344	5560	553!
62	63	Daejeon	South Korea	7018	4192	4304	2250
0	1	Oslo	Norway	6480	6512	7516	456!
6	7	Vienna	Austria	5683	7608	6232	541!
5	6	Montreal	Canada	7490	4848	6624	6180
...	...	...	...	...	...	...	...
59	60	Bordeaux	France	5550	7586	6832	7720
58	59	Osaka	Japan	7233	4758	5732	210!
93	94	Kuala Lumpur	Malaysia	4347	2076	5890	2720
76	77	Bayreuth	Germany	3932	7696	5630	639!
88	89	Ancona	Italy	4208	7496	5826	7910

102 rows × 11 columns

In [23]:

```
Query_string= """ select * from dfpsql order by Smart_People Desc,City Desc;"""
sqldf(Query_string, globals())
```

Out[23]:

	<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Mobility</b>	<b>Smart_Environment</b>	<b>Smart_Government</b>	<b>Smart_Econom</b>
0	9	Singapore	Singapore	5790	4344	5560	553!
1	63	Daejeon	South Korea	7018	4192	4304	2250
2	1	Oslo	Norway	6480	6512	7516	456!
3	7	Vienna	Austria	5683	7608	6232	541!
4	6	Montreal	Canada	7490	4848	6624	6180
...	...	...	...	...	...	...	...
97	60	Bordeaux	France	5550	7586	6832	7720
98	59	Osaka	Japan	7233	4758	5732	210!
99	94	Kuala Lumpur	Malaysia	4347	2076	5890	2720
100	77	Bayreuth	Germany	3932	7696	5630	639!
101	89	Ancona	Italy	4208	7496	5826	7910

102 rows × 11 columns

In [24]:

```
df.agg({'Smart_Living': ['count','min', 'max', 'mean']})
```

Out[24]:

Smart_Living	
<b>count</b>	102.000000
<b>min</b>	1980.000000
<b>max</b>	10000.000000
<b>mean</b>	6377.039216

```
In [25]: Query_string= """ select count(Smart_Living) as count, min(Smart_Living) as min, max(S  
sqlsdf(Query_string, globals())
```

Out[25]:

	count	min	max	mean
<b>0</b>	102	1980	10000	6377.039216

```
In [26]: df.groupby('Id')[['City']].sum()
```

Out[26]:

```
Id  
1                 Oslo  
2                 Bergen  
3            Amsterdam  
4            Copenhagen  
5             Stockholm  
...  
98                Riga  
99                Beijing  
100           St Petersburg  
101              Calgary  
102            Edmonton  
Name: City, Length: 102, dtype: object
```

```
In [27]: Query_string= """ select Id , sum(City) as Sum_City from df group by Id;"""  
sqlsdf(Query_string, globals())
```

Out[27]:

	<b>Id</b>	<b>Sum_City</b>
<b>0</b>	1	0.0
<b>1</b>	2	0.0
<b>2</b>	3	0.0
<b>3</b>	4	0.0
<b>4</b>	5	0.0
...	...	...
<b>97</b>	98	0.0
<b>98</b>	99	0.0
<b>99</b>	100	0.0
<b>100</b>	101	0.0
<b>101</b>	102	0.0

102 rows × 2 columns

In [28]: `df.groupby('Country').agg({'Smart_People':[ 'mean', 'min', 'max' ]})`

Out[28]:

Country	Smart_People		
	mean	min	max
Australia	5495.250000	4678	7305
Austria	8580.000000	8580	8580
Belgium	6680.000000	6680	6680
Canada	7112.500000	6200	8465
China	4663.666667	4385	5183
Czech Republic	5315.000000	5315	5315
Denmark	6377.000000	5780	6955
Estonia	4955.000000	4955	4955
Finland	5980.000000	4385	7523
France	5400.857143	3425	7075
Germany	5308.600000	3108	6853
Hungary	7260.000000	7260	7260
Iceland	7285.000000	7285	7285
Ireland	5425.000000	5425	5425
Israel	4113.000000	4113	4113
Italy	4660.181818	2825	6185
Japan	3780.500000	3403	4158
Latvia	3745.000000	3745	3745
Luxembourg	4723.000000	4723	4723
Malaysia	3313.000000	3313	3313
Netherlands	7098.000000	7098	7098
New Zealand	7855.000000	7855	7855
Norway	7955.250000	7558	8618
Portugal	3560.000000	3560	3560
Russia	4766.500000	4143	5390
Singapore	9695.000000	9695	9695
Slovakia	5485.000000	5485	5485
Slovenia	5503.000000	5503	5503
South Korea	8247.500000	7860	8635
Spain	5725.000000	3635	7815
Sweden	6132.666667	4590	7065

Country	Smart_People		
	mean	min	max
<b>Switzerland</b>	6425.000000	5650	7200
<b>Taiwan</b>	7858.000000	7858	7858
<b>United Arab Emirates</b>	6529.000000	5933	7125
<b>United Kingdom</b>	5757.000000	4548	6363
<b>United States</b>	6594.142857	5495	7498

In [29]:

```
Query_string= """ select Country , avg(Smart_People) as mean, min(Smart_People) as min, max(Smart_People) as max
sql(df(Query_string, globals()))
```

Out[29]:

	Country	mean	min	max
<b>0</b>	Australia	5495.250000	4678	7305
	Austria	8580.000000	8580	
<b>2</b>	Belgium	6680.000000	6680	
<b>3</b>	Canada	7112.500000	6200	
<b>4</b>	China	4663.666667	4385	
<b>5</b>	Czech Republic	5315.000000	5315	
<b>6</b>	Denmark	6377.000000	5780	
<b>7</b>	Estonia	4955.000000	4955	
<b>8</b>	Finland	5980.000000	4385	
<b>9</b>	France	5400.857143	3425	
			3108	
	Germany	5308.600000		
	Hungary	7260.000000	7260	
	Iceland	7285.000000	7285	
	Ireland	5425.000000	5425	
	Israel			
		4113.000000	4113	
	Italy	4660.181818	2825	
	Japan	3780.500000	3403	
	Latvia			
		3745.000000	3745	
	Luxembourg	4723.000000	4723	
	Malaysia	3313.000000	3313	
	Netherlands			
		7098.000000	7098	
	New Zealand			
		7855.000000	7855	
	Norway	7955.250000	7558	
	Portugal	3560.000000	3560	
	Russia			
		4143		
		4766.500000		
	Singapore	9695.000000	9695	
	Slovakia	5485.000000	5485	

Slovenia	5503.000000	5503
South Korea	8247.500000	7860
Spain	5725.000000	3635
Sweden	6132.666667	4590
Switzerland		
Sweden	6425.000000	5650
	7858.000000	7858

	Country	mean	min	max
33	United Arab Emirates	6529.000000	5933	7125
34	United Kingdom	5757.000000	4548	6363
35	United States	6594.142857	5495	7498

```
In [30]: import pandas as pd
df = pd.read_csv("D:\Lab setup\Smart_City_index_headers.csv")
```

```
In [54]: SmartCity_df= pd.DataFrame(df, columns= ['Id','City', 'Country', 'Smart_Index','SmartCity_Index_relative_Edmonton']
SmartCity_df
```

Out[54]:

	<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Index</b>	<b>SmartCity_Index_relative_Edmonton</b>
0	1	Oslo	Norway	NaN	666
1	2	Bergen	Norway	NaN	823
2	3	Amsterdam	Netherlands	NaN	839
3	4	Copenhagen	Denmark	NaN	698
4	5	Stockholm	Sweden	NaN	340
...	...	...	...	...	...
97	98	Riga	Latvia	NaN	-1760
98	99	Beijing	China	NaN	-2023
99	100	St Petersburg	Russia	NaN	-2281
100	101	Calgary	Canada	NaN	206
101	102	Edmonton	Canada	NaN	0

102 rows × 5 columns

```
In [55]: import pandas as pd
df = pd.read_csv("D:\Lab setup\Smart_City_index_headers.csv")
Smart_Living_df= pd.DataFrame(df, columns= ['Id','City', 'Country', 'Smart_Living'])
Smart_Living_df
```

Out[55]:

	<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Living</b>
0	1	Oslo	Norway	9090
1	2	Bergen	Norway	9090
2	3	Amsterdam	Netherlands	7280
3	4	Copenhagen	Denmark	7200
4	5	Stockholm	Sweden	7730
...	...	...	...	...
97	98	Riga	Latvia	4330
98	99	Beijing	China	1980
99	100	St Petersburg	Russia	4100
100	101	Calgary	Canada	8657
101	102	Edmonton	Canada	8141

102 rows × 4 columns

In [56]: `pd.merge(SmartCity_df, Smart_Living_df, left_on='Id', right_on='Id', how='inner')`

Out[56]:

	<b>Id</b>	<b>City_x</b>	<b>Country_x</b>	<b>Smart_Index</b>	<b>SmartCity_Index_relative_Edmonton</b>	<b>City_y</b>	<b>C</b>
0	1	Oslo	Norway	NaN	666	Oslo	
1	2	Bergen	Norway	NaN	823	Bergen	
2	3	Amsterdam	Netherlands	NaN	839	Amsterdam	Ne
3	4	Copenhagen	Denmark	NaN	698	Copenhagen	
4	5	Stockholm	Sweden	NaN	340	Stockholm	
...	...	...	...	...	...	...	...
97	98	Riga	Latvia	NaN	-1760	Riga	
98	99	Beijing	China	NaN	-2023	Beijing	
99	100	St Petersburg	Russia	NaN	-2281	St Petersburg	
100	101	Calgary	Canada	NaN	206	Calgary	
101	102	Edmonton	Canada	NaN	0	Edmonton	

102 rows × 8 columns

In [57]: `Query_string = """ select * from SmartCity_df a INNER JOIN Smart_Living_df b ON a.Id = sqldf(Query_string, globals())`

<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Index</b>	<b>SmartCity_Index_relative_Edmonton</b>	<b>Id</b>	<b>City</b>
-----------	-------------	----------------	--------------------	--	-----------	-------------

Out[57]:

0	1	Oslo	Norway	None	666	1	Oslo
1	2	Bergen	Norway	None	823	2	Bergen
2	3	Amsterdam	Netherlands	None	839	3	Amsterdam
3	4	Copenhagen	Denmark	None	698	4	Copenhagen
4	5	Stockholm	Sweden	None	340	5	Stockholm
...	...	...	...	...	...	...	...
97	98	Riga	Latvia	None	-1760	98	Riga
98	99	Beijing	China	None	-2023	99	Beijing
99	100	St Petersburg	Russia	None	-2281	100	St Petersburg
100	101	Calgary	Canada	None	206	101	Calgary
101	102	Edmonton	Canada	None	0	102	Edmonton

102 rows × 9 columns

In [58]: `pd.merge(SmartCity_df, Smart_Living_df, left_on='Id', right_on='Id', how='left')`

<b>Id</b>	<b>City_x</b>	<b>Country_x</b>	<b>Smart_Index</b>	<b>SmartCity_Index_relative_Edmonton</b>	<b>City_y</b>	<b>Country_y</b>
0	1	Oslo	Norway	NaN	666	Oslo
1	2	Bergen	Norway	NaN	823	Bergen
2	3	Amsterdam	Netherlands	NaN	839	Amsterdam
3	4	Copenhagen	Denmark	NaN	698	Copenhagen
4	5	Stockholm	Sweden	NaN	340	Stockholm
...	...	...	...	...	...	...
97	98	Riga	Latvia	NaN	-1760	Riga
98	99	Beijing	China	NaN	-2023	Beijing
99	100	St Petersburg	Russia	NaN	-2281	St Petersburg
100	101	Calgary	Canada	NaN	206	Calgary
101	102	Edmonton	Canada	NaN	0	Edmonton

```
Query_string= """ select * from SmartCity_df
sqldf(Query_string, globals())
```

<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Index</b>	<b>SmartCity_Index_relative_Edmonton</b>	<b>Id</b>	<b>Cit</b>
-----------	-------------	----------------	--------------------	--	-----------	------------

In [59]:

```
LEFT JOIN Smart_Living_df b ON a.Id = b.Id
```

102 rows × 8 columns

```
Query_string= """ select * from SmartCity_df a  
sqldf(Query_string, globals())
```

<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Index</b>	<b>SmartCity_Index_relative_Edmonton</b>	<b>Id</b>	<b>City</b>
-----------	-------------	----------------	--------------------	--	-----------	-------------

Out[59]:

0	1	Oslo	Norway	None	666	1	Oslo
1	2	Bergen	Norway	None	823	2	Bergen
2	3	Amsterdam	Netherlands	None	839	3	Amsterdam
3	4	Copenhagen	Denmark	None	698	4	Copenhagen
4	5	Stockholm	Sweden	None	340	5	Stockholm
...	...	...	...	...	...	...	...
97	98	Riga	Latvia	None	-1760	98	Riga
98	99	Beijing	China	None	-2023	99	Beijing
99	100	St Petersburg	Russia	None	-2281	100	St Petersburg
100	101	Calgary	Canada	None	206	101	Calgary
101	102	Edmonton	Canada	None	0	102	Edmonton

102 rows × 9 columns

In [60]: `pd.merge(SmartCity_df, Smart_Living_df, left_on='Id', right_on='Id', how='right')`

<b>Id</b>	<b>City_x</b>	<b>Country_x</b>	<b>Smart_Index</b>	<b>SmartCity_Index_relative_Edmonton</b>	<b>City_y</b>	<b>Country_y</b>
0	1	Oslo	Norway	NaN	666	Oslo
1	2	Bergen	Norway	NaN	823	Bergen
2	3	Amsterdam	Netherlands	NaN	839	Amsterdam
3	4	Copenhagen	Denmark	NaN	698	Copenhagen
4	5	Stockholm	Sweden	NaN	340	Stockholm
...	...	...	...	...	...	...
97	98	Riga	Latvia	NaN	-1760	Riga
98	99	Beijing	China	NaN	-2023	Beijing
99	100	St Petersburg	Russia	NaN	-2281	St Petersburg
100	101	Calgary	Canada	NaN	206	Calgary
101	102	Edmonton	Canada	NaN	0	Edmonton

```
Query_string= """ select * from SmartCity_df
sqldf(Query_string, globals())
```

<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Index</b>	<b>SmartCity_Index_relative_Edmonton</b>	<b>Id</b>	<b>Cit</b>
-----------	-------------	----------------	--------------------	--	-----------	------------

In [61]:

RIGHT JOIN Smart\_Living\_df b ON a.Id =

102 rows × 8 columns

```
Query_string= """ select * from SmartCity_df a
sql(df(Query_string, globals()))
```

<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Index</b>	<b>SmartCity_Index_relative_Edmonton</b>	<b>Id</b>	<b>City</b>
-----------	-------------	----------------	--------------------	--	-----------	-------------

Out[61]:

0	1	Oslo	Norway	None	666	1	Oslo
1	2	Bergen	Norway	None	823	2	Bergen
2	3	Amsterdam	Netherlands	None	839	3	Amsterdam
3	4	Copenhagen	Denmark	None	698	4	Copenhagen
4	5	Stockholm	Sweden	None	340	5	Stockholm
...	...	...	...	...	...	...	...
97	98	Riga	Latvia	None	-1760	98	Riga
98	99	Beijing	China	None	-2023	99	Beijing
99	100	St Petersburg	Russia	None	-2281	100	St Petersburg
100	101	Calgary	Canada	None	206	101	Calgary
101	102	Edmonton	Canada	None	0	102	Edmonton

102 rows × 9 columns

In [62]: `pd.merge(SmartCity_df, Smart_Living_df, left_on='Id', right_on='Id', how='outer')`

<b>Id</b>	<b>City_x</b>	<b>Country_x</b>	<b>Smart_Index</b>	<b>SmartCity_Index_relative_Edmonton</b>	<b>City_y</b>	<b>Country_y</b>	
0	1	Oslo	Norway	NaN	666	Oslo	
1	2	Bergen	Norway	NaN	823	Bergen	
2	3	Amsterdam	Netherlands	NaN	839	Amsterdam	Ne
3	4	Copenhagen	Denmark	NaN	698	Copenhagen	
4	5	Stockholm	Sweden	NaN	340	Stockholm	
...	...	...	...	...	...	...	
97	98	Riga	Latvia	NaN	-1760	Riga	
98	99	Beijing	China	NaN	-2023	Beijing	
99	100	St Petersburg	Russia	NaN	-2281	St Petersburg	
100	101	Calgary	Canada	NaN	206	Calgary	
101	102	Edmonton	Canada	NaN	0	Edmonton	

```
Query_string= """ select * from SmartCity_df
sqldf(Query_string, globals())
```

<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Index</b>	<b>SmartCity_Index_relative_Edmonton</b>	<b>Id</b>	<b>Cit</b>
-----------	-------------	----------------	--------------------	--	-----------	------------

In [63]:

```
left OUTER JOIN Smart_Living_df b ON a.
```

102 rows × 8 columns

```
Query_string= """ select * from SmartCity_df a  
sqldf(Query_string, globals())
```

<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Index</b>	<b>SmartCity_Index_relative_Edmonton</b>	<b>Id</b>	<b>Cit</b>
-----------	-------------	----------------	--------------------	--	-----------	------------

Out[63]:

0	1	Oslo	Norway	None	666	1	Oslo
1	2	Bergen	Norway	None	823	2	Bergen
2	3	Amsterdam	Netherlands	None	839	3	Amsterdam
3	4	Copenhagen	Denmark	None	698	4	Copenhagen
4	5	Stockholm	Sweden	None	340	5	Stockholm
...	...	...	...	...	...	...	...
97	98	Riga	Latvia	None	-1760	98	Riga
98	99	Beijing	China	None	-2023	99	Beijing
99	100	St Petersburg	Russia	None	-2281	100	S Petersburg
100	101	Calgary	Canada	None	206	101	Calgary
101	102	Edmonton	Canada	None	0	102	Edmonton

102 rows × 9 columns

In [64]: `import pandas as pd  
df = pd.read_csv("D:\\Lab setup\\Smart_City_index_headers.csv")`

In [65]: `df.head(5)`

<b>Id</b>	<b>City</b>	<b>Country</b>	<b>Smart_Mobility</b>	<b>Smart_Environment</b>	<b>Smart_Government</b>	<b>Smart_Econ</b>
0	1	Oslo	6480	6512	7516	4
1	2	Bergen	7097	6876	7350	4
2	3	Amsterdam	7540	5558	8528	8
3	4	Copenhagen	7490	7920	8726	5
4	5	Stockholm	6122	7692	8354	4

In [66]: `df.describe()`

date

	<b>Id</b>	<b>Smart_Mobility</b>	<b>Smart_Environment</b>	<b>Smart_Government</b>	<b>Smart_Economy</b>	<b>Smart_Pe</b>
<b>count</b>	102.000000	102.000000	102.000000	102.000000	102.000000	102.000000
<b>mean</b>	51.500000	5759.401961	5943.500000	5893.803922	6131.803922	5874.04
<b>std</b>	29.588849	1214.030137	1724.032171	1153.375297	1801.555148	1449.09
<b>min</b>	1.000000	3175.000000	1850.000000	2806.000000	1490.000000	2825.00
<b>25%</b>	26.250000	4809.750000	4530.500000	5143.000000	5007.500000	4724.75
<b>50%</b>	51.500000	5651.500000	6495.000000	5911.000000	6432.500000	5747.50
<b>75%</b>	76.750000	6763.750000	7310.000000	6581.500000	7492.500000	7061.25
<b>max</b>	102.000000	8110.000000	8844.000000	8726.000000	9225.000000	9695.00

◀ ▶

```
In [68]: import pandas as pd
df= pd.read_csv("D:\Lab setup\daily-total-female-births-CA.csv",index_col=0,parse_date
df.head(5)
```

Out[68]:

**births**

<b>date</b>	
<b>1959-01-01</b>	35
<b>1959-01-02</b>	32
<b>1959-01-03</b>	30
<b>1959-01-04</b>	31
<b>1959-01-05</b>	44

```
In [69]: df.births.resample('M').mean()
```

```
Out[69]: date
1959-01-31    39.129032
1959-02-28    41.000000
1959-03-31    39.290323
1959-04-30    39.833333
1959-05-31    38.967742
1959-06-30    40.400000
1959-07-31    41.935484
1959-08-31    43.580645
1959-09-30    48.200000
1959-10-31    44.129032
1959-11-30    45.000000
1959-12-31    42.387097
Freq: M, Name: births, dtype: float64
```

```
In [70]: df.births.resample('Q').mean()
```

```
Out[70]: 1959-03-31    39.766667  
1959-06-30    39.725275  
1959-09-30    44.532609  
1959-12-31    43.826087  
Freq: Q-DEC, Name: births, dtype: float64
```

```
In [71]: df.births.resample('Y').mean()
```

```
Out[71]: date  
1959-12-31    41.980822  
Freq: A-DEC, Name: births, dtype: float64
```

```
In [72]: df.births.resample('W').mean()
```

date

```
Out[72]:
```

1959-01-04	32.000000
1959-01-11	37.714286
1959-01-18	44.285714
1959-01-25	41.142857
1959-02-01	35.142857
1959-02-08	40.428571
1959-02-15	42.857143
1959-02-22	42.428571
1959-03-01	40.000000
1959-03-08	39.428571
1959-03-15	36.571429
1959-03-22	40.857143
1959-03-29	39.142857
1959-04-05	41.142857
1959-04-12	37.857143
1959-04-19	37.285714
1959-04-26	40.142857
1959-05-03	42.285714
1959-05-10	39.000000
1959-05-17	37.428571
1959-05-24	42.142857
1959-05-31	39.000000
1959-06-07	42.142857
1959-06-14	39.857143
1959-06-21	37.714286
1959-06-28	39.142857
1959-07-05	44.000000
1959-07-12	44.285714
1959-07-19	41.142857
1959-07-26	40.285714
1959-08-02	43.000000
1959-08-09	44.571429
1959-08-16	43.285714
1959-08-23	40.857143
1959-08-30	45.285714
1959-09-06	46.714286
1959-09-13	45.571429
1959-09-20	45.857143
1959-09-27	52.714286
1959-10-04	51.428571
1959-10-11	45.000000
1959-10-18	46.714286
1959-10-25	40.857143
1959-11-01	40.285714
1959-11-08	48.000000
1959-11-15	39.571429
1959-11-22	42.285714
1959-11-29	48.714286
1959-12-06	40.428571
1959-12-13	39.857143
1959-12-20	44.857143
1959-12-27	40.714286
1960-01-03	51.250000

Freq: W-SUN, Name: births, dtype: float64

In [73]: df.births.resample('SM').mean()

	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EN
--	-----	------	----	-----	------	--------	---------	----	----

**date**

```
Out[73]:
```

1958-12-31	37.642857
1959-01-15	41.375000
1959-01-31	38.533333
1959-02-15	43.384615
1959-02-28	38.000000
1959-03-15	39.812500
1959-03-31	38.333333
1959-04-15	40.666667
1959-04-30	39.133333
1959-05-15	39.625000
1959-05-31	40.800000
1959-06-15	38.600000
1959-06-30	43.733333
1959-07-15	41.375000
1959-07-31	43.733333
1959-08-15	43.250000
1959-08-31	45.333333
1959-09-15	50.266667
1959-09-30	47.800000
1959-10-15	41.312500
1959-10-31	44.466667
1959-11-15	45.133333
1959-11-30	40.933333
1959-12-15	43.875000
1959-12-31	50.000000

Freq: SM-15, Name: births, dtype: float64

```
In [77]:
```

```
import pandas as pd
dfExcelwin= pd.read_excel("D:\Lab setup\istambul_stock_exchange.xlsx", sheet_name= 'Da
dfExcelwin.head(5)
```

C:\Users\djsce.student\AppData\Local\anaconda3\Lib\site-packages\openpyxl\worksheet\\_read\_only.py:79: UserWarning: Unknown extension is not supported and will be removed  
for idx, row in parser.parse():

```
Out[77]:
```

	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EN
<b>date</b>									
<b>2009-01-05</b>	0.035754	0.038376	-0.004679	0.002193	0.003894	0.000000	0.031190	0.012698	0.028524
<b>2009-01-06</b>	0.025426	0.031813	0.007787	0.008455	0.012866	0.004162	0.018920	0.011341	0.008773
<b>2009-01-07</b>	-0.028862	-0.026353	-0.030469	-0.017833	-0.028735	0.017293	-0.035899	-0.017073	-0.020015
<b>2009-01-08</b>	-0.062208	-0.084716	0.003391	-0.011726	-0.000466	-0.040061	0.028283	-0.005561	-0.019424
<b>2009-01-09</b>	0.009860	0.009658	-0.021533	-0.019873	-0.012710	-0.004474	-0.009764	-0.010989	-0.007802

```
In [78]: dfExcelwin.rolling(window=4).mean().head(10)
```

date

Out[78]:

<b>2009-01-05</b>	NaN								
<b>2009-01-06</b>	NaN								
<b>2009-01-07</b>	NaN								
<b>2009-01-08</b>	-0.007473	-0.010220	-0.005993	-0.004728	-0.003110	-0.004651	0.010624	0.000351	-0.000536
<b>2009-01-09</b>	-0.013946	-0.017400	-0.010206	-0.010244	-0.007261	-0.005770	0.000385	-0.005570	-0.009617
<b>2009-01-12</b>	-0.027600	-0.035943	-0.017858	-0.015739	-0.011734	-0.019070	-0.017807	-0.011518	-0.017468
<b>2009-01-13</b>	-0.016523	-0.029423	-0.009802	-0.015700	-0.006086	-0.023393	-0.007940	-0.010305	-0.013671
<b>2009-01-14</b>	-0.011263	-0.017132	-0.019158	-0.024614	-0.018705	-0.012650	-0.025086	-0.020220	-0.010984
<b>2009-01-15</b>	-0.013563	-0.023863	-0.013443	-0.024533	-0.019112	-0.024143	-0.015066	-0.020491	-0.014891
<b>2009-01-16</b>	-0.000756	-0.005204	-0.005854	-0.019454	-0.016283	-0.005521	-0.000387	-0.015237	-0.006504

In [79]: dfExcelwin.expanding(min\_periods=4).mean().head(10)

	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EN
--	-----	------	----	-----	------	--------	---------	----	----

Out[79]:

date	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EN
2009-01-05	NaN								
2009-01-06	NaN								
2009-01-07	NaN								
2009-01-08	-0.007473	-0.010220	-0.005993	-0.004728	-0.003110	-0.004651	0.010624	0.000351	-0.000536
2009-01-09	-0.004006	-0.006244	-0.009101	-0.007757	-0.005030	-0.004616	0.006546	-0.001917	-0.001989
2009-01-12	-0.008204	-0.012264	-0.011388	-0.008718	-0.005029	-0.012020	-0.003520	-0.003672	-0.005429
2009-01-13	-0.004825	-0.010551	-0.009510	-0.009998	-0.005188	-0.010303	-0.002507	-0.004894	-0.005343
2009-01-14	-0.009368	-0.013676	-0.012575	-0.014671	-0.010908	-0.008651	-0.007231	-0.009934	-0.005760
2009-01-15	-0.008254	-0.014075	-0.011030	-0.015213	-0.011289	-0.013295	-0.003059	-0.010172	-0.007723
2009-01-16	-0.005224	-0.009440	-0.009174	-0.013013	-0.009531	-0.009420	-0.002267	-0.008298	-0.005859

◀ ▶

In [81]: dfExcelwin.ewm(com=0.5).mean().head(10)

Out[81]:

2009-01-05	0.035754	0.038376	-0.004679	0.002193	0.003894	0.000000	0.031190	0.012698	0.028524
2009-01-06	0.028008	0.033454	0.004670	0.006890	0.010623	0.003122	0.021987	0.011680	0.013711
2009-01-07	-0.011363	-0.007951	-0.019657	-0.010226	-0.016625	0.012933	-0.018088	-0.008226	-0.009638
2009-01-08	-0.045684	-0.059767	-0.004099	-0.011239	-0.005718	-0.022838	0.013213	-0.006427	-0.016243
2009-01-09	-0.008502	-0.013292	-0.015770	-0.017019	-0.010398	-0.010545	-0.002168	-0.009481	-0.010593
2009-01-12	-0.022313	-0.032698	-0.020478	-0.014687	-0.006812	-0.036242	-0.036670	-0.011464	-0.018628
2009-01-13	0.002871	-0.011071	-0.005648	-0.016679	-0.006365	-0.012070	-0.009830	-0.011968	-0.009423
2009-01-14	-0.026493	-0.027394	-0.024574	-0.037152	-0.036090	-0.002080	-0.030147	-0.034140	-0.008925
2009-01-15	-0.008389	-0.020643	-0.007305	-0.025417	-0.021586	-0.034327	0.010162	-0.019426	-0.018595
2009-01-16	0.011896	0.014638	0.002587	-0.003945	-0.003002	0.005527	0.006632	-0.000768	0.001080

In [83]:

```
import pandas as pd
dfshift= pd.read_excel("D:\Lab setup\istambul_stock_exchange.xlsx", sheet_name= 'Data'
dfshift.head(5)
```

C:\Users\djsce.student\AppData\Local\anaconda3\Lib\site-packages\openpyxl\worksheet\\_read\_only.py:79: UserWarning: Unknown extension is not supported and will be removed  
for idx, row in parser.parse():

Out[83]:

	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EN
date									
2009-01-05	0.035754	0.038376	-0.004679	0.002193	0.003894	0.000000	0.031190	0.012698	0.028524
2009-01-06	0.025426	0.031813	0.007787	0.008455	0.012866	0.004162	0.018920	0.011341	0.008773
2009-01-07	-0.028862	-0.026353	-0.030469	-0.017833	-0.028735	0.017293	-0.035899	-0.017073	-0.020015
2009-01-08	-0.062208	-0.084716	0.003391	-0.011726	-0.000466	-0.040061	0.028283	-0.005561	-0.019424
2009-01-09	0.009860	0.009658	-0.021533	-0.019873	-0.012710	-0.004474	-0.009764	-0.010989	-0.007802

ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EN
-----	------	----	-----	------	--------	---------	----	----

**date**

```
In [84]: dfshift.shift(periods=3).head(7)
```

Out[84]:

ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EN
-----	------	----	-----	------	--------	---------	----	----

**date**

<b>2009-01-05</b>	NaN								
<b>2009-01-06</b>	NaN								
<b>2009-01-07</b>	NaN								
<b>2009-01-08</b>	0.035754	0.038376	-0.004679	0.002193	0.003894	0.000000	0.031190	0.012698	0.028524
<b>2009-01-09</b>	0.025426	0.031813	0.007787	0.008455	0.012866	0.004162	0.018920	0.011341	0.008773
<b>2009-01-12</b>	-0.028862	-0.026353	-0.030469	-0.017833	-0.028735	0.017293	-0.035899	-0.017073	-0.020015
<b>2009-01-13</b>	-0.062208	-0.084716	0.003391	-0.011726	-0.000466	-0.040061	0.028283	-0.005561	-0.019424

```
In [85]: dfshift.shift(periods=-1).head(7)
```

Out[85]:

ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EN
-----	------	----	-----	------	--------	---------	----	----

**date**

<b>2009-01-05</b>	0.025426	0.031813	0.007787	0.008455	0.012866	0.004162	0.018920	0.011341	0.008773
<b>2009-01-06</b>	-0.028862	-0.026353	-0.030469	-0.017833	-0.028735	0.017293	-0.035899	-0.017073	-0.020015
<b>2009-01-07</b>	-0.062208	-0.084716	0.003391	-0.011726	-0.000466	-0.040061	0.028283	-0.005561	-0.019424
<b>2009-01-08</b>	0.009860	0.009658	-0.021533	-0.019873	-0.012710	-0.004474	-0.009764	-0.010989	-0.007802
<b>2009-01-09</b>	-0.029191	-0.042361	-0.022823	-0.013526	-0.005026	-0.049039	-0.053849	-0.012451	-0.022630
<b>2009-01-12</b>	0.015445	-0.000272	0.001757	-0.017674	-0.006141	0.000000	0.003572	-0.012220	-0.004827
<b>2009-01-13</b>	-0.041168	-0.035552	-0.034032	-0.047383	-0.050945	0.002912	-0.040302	-0.045220	-0.008677

```
In [86]: dfshift.shift(periods=3, axis=1).head(7)
```

Out[86]:

	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EM
date									
2009-01-05	NaN	NaN	NaN	0.035754	0.038376	-0.004679	0.002193	0.003894	0.000000
2009-01-06	NaN	NaN	NaN	0.025426	0.031813	0.007787	0.008455	0.012866	0.004162
2009-01-07	NaN	NaN	NaN	-0.028862	-0.026353	-0.030469	-0.017833	-0.028735	0.017293
2009-01-08	NaN	NaN	NaN	-0.062208	-0.084716	0.003391	-0.011726	-0.000466	-0.040061
2009-01-09	NaN	NaN	NaN	0.009860	0.009658	-0.021533	-0.019873	-0.012710	-0.004474
2009-01-12	NaN	NaN	NaN	-0.029191	-0.042361	-0.022823	-0.013526	-0.005026	-0.049039
2009-01-13	NaN	NaN	NaN	0.015445	-0.000272	0.001757	-0.017674	-0.006141	0.000000

In [87]: `dfshift.shift(periods=3, fill_value=0).head(7)`

Out[87]:

	ISE	ISED	SP	DAX	FTSE	NIKKEI	BOVESPA	EU	EM
date									
2009-01-05	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
2009-01-06	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
2009-01-07	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
2009-01-08	0.035754	0.038376	-0.004679	0.002193	0.003894	0.000000	0.031190	0.012698	0.028524
2009-01-09	0.025426	0.031813	0.007787	0.008455	0.012866	0.004162	0.018920	0.011341	0.008773
2009-01-12	-0.028862	-0.026353	-0.030469	-0.017833	-0.028735	0.017293	-0.035899	-0.017073	-0.020015
2009-01-13	-0.062208	-0.084716	0.003391	-0.011726	-0.000466	-0.040061	0.028283	-0.005561	-0.019424

In [88]:

```
import pandas as pd
dfmiss= pd.read_csv("D:\Lab setup\daily-total-female-births-CA-with_nulls.csv",index_col=0)
dfmiss
```

Out[88]:

**births**

<b>date</b>	
<b>1959-01-01</b>	35.0
<b>1959-01-02</b>	32.0
<b>1959-01-03</b>	30.0
<b>1959-01-04</b>	31.0
<b>1959-01-05</b>	44.0
...	...
<b>1959-12-27</b>	37.0
<b>1959-12-28</b>	52.0
<b>1959-12-29</b>	48.0
<b>1959-12-30</b>	55.0
<b>1959-12-31</b>	50.0

365 rows × 1 columns

In [89]: `dfmiss.isnull().sum()`Out[89]: `births 16`  
`dtype: int64`In [90]: `dfmiss.bfill()`

Out[90]:

**births**

<b>date</b>	
<b>1959-01-01</b>	35.0
<b>1959-01-02</b>	32.0
<b>1959-01-03</b>	30.0
<b>1959-01-04</b>	31.0
<b>1959-01-05</b>	44.0
...	...
<b>1959-12-27</b>	37.0
<b>1959-12-28</b>	52.0
<b>1959-12-29</b>	48.0
<b>1959-12-30</b>	55.0
<b>1959-12-31</b>	50.0

365 rows × 1 columns

```
In [91]: dfmiss.ffill()
```

```
Out[91]:          births
```

	date
1	1959-01-01 35.0
2	1959-01-02 32.0
3	1959-01-03 30.0
4	1959-01-04 31.0
5	1959-01-05 44.0
	...
365	1959-12-27 37.0
366	1959-12-28 52.0
367	1959-12-29 48.0
368	1959-12-30 55.0
369	1959-12-31 50.0

365 rows × 1 columns

```
In [92]: dfmiss.fillna(10)
```

```
Out[92]:          births
```

	date
1	1959-01-01 35.0
2	1959-01-02 32.0
3	1959-01-03 30.0
4	1959-01-04 31.0
5	1959-01-05 44.0
	...
365	1959-12-27 37.0
366	1959-12-28 52.0
367	1959-12-29 48.0
368	1959-12-30 55.0
369	1959-12-31 50.0

365 rows × 1 columns

```
localhost:8888 In [93]: dfmiss.interpolate(method='linear',limit_direction='forward')
```

Out[93]:

**births**

<b>date</b>	
<b>1959-01-01</b>	35.0
<b>1959-01-02</b>	32.0
<b>1959-01-03</b>	30.0
<b>1959-01-04</b>	31.0
<b>1959-01-05</b>	44.0
...	...
<b>1959-12-27</b>	37.0
<b>1959-12-28</b>	52.0
<b>1959-12-29</b>	48.0
<b>1959-12-30</b>	55.0
<b>1959-12-31</b>	50.0

365 rows × 1 columns

In [ ]: