

# Efficient Data Cleaning for Customer Sales Analysis using Pandas

## Importing the library

In [288].

```
import pandas as pd
```

## Reading the dataset

In [289].

```
df = pd.read_excel("C:\\Users\\HP\\Desktop\\EDA\\msh\\Customer_Call_List.xlsx")
```

In [290].

CustomerID	First_Name	Last_Name	Phone_Number	Address	Paying Customer	Do_Not_Contact	Not_Useful_Column
1001	Frodo	Baggins	123-545-5421	123 Shire Lane, Shire	Yes	No	True
1002	Abed	Nadri	1236439775	93 West Main Street	No	Yes	False
1003	Walter	White	7066950392	298 Drugs Driveway	N	NaN	True
1004	Dwight	Schrute	123-543-2345	980 Paper Avenue, Pennsylvania, 18503	Yes	Y	True
1005	Jon	Snow	8766783469	123 Dragons Road	Y	No	True
1006	Ron	Swanson	304-762-2467	768 City Parkway	Yes	Yes	True
1007	Jeff	Winger	NaN	1209 South Street	No	No	False
1008	Sherlock	Holmes	8766783469	98 Clue Drive	N	No	False
1009	Gandalf	NaN	N/A	123 Middle Earth	Yes	NaN	False
1010	Peter	Parker	123-545-5421	25th Main Street, New York	Yes	No	True
1011	Samwise	Gamgee	NaN	612 Shire Lane, Shire	Yes	No	True
1012	Harry	Potter	7066950392	2394 Hogwarts Avenue	Y	NaN	True
1013	Don	Draper	123-543-2345	2039 Main Street	Yes	N	False
1014	Leslie	Knope	8766783469	343 City Parkway	Yes	No	False
1015	Toby	Flenderson	304-762-2467	214 HR Avenue	N	No	False
1016	Ron	Weasley	123-545-5421	2395 Hogwarts Avenue	No	N	False
1017	Michael	Scott	1236439775	121 Paper Avenue, Pennsylvania	Yes	No	False
1018	Clark	Kent	7066950392	3498 Super Lane	Y	NaN	True
1019	Creed	Braton	N/A	N/A	N/A	Yes	True
1020	Anakin	Skywalker	8766783469	910 Tatooine Road, Tatooine	Yes	N	True
1020	Anakin	Skywalker	8766783469	910 Tatooine Road, Tatooine	Yes	N	True

df.drop\_duplicates(inplace=True)

## To drop the duplicates in the entire column

In [291].

```
df=df.drop_duplicates()
df
```

Out[291].

	CustomerID	First_Name	Last_Name	Phone_Number	Address	Paying Customer	Do_Not_Contact	Not_Useful_Column
0	1001	Frodo	Baggins	123-545-5421	123 Shire Lane, Shire	Yes	No	True
1	1002	Abed	Nadri	1236439775	93 West Main Street	No	Yes	False
2	1003	Walter	White	7066950392	298 Drugs Driveway	N	NaN	True
3	1004	Dwight	Schrute	123-543-2345	980 Paper Avenue, Pennsylvania, 18503	Yes	Y	True
4	1005	Jon	Snow	8766783469	123 Dragons Road	Y	No	True
5	1006	Ron	Swanson	304-762-2467	768 City Parkway	Yes	Yes	True
6	1007	Jeff	Winger	NaN	1209 South Street	No	No	False
7	1008	Sherlock	Holmes	8766783469	98 Clue Drive	N	No	False
8	1009	Gandalf	NaN	N/A	123 Middle Earth	Yes	NaN	False
9	1010	Peter	Parker	123-545-5421	25th Main Street, New York	Yes	No	True
10	1011	Samwise	Gamgee	NaN	612 Shire Lane, Shire	Yes	No	True
11	1012	Harry	Potter	7066950392	2394 Hogwarts Avenue	Y	NaN	True
12	1013	Don	Draper	123-543-2345	2039 Main Street	Yes	N	False
13	1014	Leslie	Knope	8766783469	343 City Parkway	Yes	No	False
14	1015	Toby	Flenderson	304-762-2467	214 HR Avenue	N	No	False
15	1016	Ron	Weasley	123-545-5421	2395 Hogwarts Avenue	No	N	False
16	1017	Michael	Scott	1236439775	121 Paper Avenue, Pennsylvania	Yes	No	False
17	1018	Clark	Kent	7066950392	3498 Super Lane	Y	NaN	True
18	1019	Creed	Braton	N/A	N/A	N/A	Yes	True
19	1020	Anakin	Skywalker	8766783469	910 Tatooine Road, Tatooine	Yes	N	True

## To remove the unwanted columns

In [292].

```
df=df.drop(columns = "Not_Useful_Column")
df
```

Out[292].

	CustomerID	First_Name	Last_Name	Phone_Number	Address	Paying Customer	Do_Not_Contact
0	1001	Frodo	Baggins	123-545-5421	123 Shire Lane, Shire	Yes	No
1	1002	Abed	Nadri	1236439775	93 West Main Street	No	Yes
2	1003	Walter	White	7066950392	298 Drugs Driveway	N	NaN
3	1004	Dwight	Schrute	123-543-2345	980 Paper Avenue, Pennsylvania, 18503	Yes	Y
4	1005	Jon	Snow	8766783469	123 Dragons Road	Y	No
5	1006	Ron	Swanson	304-762-2467	768 City Parkway	Yes	Yes
6	1007	Jeff	Winger	NaN	1209 South Street	No	No
7	1008	Sherlock	Holmes	8766783469	98 Clue Drive	N	No
8	1009	Gandalf	NaN	N/A	123 Middle Earth	Yes	NaN
9	1010	Peter	Parker	123-545-5421	25th Main Street, New York	Yes	No
10	1011	Samwise	Gamgee	NaN	612 Shire Lane, Shire	Yes	No
11	1012	Harry	Potter	7066950392	2394 Hogwarts Avenue	Y	NaN
12	1013	Don	Draper	123-543-2345	2039 Main Street	Yes	N
13	1014	Leslie	Knope	8766783469	343 City Parkway	Yes	No
14	1015	Toby	Flenderson	304-762-2467	214 HR Avenue	N	No
15	1016	Ron	Weasley	123-545-5421	2395 Hogwarts Avenue	No	N
16	1017	Michael	Scott	1236439775	121 Paper Avenue, Pennsylvania	Yes	No
17	1018	Clark	Kent	7066950392	3498 Super Lane	Y	NaN
18	1019	Creed	Braton	N/A	N/A	N/A	Yes
19	1020	Anakin	Skywalker	8766783469	910 Tatooine Road, Tatooine	Yes	N

## Using 'Strip' command to remove the carbages in a particular column

In [293].

```
#df["Last_Name"]=df["Last_Name"].str.strip() # Left strip is used here
#df["Last_Name"]=df["Last_Name"].str.strip(" ") # Left strip is used here
#df["Last_Name"]=df["Last_Name"].str.strip("/") # Left strip is used here
#df["Last_Name"]=df["Last_Name"].str.rstrip(" ") # right strip is used here
df["Last_Name"]=df["Last_Name"].str.strip("123_/_") # Removing all carbages in one commad using 'strip' command
df
```

Out[293].

	CustomerID	First_Name	Last_Name	Phone_Number	Address	Paying Customer	Do_Not_Contact
0	1001	Frodo	Baggins	123-545-5421	123 Shire Lane, Shire	Yes	No
1	1002	Abed	Nadri	1236439775	93 West Main Street	No	Yes
2	1003	Walter	White	7066950392	298 Drugs Driveway	N	NaN
3	1004	Dwight	Schrute	123-543-2345	980 Paper Avenue, Pennsylvania, 18503	Yes	Y
4	1005	Jon	Snow	8766783469	123 Dragons Road	Y	No
5	1006	Ron	Swanson	304-762-2467	768 City Parkway	Yes	Yes
6	1007	Jeff	Winger	NaN	1209 South Street	No	No
7	1008	Sherlock	Holmes	8766783469	98 Clue Drive	N	No
8	1009	Gandalf	NaN	N/A	123 Middle Earth	Yes	NaN
9	1010	Peter	Parker	123-545-5421	25th Main Street, New York	Yes	No
10	1011	Samwise	Gamgee	NaN	612 Shire Lane, Shire	Yes	No
11	1012	Harry	Potter	7066950392	2394 Hogwarts Avenue	Y	NaN
12	1013	Don	Draper	123-543-2345	2039 Main Street	Yes	N
13	1014	Leslie	Knope	8766783469	343 City Parkway	Yes	No
14	1015	Toby	Flenderson	304-762-2467	214 HR Avenue	N	No
15	1016	Ron	Weasley	123-545-5421	2395 Hogwarts Avenue	No	N
16	1017	Michael	Scott	1236439775	121 Paper Avenue, Pennsylvania	Yes	No
17	1018	Clark	Kent	7066950392	3498 Super Lane	Y	NaN
18	1019	Creed	Braton	N/A	N/A	N/A	Yes
19	1020	Anakin	Skywalker	8766783469	910 Tatooine Road, Tatooine	Yes	N

In [294].

1020

Anakin

Skywalker

8766783469

910 Tatooine Road, Tatooine

Yes

N

eaning the phone number

f["Phone\_Number"]=df["Phone\_Number"].str.replace('[^a-zA-Z0-9]','')

df["Phone\_Number"]=df["Phone\_Number"].apply(lambda x: x[0:3]+'-' + x[3:6]+'-' + x[6:10])

f["Phone\_Number"] = df["Phone\_Number"].apply(lambda x: str(x))

f["Phone\_Number"] = df["Phone\_Number"].apply(lambda x: x[0:3]+'-' + x[3:6]+'-' + x[6:10])

f["Phone\_Number"]=df["Phone\_Number"].str.replace('nan--','')

f["Phone\_Number"]=df["Phone\_Number"].str.replace('Na--','')

f

Users\HP\AppData\Local\Temp\ipykernel\_9892\4225834453.py:1: FutureWarning: The default value of regex will change from True to False in a future version.

df["Phone\_Number"]=df["Phone\_Number"].str.replace('[^a-zA-Z0-9]','')

CustomerID

First\_Name

Last\_Name

Phone\_Number

Address

Paying Customer

Do\_Not\_Contact

0

1001

Frodo

Baggins

123-545-5421

123 Shire Lane, Shire

Yes

No

1

1002

Abed

Nadir

123-643-9775

93 West Main Street

No

Yes

2

1003

Walter

White

706-695-0392

298 Drugs Driveway

N

NaN

3

1004

Dwight

Schrute

123-543-2345

980 Paper Avenue, Pennsylvania, 18503

Yes

Y

4

1005

Jon

Snow

876-678-3469

123 Dragons Road

Y

No

5

1006

Ron

Swanson

304-762-2467

768 City Parkway

Yes

Yes

6

1007

Jeff

Winger

NaN

1209 South Street

No

No

7

1008

Sherlock

Holmes

8766783469

98 Clue Drive

N

No

8

1009

Gandalf

NaN

N/A

123 Middle Earth

Yes

NaN

9

1010

Peter

Parker

123-545-5421

25th Main Street, New York

Yes

No

10

1011

Samwise

Gamgee

NaN

612 Shire Lane, Shire

Yes

No

11

1012

Harry

Potter

7066950392

2394 Hogwarts Avenue

Y

NaN

12

1013

Don

Draper

123-543-2345

2039 Main Street

Yes

N

13

1014

Leslie

Knope

8766783469

343 City Parkway

Yes

No

14

1015

Toby

Flenderson

304-762-2467

214 HR Avenue

N

No

15

1016

Ron

Weasley

123-545-5421

2395 Hogwarts Avenue

No

N

16

1017

Michael

Scott

1236439775

121 Paper Avenue, Pennsylvania

Yes

No

17

1018

Clark

Kent

7066950392

3498 Super Lane

Y

NaN

18

1019

Creed

Braton

N/A

N/A

N/A

Yes

19

1020

Anakin

Skywalker

8766783469

910 Tatooine Road, Tatooine

Yes

N

## Cleaning the phone number

In [295].

```
df["Phone_Number"]=df["Phone_Number"].str.replace('[^a-zA-Z0-9]', '')
#df["Phone_Number"]=df["Phone_Number"].apply(lambda x: x[0:3]+'-'+x[3:6]+'-'+x[6:10])
df["Phone_Number"] = df["Phone_Number"].apply(lambda x: str(x))
df["Phone_Number"] = df["Phone_Number"].apply(lambda x: x[0:3]+'-'+x[3:6]+'-'+x[6:10])
df["Phone_Number"]=df["Phone_Number"].str.replace('NaN--','')
df["Phone_Number"]=df["Phone_Number"].str.replace('NaN--','')
df
C:\Users\HP\AppData\Local\Temp\ipykernel_9892\4225834453.py:1: FutureWarning: The default value of regex will change from True to False in a future version.
df["Phone_Number"]=df["Phone_Number"].str.replace('[^a-zA-Z0-9]', '')
```

Out[295].

	CustomerID	First_Name	Last_Name	Phone_Number	Address	Paying Customer	Do_Not_Contact
0	1001	Frodo	Baggins	123-545-5421	123 Shire Lane, Shire	Yes	No
1	1002	Abed	Nadri	123-643-9775	93 West Main Street	No	Yes
2	1003	Walter	White	7066950392	298 Drugs Driveway	N	NaN
3	1004	Dwight	Schrute	123-543-2345	980 Paper Avenue, Pennsylvania, 18503	Yes	Y
4	1005	Jon	Snow	876-678-3469	123 Dragons Road	Y	No
5	1006	Ron	Swanson	304-762-2467	768 City Parkway	Yes	Yes
6	1007	Jeff	Winger	NaN	1209 South Street	No	No
7	1008	Sherlock	Holmes	876-678-3469	98 Clue Drive	N	No
8	1009	Gandalf	NaN	N/A	123 Middle Earth	Yes	NaN
9	1010	Peter	Parker	123-545-5421	25th Main Street, New York	Yes	No
10	1011	Samwise	Gamgee	NaN	612 Shire Lane, Shire	Yes	No
11	1012	Harry	Potter	7066950392	2394 Hogwarts Avenue	Y	NaN
12	1013	Don	Draper	123-543-2345	2039 Main Street	Yes	N
13	1014	Leslie	Knope	876-678-3469	343 City Parkway	Yes	No
14	1015	Toby	Flenderson	304-762-2467	214 HR Avenue	N	No
15	1016	Ron	Weasley	123-545-5421	2395 Hogwarts Avenue	No	N
16	1017	Michael	Scott	123-643-9775	121 Paper Avenue, Pennsylvania	Yes	No
17	1018	Clark	Kent	7066950392	3498 Super Lane	Y	NaN
18	1019	Creed	Braton	N/A	N/A	N/A	Yes
19	1020	Anakin	Skywalker	876-678-3469	910 Tatooine Road, Tatooine	Yes	N

## To expand the address column

In [296].

```
#df["Address"].str.split(',') # didnt make any change
#df["Address"].str.split(',') # didnt make any change
df["Address"].str.split(',') # for creating one column
df["Address"].str.split(',') # for creating two columns
```

Out[296].

	0	1	2
0	123 Shire Lane	Shire	None
1	93 West Main Street	None	None
2	298 Drugs Driveway	None	None
3	980 Paper Avenue Pennsylvania 18503	None	None
4	123 Dragons Road	None	None
5	768 City Parkway	None	None
6	1209 South Street	None	None
7	98 Clue Drive	None	None
8	123 Middle Earth	None	None
9	25th Main Street New York	None	None
10	612 Shire Lane Shire	None	None
11	2394 Hogwarts Avenue	None	None
12	2039 Main Street	None	None
13	343 City Parkway	None	None
14	214 HR Avenue	None	None
15	2395 Hogwarts Avenue	None	None
16	121 Paper Avenue Pennsylvania	None	None
17	3498 Super Lane	None	None
18	N/A	None	None
19	910 Tatooine Road	Tatooine	None

## Assigning a column name for the expanded address

In [297].

```
df[["Street_Address","State","Zip_Code"]] = df["Address"].str.split(',') # for creating one column
df
```

Out[297].

1017	Michael	Scott	123-643-9775	121 Paper Avenue, Pennsylvania	Yes	No	121 Paper Avenue	Pennsylvania	None
1018	Clark	Kent		3498 Super Lane	Y	NaN	3498 Super Lane	None	None
1019	Creed	Braton		N/A	N/A	Yes	N/A	None	None
1020	Anakin	Skywalker	876-678-3469	910 Tatooine Road, Tatooine	Yes	N	910 Tatooine Road	Tatooine	None

Working on paying customer

```
f["Paying Customer"] = df["Paying Customer"].str.replace('Yes','Y')

f["Paying Customer"] = df["Paying Customer"].str.replace('No','N')

f
```

CustomerID	First_Name	Last_Name	Phone_Number	Address	Paying Customer	Do_Not_Contact	Street_Address	State	Zip_Code
1001	Frodo	Baggins	123-545-5421	123 Shire Lane, Shire	Y	No	123 Shire Lane	Shire	None
1002	Abed	Nadri	123-643-9775	93 West Main Street	N	Yes	93 West Main Street	None	None
1003	Walter	White		298 Drugs Driveway	N	NaN	298 Drugs Driveway	None	None
1004	Dwight	Schrute	123-543-2345	980 Paper Avenue, Pennsylvania, 18503	Y	Y	980 Paper Avenue	Pennsylvania	18503
1005	Jon	Snow	876-678-3469	123 Dragons Road	Y	No	123 Dragons Road	None	None
1006	Ron	Swanson	304-762-2467	768 City Parkway	Y	Yes	768 City Parkway	None	None
1007	Jeff	Winger		1209 South Street	N	No	1209 South Street	None	None
1008	Sherlock	Holmes	876-678-3469	98 Clue Drive	N	No	98 Clue Drive	None	None
1009	Gandalf	NaN		123 Middle Earth	Y	NaN	123 Middle Earth	None	None
1010	Peter	Parker	123-545-5421	25th Main Street, New York	Y	No	25th Main Street	New York	None
1011	Sarmise	Gamagee		612 Shire Lane, Shire	Y	No	612 Shire Lane	Shire	None
1012	Harry	Potter		2394 Hogwarts Avenue	Y	NaN	2394 Hogwarts Avenue	None	None
1013	Don	Draper	123-543-2345	2039 Main Street	Y	N	2039 Main Street	None	None