

DSc

March 7, 2024

1 Initialisation

```
[1]: import pandas as pd
import numpy as np
from IPython.display import HTML, Markdown, Latex

def display_df(tp_df=None, index=False):
    tp_df = tp_df if tp_df is not None else df
    display(Markdown(tp_df.to_markdown(index=index)))
```

2 15/02/24

Given a dataset, print the following:

- 1) Records of index 1 & 3
- 2) Records where age ≥ 15
- 3) Records where age ≥ 12 and gender = Male
- 4) City and gender of people with age ≥ 12

```
[2]: data = {
    'age':          [10,22,13,21,12,11,17],
    'section':      ['A','B','C','B','B','A','A'],
    'city':         ['Gurgaon','Delhi','Mumbai','Delhi','Mumbai','Delhi','Mumbai'],
    'gender':       ['M','F','F','M','M','M','F'],
    'favorite_color': ['red','black','yellow','pink','black','green','red']
}
df = pd.DataFrame(data)

print('\nOriginal data:')
display_df()

print('\n1) Records of index 1 & 3')
display_df( df.iloc[ [1,3] , : ] )

print('\n2) Records where age  $\geq 15$ :')
display_df( df.query('age  $\geq 15$ ') )
```

```

print('\n3) Records where age >= 12 and gender = Male:')
display_df( df.query('age >= 12 and gender == "M"') )

print('\n4) City and gender of people with age >= 12:')
display_df( df.query('age >= 12')[['city','gender']] )

```

Original data:

age	section	city	gender	favorite_color
10	A	Gurgaon	M	red
22	B	Delhi	F	black
13	C	Mumbai	F	yellow
21	B	Delhi	M	pink
12	B	Mumbai	M	black
11	A	Delhi	M	green
17	A	Mumbai	F	red

1) Records of index 1 & 3

age	section	city	gender	favorite_color
22	B	Delhi	F	black
21	B	Delhi	M	pink

2) Records where age >= 15:

age	section	city	gender	favorite_color
22	B	Delhi	F	black
21	B	Delhi	M	pink
17	A	Mumbai	F	red

3) Records where age >= 12 and gender = Male:

age	section	city	gender	favorite_color
21	B	Delhi	M	pink
12	B	Mumbai	M	black

4) City and gender of people with age >= 12:

city	gender
Delhi	F
Mumbai	F
Delhi	M
Mumbai	M
Mumbai	F

3 22/02/24

Create a dataframe to store data of 10 students, with the columns being “Name”, “Age”, “Semester I marks out of 600”, “Semester II marks out of 500”, and “Attendance”

- 1) Display details of students who scored more than 560 marks in sem I
- 2) Display details of students who scored less than 250 marks in sem II
- 3) Display details of student who scored minimum marks in sem II
- 4) Display details of student who scored maximum marks in sem II
- 5) Display details of students whose attendance is more than 75
- 6) Display details of students whose attendance is less than 50
- 7) Insert 2 new records
- 8) Add a column corresponding to percentage of marks of both semesters
- 9) Add a new column corresponding to grades:

Percentage	Grade
>=90	O
>=75 and <90	A+
>=60 and <75	A
>=50 and <60	B+
>=40 and <50	B
>40	F

```
[3]: data = {
    'Name': ['A', 'B', 'C', 'D', 'E', 'F', 'G', 'H', 'I', 'J'],
    'Age': [20, 21, 20, 22, 23, 20, 21, 22, 20, 21],
    'Semester I marks out of 600': [213, 31, 57, 406, 417, 45, 217, 200, 588, 319],
    'Semester II marks out of 500': [198, 378, 133, 450, 283, 485, 193, 283, 236, 191],
    'Attendance': [76, 26, 53, 32, 50, 67, 92, 62, 44, 85]
}
df = pd.DataFrame(data)

print('\nOriginal data:')
display_df()
```

```

print('\n1) Students who scored more than 560 marks in sem I:')
ans = df.query('`Semester I marks out of 600` > 560')
display_df(ans, index=True)

print('\n2) Students who scored less than 250 marks in sem II:')
ans = df.query('`Semester II marks out of 500` < 250')
display_df(ans, index=True)

print('\n3) Student who scored minimum marks in sem II:')
min_marks = min(df['Semester II marks out of 500'])
ans = df.query('`Semester II marks out of 500` == @min_marks')
display_df( ans , index=True )

print('\n4) Student who scored maximum marks in sem II:')
ans = df.sort_values(by='Semester II marks out of 500',ascending=False).head(1)
display_df(ans, index=True)

print('\n5) Students whose attendance is more than 75:')
ans = df.query('Attendance > 75')
display_df(ans, index=True)

print('\n6) Students whose attendance is less than 50:')
ans = df.query('Attendance < 50')
display_df(ans, index=True)

print('\n7) Inserted two new records:')
new_data = {
    'Name':          ['K', 'L'],
    'Age':            [22, 23],
    'Semester I marks out of 600': [300, 400],
    'Semester II marks out of 500': [400, 300],
    'Attendance':     [80, 40]
}
new_df = pd.DataFrame(new_data)
df = pd.concat([df,new_df], ignore_index=True)
display_df(index=True)

print('\n8) Added the percentage column:')
df['Percentage'] = (df['Semester I marks out of 600'] + df['Semester II marks_
out of 500']) / 11
df['Percentage'] = df['Percentage'].apply(lambda x: round(x,2))
display_df()

print('\n9) Added the grade column:')
def get_grade(x: float):
    if x >= 90: return 'O'
    elif x >= 75: return 'A+'

```

```

elif x >= 60: return 'A'
elif x >= 50: return 'B+'
elif x >= 40: return 'B'
else: return 'F'
df['Grade'] = df['Percentage'].apply(get_grade)
display_df()

```

Original data:

Name	Age	Semester I marks out of 600	Semester II marks out of 500	Attendance
A	20	213	198	76
B	21	31	378	26
C	20	57	133	53
D	22	406	450	32
E	23	417	283	50
F	20	45	485	67
G	21	217	193	92
H	22	200	283	62
I	20	588	236	44
J	21	319	191	85

1) Students who scored more than 560 marks in sem I:

Name	Age	Semester I marks out of 600	Semester II marks out of 500	Attendance
8 I	20	588	236	44

2) Students who scored less than 250 marks in sem II:

Name	Age	Semester I marks out of 600	Semester II marks out of 500	Attendance
0 A	20	213	198	76
2 C	20	57	133	53
6 G	21	217	193	92
8 I	20	588	236	44
9 J	21	319	191	85

3) Student who scored minimum marks in sem II:

	Name	Age	Semester I marks out of 600	Semester II marks out of 500	Attendance
2	C	20	57	133	53

4) Student who scored maximum marks in sem II:

	Name	Age	Semester I marks out of 600	Semester II marks out of 500	Attendance
5	F	20	45	485	67

5) Students whose attendance is more than 75:

	Name	Age	Semester I marks out of 600	Semester II marks out of 500	Attendance
0	A	20	213	198	76
6	G	21	217	193	92
9	J	21	319	191	85

6) Students whose attendance is less than 50:

	Name	Age	Semester I marks out of 600	Semester II marks out of 500	Attendance
1	B	21	31	378	26
3	D	22	406	450	32
8	I	20	588	236	44

7) Inserted two new records:

	Name	Age	Semester I marks out of 600	Semester II marks out of 500	Attendance
0	A	20	213	198	76
1	B	21	31	378	26
2	C	20	57	133	53
3	D	22	406	450	32
4	E	23	417	283	50
5	F	20	45	485	67
6	G	21	217	193	92
7	H	22	200	283	62
8	I	20	588	236	44
9	J	21	319	191	85

	Name	Age	Semester I marks out of 600	Semester II marks out of 500	Attendance
10	K	22	300	400	80
11	L	23	400	300	40

8) Added the percentage column:

	Name	Age	Semester I marks out of 600	Semester II marks out of 500	Attendance	Percentage
	A	20	213	198	76	37.36
	B	21	31	378	26	37.18
	C	20	57	133	53	17.27
	D	22	406	450	32	77.82
	E	23	417	283	50	63.64
	F	20	45	485	67	48.18
	G	21	217	193	92	37.27
	H	22	200	283	62	43.91
	I	20	588	236	44	74.91
	J	21	319	191	85	46.36
	K	22	300	400	80	63.64
	L	23	400	300	40	63.64

9) Added the grade column:

	Name	Age	Semester I marks out of 600	Semester II marks out of 500	Attendance	Percentage	Grade
	A	20	213	198	76	37.36	F
	B	21	31	378	26	37.18	F
	C	20	57	133	53	17.27	F
	D	22	406	450	32	77.82	A+
	E	23	417	283	50	63.64	A
	F	20	45	485	67	48.18	B
	G	21	217	193	92	37.27	F
	H	22	200	283	62	43.91	B
	I	20	588	236	44	74.91	A
	J	21	319	191	85	46.36	B
	K	22	300	400	80	63.64	A
	L	23	400	300	40	63.64	A

4 29/02/24

1. Create a DataFrame based on E-Commerce data and generate mean, mode, and median
2. Write a program to implement pivot() and pivot-table() on a DataFrame

3. Write a Program to read a CSV file and create its DataFrame
4. Consider the DataFrame QtrSales where each row contains the item category, item name and expenditure and group the rows by category, and print the average expenditure per category
5. Create a DataFrame having age, name, weight of five students. Write a program to display only the weight of first and fourth rows
6. Write a program to create a DataFrame to store weight, age and name of three people. Print the DataFrame and its transpose

```
[4]: # 1
data = {
    'Order_ID':    ['101', '102', '103', '104', '105'],
    'Price':       [50, 20, 40, 50, 45],
    'Quantity':    [2, 3, 1, 2, 1]
}
df = pd.DataFrame(data)
print('\n\nn1.1) Original dataframe:')
display_df(df)

mean_df    = df.mean(numeric_only=True)
mode_df     = df.mode(numeric_only=True)
median_df  = df.median(numeric_only=True)

print('\nn1.2) Means: ')
display_df(mean_df, index=True)
print('\nn1.3) Modes: ')
display_df(mode_df, index=True)
print('\nn1.4) Medians: ')
display_df(median_df, index=True)

# 2
data = {
    'Day':         ['Monday', 'Monday', 'Tuesday', 'Tuesday', 'Wednesday', 'Wednesday'],
    'City':         ['Delhi', 'Mumbai', 'Delhi', 'Mumbai', 'Delhi', 'Mumbai'],
    'Temperature': [32, 34, 33, 35, 34, 36],
}
df = pd.DataFrame(data)
print('\n\nn2.1) Dataframe for pivot(): ')
display_df(df)

pivot_df = df.pivot(index='Day', columns='City', values='Temperature')
print('\nn2.2) pivot(): ')
display_df(pivot_df, index=True)

data = {
```



```

        'Day':          ['Monday', 'Monday', 'Monday', 'Tuesday', 'Tuesday',
↪ 'Tuesday'],
        'City':         ['Delhi', 'Delhi', 'Mumbai', 'Delhi', 'Mumbai', 'Mumbai'],
        'Temperature': [32, 33, 36, 33, 36, 37],
    }
    df = pd.DataFrame(data)
    print('\n2.3) Dataframe for pivot_table(): ')
    display_df(df)

    pivot_table_df = df.pivot_table(index='Day', columns='City',
↪ values='Temperature', aggfunc='count')
    print('\n2.4) pivot_table(): ')
    display_df(pivot_table_df, index=True)

# 3
filename = 'data.csv'
print(f'\n\n3.1) Contents of {filename}:')
with open(filename) as f:
    print(f.read())

df = pd.read_csv(filename)
print('\n3.2) Dataframe:')
display_df(df)

# 4
QtrSales = pd.DataFrame({
    'category':        ['Electronics', 'Electronics', 'Fashion', 'Fashion',
↪ 'Electronics', 'Fashion'],
    'item_name':        ['Laptop', 'Headphones', 'T-Shirt', 'Jeans', 'Smartphone',
↪ 'Shoes'],
    'expenditure': [1200, 100, 31, 50, 800, 60]
})
print('\n\n4.1) Original dataframe:')
display_df(QtrSales)

grouped = QtrSales.groupby(by='category')['expenditure']
mean_df = grouped.mean()
print('\n4.2) Average expenditure per category:')
display_df(mean_df, index=True)

# 5

```

```

data = {
    'name':    ['John', 'Emma', 'Michael', 'Sophia', 'William'],
    'age':     [20, 21, 22, 20, 23],
    'weight':  [70, 65, 75, 68, 72]
}
df = pd.DataFrame(data)
print('\n\n5.1) Original dataframe:')
display_df(df, index=True)

weight_df = df.loc[ [0,3] , ['weight'] ]
print('\n5.2) Weight of the first and fourth rows:')
display_df(weight_df, index=True)

# 6
data = {
    'name':    ['John', 'Emma', 'Michael', 'Sophia', 'William'],
    'age':     [20, 21, 22, 20, 23],
    'weight':  [70, 65, 75, 68, 72]
}
df = pd.DataFrame(data)
print('\n\n6.1) Original dataframe:')
display_df(df, index=True)

print('\n6.2) Transpose of the dataframe:')
display_df(df.T, index=True)

```

1.1) Original dataframe:

Order_ID	Price	Quantity
101	50	2
102	20	3
103	40	1
104	50	2
105	45	1

1.2) Means:

	0
Price	41
Quantity	1.8

1.3) Modes:

	Price	Quantity
0	50	1
1	nan	2

1.4) Medians:

	0
Price	45
Quantity	2

2.1) Dataframe for pivot():

Day	City	Temperature
Monday	Delhi	32
Monday	Mumbai	34
Tuesday	Delhi	33
Tuesday	Mumbai	35
Wednesday	Delhi	34
Wednesday	Mumbai	36

2.2) pivot():

Day	Delhi	Mumbai
Monday	32	34
Tuesday	33	35
Wednesday	34	36

2.3) Dataframe for pivot_table():

Day	City	Temperature
Monday	Delhi	32
Monday	Delhi	33
Monday	Mumbai	36
Tuesday	Delhi	33

Day	City	Temperature
Tuesday	Mumbai	36
Tuesday	Mumbai	37

2.4) `pivot_table()`:

Day	Delhi	Mumbai
Monday	2	1
Tuesday	1	2

3.1) Contents of `data.csv`:

Name, Age, Gender

Ram, 16, M

Manish, 18, M

Sahil, 15, M

Amrit, 20, F

Mark, 19, M

3.2) Dataframe:

Name	Age	Gender
Ram	16	M
Manish	18	M
Sahil	15	M
Amrit	20	F
Mark	19	M

4.1) Original dataframe:

category	item_name	expenditure
Electronics	Laptop	1200
Electronics	Headphones	100
Fashion	T-Shirt	31
Fashion	Jeans	50
Electronics	Smartphone	800
Fashion	Shoes	60

4.2) Average expenditure per category:

category	expenditure
Electronics	700
Fashion	47

5.1) Original dataframe:

	name	age	weight
0	John	20	70
1	Emma	21	65
2	Michael	22	75
3	Sophia	20	68
4	William	23	72

5.2) Weight of the first and fourth rows:

	weight
0	70
3	68

6.1) Original dataframe:

	name	age	weight
0	John	20	70
1	Emma	21	65
2	Michael	22	75
3	Sophia	20	68
4	William	23	72

6.2) Transpose of the dataframe:

	0	1	2	3	4
name	John	Emma	Michael	Sophia	William

	0	1	2	3	4
age	20	21	22	20	23
weight	70	65	75	68	72