

# Lab+3+-+Data+Preparation

February 5, 2025

## 1 Lab 3: Data Preparation

CPE232 Data Models

---

## 2 [1] Reviews on Pandas

1.1) Discover

- methods to explore and understand your DataFrame

```
[121]: import pandas as pd
```

```
df = pd.read_csv('nss15.csv')
```

```
[122]: # see the shape of the dataframe
print(df.shape)
```

```
(334839, 12)
```

```
[123]: # seeing the summary of the dataframe
print(df.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 334839 entries, 0 to 334838
Data columns (total 12 columns):
#   Column          Non-Null Count  Dtype
---  -
0   caseNumber      334839 non-null  int64
1   treatmentDate   334839 non-null  object
2   statWeight      334839 non-null  float64
3   stratum         334839 non-null  object
4   age             334839 non-null  int64
5   sex             334837 non-null  object
6   race            205014 non-null  object
7   diagnosis       334839 non-null  int64
8   bodyPart        334839 non-null  int64
9   disposition     334839 non-null  int64
10  location        334839 non-null  int64
```

```

11 product          334839 non-null int64
dtypes: float64(1), int64(7), object(4)
memory usage: 30.7+ MB
None

```

```

[124]: # seeing the stats of the column in dataframe
print(df.describe())

```

```

      caseNumber  statWeight  age  diagnosis \
count  3.348390e+05  334839.000000  334839.000000  334839.000000
mean    1.510271e+08    39.343028    31.385451    60.154591
std     1.720330e+06    34.142933    26.105098     6.170699
min     1.501032e+08     4.965500     0.000000    41.000000
25%     1.504405e+08    15.059100    10.000000    57.000000
50%     1.507358e+08    15.776200    23.000000    59.000000
75%     1.510231e+08    74.881300    51.000000    64.000000
max     1.603418e+08    97.923900   107.000000    74.000000

      bodyPart  disposition  location  product
count  334839.000000  334839.000000  334839.000000  334839.000000
mean     64.374192     1.307930     2.485451    2098.900854
std     24.002331     0.977627     3.217617    1332.222670
min       0.000000     1.000000     0.000000    106.000000
25%     35.000000     1.000000     0.000000    1211.000000
50%     75.000000     1.000000     1.000000    1807.000000
75%     82.000000     1.000000     5.000000    3265.000000
max     94.000000     9.000000     9.000000    5555.000000

```

```

[125]: # seeing the first 5 rows of the dataframe
print(df.head())

```

```

      caseNumber  treatmentDate  statWeight  stratum  age  sex  race \
0   150733174    7/11/2015    15.7762      V      5   Male  NaN
1   150734723    7/6/2015    83.2157      S     36   Male  White
2   150817487    8/2/2015    74.8813      L     20  Female  NaN
3   150717776    6/26/2015    15.7762      V     61   Male  NaN
4   150721694    7/4/2015    74.8813      L     88  Female  Other

      diagnosis  bodyPart  disposition  location  product
0           57        33           1         9      1267
1           57        34           1         1      1439
2           71        94           1         0      3274
3           71        35           1         0        611
4           62        75           1         0      1893

```

```

[126]: # seeing the last 5 rows of the dataframe
print(df.tail())

```

```

      caseNumber  treatmentDate  statWeight  stratum  age  sex  race \

```

334834	150739278	5/31/2015	15.0591	V	7	Male	NaN
334835	150733393	7/11/2015	5.6748	C	3	Female	Black
334836	150819286	7/24/2015	15.7762	V	38	Male	NaN
334837	150823002	8/8/2015	97.9239	M	38	Female	White
334838	150723074	6/20/2015	49.2646	M	5	Female	White

	diagnosis	bodyPart	disposition	location	product
334834	59	76	1	1	1864
334835	68	85	1	0	1931
334836	71	79	1	0	3250
334837	59	82	1	1	464
334838	57	34	1	9	3273

```
[127]: # seeing the list of columns in the dataframe
print(df.columns)
```

```
Index(['caseNumber', 'treatmentDate', 'statWeight', 'stratum', 'age', 'sex',
      'race', 'diagnosis', 'bodyPart', 'disposition', 'location', 'product'],
      dtype='object')
```

## 1.2) Selecting variables

- select specific columns from the DataFrame to create a new DataFrame with only those columns

```
[128]: df['age']
```

```
[128]: 0      5
      1     36
      2     20
      3     61
      4     88
      ..
334834    7
334835    3
334836   38
334837   38
334838    5
Name: age, Length: 334839, dtype: int64
```

```
[129]: df['age'].head()
```

```
[129]: 0      5
      1     36
      2     20
      3     61
      4     88
Name: age, dtype: int64
```

```
[130]: df[['caseNumber', 'age']]
```

```
[130]:      caseNumber  age
0      150733174    5
1      150734723   36
2      150817487   20
3      150717776   61
4      150721694   88
...
334834  150739278    7
334835  150733393    3
334836  150819286   38
334837  150823002   38
334838  150723074    5
```

[334839 rows x 2 columns]

```
[131]: # select columns based on the data type
df.select_dtypes(include=['number'])
```

```
[131]:      caseNumber  statWeight  age  diagnosis  bodyPart  disposition \
0      150733174    15.7762    5         57         33           1
1      150734723    83.2157   36         57         34           1
2      150817487    74.8813   20         71         94           1
3      150717776    15.7762   61         71         35           1
4      150721694    74.8813   88         62         75           1
...
334834  150739278    15.0591    7         59         76           1
334835  150733393     5.6748    3         68         85           1
334836  150819286    15.7762   38         71         79           1
334837  150823002    97.9239   38         59         82           1
334838  150723074    49.2646    5         57         34           1
```

```
      location  product
0             9    1267
1             1    1439
2             0    3274
3             0     611
4             0    1893
...
334834      ...      ...
334834      1    1864
334835      0    1931
334836      0    3250
334837      1     464
334838      9    3273
```

[334839 rows x 8 columns]

```
[132]: # select row by .loc
df.loc[0]
```

```
[132]: caseNumber      150733174
      treatmentDate   7/11/2015
      statWeight      15.7762
      stratum         V
      age             5
      sex             Male
      race            NaN
      diagnosis       57
      bodyPart        33
      disposition     1
      location        9
      product         1267
      Name: 0, dtype: object
```

```
[133]: # select column by .loc
df.loc[:, 'treatmentDate': 'diagnosis']
```

```
[133]:  treatmentDate  statWeight  stratum  age  sex  race  diagnosis
0      7/11/2015      15.7762      V    5  Male  NaN        57
1      7/6/2015      83.2157      S   36  Male  White       57
2      8/2/2015      74.8813      L   20  Female  NaN       71
3      6/26/2015      15.7762      V   61  Male  NaN       71
4      7/4/2015      74.8813      L   88  Female  Other       62
5      7/2/2015       5.6748      C    1  Female  White       71
6      6/8/2015      15.7762      V   25  Male  Black       51
```

```
[134]: df.loc[df['age']>80, ['treatmentDate', 'age']]
```

```
[134]:  treatmentDate  age
4      7/4/2015   88
8      7/16/2015  98
39     5/3/2015   88
46     4/15/2015  91
63     1/12/2015  97
...
334701  4/27/2015  86
334784  7/7/2015  82
334785  7/11/2015  86
334815  10/28/2015 85
334819  1/13/2015 85
```

```
[20422 rows x 2 columns]
```

```
[135]: # select row by .iloc
df.iloc[0:5]
```

```
[135]:
```

	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	\
0	150733174	7/11/2015	15.7762	V	5	Male	NaN	
1	150734723	7/6/2015	83.2157	S	36	Male	White	
2	150817487	8/2/2015	74.8813	L	20	Female	NaN	
3	150717776	6/26/2015	15.7762	V	61	Male	NaN	
4	150721694	7/4/2015	74.8813	L	88	Female	Other	

  

	diagnosis	bodyPart	disposition	location	product
0	57	33	1	9	1267
1	57	34	1	1	1439
2	71	94	1	0	3274
3	71	35	1	0	611
4	62	75	1	0	1893

```
[136]: # select column by .iloc
df.iloc[:, [0,1,2,3,4]]
```

```
[136]:
```

	caseNumber	treatmentDate	statWeight	stratum	age
0	150733174	7/11/2015	15.7762	V	5
1	150734723	7/6/2015	83.2157	S	36
2	150817487	8/2/2015	74.8813	L	20
3	150717776	6/26/2015	15.7762	V	61
4	150721694	7/4/2015	74.8813	L	88
...	...	...	...	...	...
334834	150739278	5/31/2015	15.0591	V	7
334835	150733393	7/11/2015	5.6748	C	3
334836	150819286	7/24/2015	15.7762	V	38
334837	150823002	8/8/2015	97.9239	M	38
334838	150723074	6/20/2015	49.2646	M	5

[334839 rows x 5 columns]

### 1.3) Filtering the data

```
[137]: # filter rows based on the condition
df[df['age'] > 50]
```

```
[137]:
```

	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	\
3	150717776	6/26/2015	15.7762	V	61	Male	NaN	
4	150721694	7/4/2015	74.8813	L	88	Female	Other	
7	150704114	6/14/2015	83.2157	S	53	Male	White	
8	150736558	7/16/2015	83.2157	S	98	Male	Black	
16	150901411	8/27/2015	83.2157	S	65	Female	White	
...	...	...	...	...	...	...	...	
334811	150702215	6/27/2015	15.7762	V	51	Female	NaN	
334815	151100368	10/28/2015	83.2157	S	85	Female	NaN	
334819	150528367	1/13/2015	49.2646	M	85	Female	NaN	
334826	150648619	6/17/2015	15.7762	V	52	Female	White	

334829	150633526	4/4/2015	49.2646	M	51	Female	NaN
--------	-----------	----------	---------	---	----	--------	-----

	diagnosis	bodyPart	disposition	location	product
3	71	35	1	0	611
4	62	75	1	0	1893
7	57	30	1	0	5040
8	59	76	1	1	1807
16	59	83	1	1	1817
...	...	...	...	...	...
334811	53	83	1	1	1426
334815	57	80	4	1	1807
334819	57	79	5	1	676
334826	64	30	1	1	1842
334829	56	92	1	1	1616

[85235 rows x 12 columns]

```
[138]: # filter coloum based on column name
df.filter(like='age')
```

```
[138]:      age
0      5
1     36
2     20
3     61
4     88
...
334834  7
334835  3
334836 38
334837 38
334838  5
```

[334839 rows x 1 columns]

#### 1.4) Sorting

- Sort the DataFrame by its index based on column

```
[139]: # sort the dataframe based on column name and ascending order
df.sort_values(by='statWeight', ascending=False)
```

```
[139]:      caseNumber  treatmentDate  statWeight  stratum  age  sex  race  \
275174  150343700      3/9/2015      97.9239         M   48  Male   NaN
36      151029422     10/6/2015      97.9239         M   37  Male  White
334806  150612491     5/29/2015      97.9239         M   18  Female White
334810  150725804     7/8/2015      97.9239         M   33  Female Black
275161  150450816     4/13/2015      97.9239         M   24  Male  White
```

...	...	...	...	...	...	...	...
44011	160222258	12/29/2015	4.9655	C	2	Female	Other
325320	151213065	11/29/2015	4.9655	C	16	Female	White
43891	160113865	12/28/2015	4.9655	C	4	Male	White
43628	151130111	11/9/2015	4.9655	C	13	Male	Black
43523	151139237	11/16/2015	4.9655	C	2	Female	Black

	diagnosis	bodyPart	disposition	location	product
275174	57	93	1	1	281
36	64	35	1	0	1267
334806	59	92	1	1	845
334810	71	94	1	0	1616
275161	71	37	1	1	3286
...	...	...	...	...	...
44011	71	92	1	1	1893
325320	62	75	1	8	3254
43891	59	76	1	1	1842
43628	53	33	1	0	5011
43523	57	80	1	0	679

[334839 rows x 12 columns]

```
[140]: # sort the index of the dataframe
df.sort_index()
```

```
[140]:
```

	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	\
0	150733174	7/11/2015	15.7762	V	5	Male	NaN	
1	150734723	7/6/2015	83.2157	S	36	Male	White	
2	150817487	8/2/2015	74.8813	L	20	Female	NaN	
3	150717776	6/26/2015	15.7762	V	61	Male	NaN	
4	150721694	7/4/2015	74.8813	L	88	Female	Other	
...	...	...	...	...	...	...	...	
334834	150739278	5/31/2015	15.0591	V	7	Male	NaN	
334835	150733393	7/11/2015	5.6748	C	3	Female	Black	
334836	150819286	7/24/2015	15.7762	V	38	Male	NaN	
334837	150823002	8/8/2015	97.9239	M	38	Female	White	
334838	150723074	6/20/2015	49.2646	M	5	Female	White	

	diagnosis	bodyPart	disposition	location	product
0	57	33	1	9	1267
1	57	34	1	1	1439
2	71	94	1	0	3274
3	71	35	1	0	611
4	62	75	1	0	1893
...	...	...	...	...	...
334834	59	76	1	1	1864
334835	68	85	1	0	1931



334836	71	79	1	0	3250
334837	59	82	1	1	464
334838	57	34	1	9	3273

[334839 rows x 12 columns]

### 1.5) Add/Remove

- This section shows how to manipulate the DataFrame's structure

```
[141]: # Dropping the column
df.drop(columns=['disposition'])
```

```
[141]:
```

	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	\
0	150733174	7/11/2015	15.7762	V	5	Male	NaN	
1	150734723	7/6/2015	83.2157	S	36	Male	White	
2	150817487	8/2/2015	74.8813	L	20	Female	NaN	
3	150717776	6/26/2015	15.7762	V	61	Male	NaN	
4	150721694	7/4/2015	74.8813	L	88	Female	Other	
...	...	...	...	...	...	...	...	
334834	150739278	5/31/2015	15.0591	V	7	Male	NaN	
334835	150733393	7/11/2015	5.6748	C	3	Female	Black	
334836	150819286	7/24/2015	15.7762	V	38	Male	NaN	
334837	150823002	8/8/2015	97.9239	M	38	Female	White	
334838	150723074	6/20/2015	49.2646	M	5	Female	White	

  

	diagnosis	bodyPart	location	product
0	57	33	9	1267
1	57	34	1	1439
2	71	94	0	3274
3	71	35	0	611
4	62	75	0	1893
...	...	...	...	...
334834	59	76	1	1864
334835	68	85	0	1931
334836	71	79	0	3250
334837	59	82	1	464
334838	57	34	9	3273

[334839 rows x 11 columns]

```
[142]: # Adding column and create into a new column
df.assign(new_column=df['diagnosis'] + df['bodyPart'])
```

```
[142]:
```

	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	\
0	150733174	7/11/2015	15.7762	V	5	Male	NaN	
1	150734723	7/6/2015	83.2157	S	36	Male	White	
2	150817487	8/2/2015	74.8813	L	20	Female	NaN	

3	150717776	6/26/2015	15.7762	V	61	Male	NaN
4	150721694	7/4/2015	74.8813	L	88	Female	Other
...	...	...	...	...	...	...	...
334834	150739278	5/31/2015	15.0591	V	7	Male	NaN
334835	150733393	7/11/2015	5.6748	C	3	Female	Black
334836	150819286	7/24/2015	15.7762	V	38	Male	NaN
334837	150823002	8/8/2015	97.9239	M	38	Female	White
334838	150723074	6/20/2015	49.2646	M	5	Female	White

	diagnosis	bodyPart	disposition	location	product	new_column
0	57	33	1	9	1267	90
1	57	34	1	1	1439	91
2	71	94	1	0	3274	165
3	71	35	1	0	611	106
4	62	75	1	0	1893	137
...	...	...	...	...	...	...
334834	59	76	1	1	1864	135
334835	68	85	1	0	1931	153
334836	71	79	1	0	3250	150
334837	59	82	1	1	464	141
334838	57	34	1	9	3273	91

[334839 rows x 13 columns]

```
[143]: # Removing the column and assigning it to a new variable
ages = df.pop('age')
```

```
[144]: df
```

```
[144]:
```

	caseNumber	treatmentDate	statWeight	stratum	sex	race	\
0	150733174	7/11/2015	15.7762	V	Male	NaN	
1	150734723	7/6/2015	83.2157	S	Male	White	
2	150817487	8/2/2015	74.8813	L	Female	NaN	
3	150717776	6/26/2015	15.7762	V	Male	NaN	
4	150721694	7/4/2015	74.8813	L	Female	Other	
...	...	...	...	...	...	...	...
334834	150739278	5/31/2015	15.0591	V	Male	NaN	
334835	150733393	7/11/2015	5.6748	C	Female	Black	
334836	150819286	7/24/2015	15.7762	V	Male	NaN	
334837	150823002	8/8/2015	97.9239	M	Female	White	
334838	150723074	6/20/2015	49.2646	M	Female	White	

  

	diagnosis	bodyPart	disposition	location	product
0	57	33	1	9	1267
1	57	34	1	1	1439
2	71	94	1	0	3274
3	71	35	1	0	611

4	62	75	1	0	1893
...	...	...	...	...	...
334834	59	76	1	1	1864
334835	68	85	1	0	1931
334836	71	79	1	0	3250
334837	59	82	1	1	464
334838	57	34	1	9	3273

[334839 rows x 11 columns]

#### 1.6) Clean missing

- to remove rows with missing values or replace missing values with a specified value

```
[145]: # replacing the missing values with a specified value
df.fillna(value=0)
```

```
[145]:
```

	caseNumber	treatmentDate	statWeight	stratum	sex	race	\
0	150733174	7/11/2015	15.7762	V	Male	0	
1	150734723	7/6/2015	83.2157	S	Male	White	
2	150817487	8/2/2015	74.8813	L	Female	0	
3	150717776	6/26/2015	15.7762	V	Male	0	
4	150721694	7/4/2015	74.8813	L	Female	Other	
...	...	...	...	...	...	...	
334834	150739278	5/31/2015	15.0591	V	Male	0	
334835	150733393	7/11/2015	5.6748	C	Female	Black	
334836	150819286	7/24/2015	15.7762	V	Male	0	
334837	150823002	8/8/2015	97.9239	M	Female	White	
334838	150723074	6/20/2015	49.2646	M	Female	White	

	diagnosis	bodyPart	disposition	location	product
0	57	33	1	9	1267
1	57	34	1	1	1439
2	71	94	1	0	3274
3	71	35	1	0	611
4	62	75	1	0	1893
...	...	...	...	...	...
334834	59	76	1	1	1864
334835	68	85	1	0	1931
334836	71	79	1	0	3250
334837	59	82	1	1	464
334838	57	34	1	9	3273

[334839 rows x 11 columns]

```
[146]: # Remove the rows with missing values
df.dropna()
```

```
[146]:
```

	caseNumber	treatmentDate	statWeight	stratum	sex	race	\
1	150734723	7/6/2015	83.2157	S	Male	White	
4	150721694	7/4/2015	74.8813	L	Female	Other	
5	150721815	7/2/2015	5.6748	C	Female	White	
6	150713483	6/8/2015	15.7762	V	Male	Black	
7	150704114	6/14/2015	83.2157	S	Male	White	
...	...	...	...	...	...	...	
334830	150628863	6/8/2015	15.7762	V	Female	White	
334831	150607637	5/22/2015	5.6748	C	Female	Black	
334835	150733393	7/11/2015	5.6748	C	Female	Black	
334837	150823002	8/8/2015	97.9239	M	Female	White	
334838	150723074	6/20/2015	49.2646	M	Female	White	

  

	diagnosis	bodyPart	disposition	location	product
1	57	34	1	1	1439
4	62	75	1	0	1893
5	71	76	1	1	1715
6	51	33	4	9	1138
7	57	30	1	0	5040
...	...	...	...	...	...
334830	64	79	1	1	1522
334831	59	94	1	0	1616
334835	68	85	1	0	1931
334837	59	82	1	1	464
334838	57	34	1	9	3273

[205014 rows x 11 columns]

### 3 [2] Data Cleaning and Preparation

#### 3.0.1 .isnull, .dropna, .fillna

2.1) checking

```
[147]: df.columns
```

```
[147]: Index(['caseNumber', 'treatmentDate', 'statWeight', 'stratum', 'sex', 'race',
        'diagnosis', 'bodyPart', 'disposition', 'location', 'product'],
        dtype='object')
```

```
[148]: # isnull checking
df.isnull().sum()
```

```
[148]: caseNumber      0
        treatmentDate  0
        statWeight     0
```

```

stratum          0
sex              2
race            129825
diagnosis        0
bodyPart         0
disposition      0
location         0
product          0
dtype: int64

```

```
[149]: # percentage of missing values for the race
df.race.isnull().sum()/df.shape[0]*100
```

```
[149]: np.float64(38.772365226272925)
```

```
[150]: df.shape[0]
```

```
[150]: 334839
```

## 2.2) Drop column

```
[151]: # remove column by using
df = df.drop(columns=['race'])
```

```
[152]: df.head()
```

```
[152]:
```

	caseNumber	treatmentDate	statWeight	stratum	sex	diagnosis	bodyPart	\
0	150733174	7/11/2015	15.7762	V	Male	57	33	
1	150734723	7/6/2015	83.2157	S	Male	57	34	
2	150817487	8/2/2015	74.8813	L	Female	71	94	
3	150717776	6/26/2015	15.7762	V	Male	71	35	
4	150721694	7/4/2015	74.8813	L	Female	62	75	

```

disposition  location  product
0            1         9    1267
1            1         1    1439
2            1         0    3274
3            1         0     611
4            1         0    1893

```

## 2.3) Data imputation

```
[153]: # fillna
# df['age'] = df['age'].fillna(df['age'].median())
```

[Q1] From the above cell, Why it showing an error?

Ans: Age col had been removed in 1st section, no key to work with

[Q2] Fix the error from Q1 problem.

```
[154]: # [Q2]

# hint: see the cell that run `df.pop()`
df["age"] = ages
# use the data that you have saved in the variable ages

# fillna again
df['age'] = df['age'].fillna(df['age'].median())

df.head()
```

```
[154]:  caseNumber treatmentDate  statWeight stratum  sex  diagnosis  bodyPart \
0    150733174      7/11/2015    15.7762      V   Male         57         33
1    150734723      7/6/2015    83.2157      S   Male         57         34
2    150817487      8/2/2015    74.8813      L  Female         71         94
3    150717776      6/26/2015    15.7762      V   Male         71         35
4    150721694      7/4/2015    74.8813      L  Female         62         75
```

```
      disposition  location  product  age
0              1         9    1267    5
1              1         1    1439   36
2              1         0    3274   20
3              1         0     611   61
4              1         0    1893   88
```

2.4) Drop row that have missing value

```
[155]: # remove column by using .dropna()
df = df.dropna()
```

```
[156]: df.isnull().sum()
```

```
[156]: caseNumber      0
      treatmentDate  0
      statWeight    0
      stratum       0
      sex           0
      diagnosis     0
      bodyPart      0
      disposition   0
      location      0
      product       0
      age           0
      dtype: int64
```

### 3.0.2 Datetime

#### 2.5) Working with the datetime format

```
[157]: df["treatmentDate"] = pd.to_datetime(df["treatmentDate"], format="%m/%d/%Y")
```

```
[158]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 334837 entries, 0 to 334838
Data columns (total 11 columns):
#   Column          Non-Null Count  Dtype
---  -
0   caseNumber      334837 non-null  int64
1   treatmentDate   334837 non-null  datetime64[ns]
2   statWeight      334837 non-null  float64
3   stratum         334837 non-null  object
4   sex             334837 non-null  object
5   diagnosis       334837 non-null  int64
6   bodyPart        334837 non-null  int64
7   disposition     334837 non-null  int64
8   location        334837 non-null  int64
9   product         334837 non-null  int64
10  age             334837 non-null  int64
dtypes: datetime64[ns](1), float64(1), int64(7), object(2)
memory usage: 30.7+ MB
```

```
[159]: df['Year'] = df['treatmentDate'].dt.year
```

```
[160]: df['Month'] = df['treatmentDate'].dt.month
```

```
[161]: df.head()
```

```
[161]:
```

	caseNumber	treatmentDate	statWeight	stratum	sex	diagnosis	bodyPart	\
0	150733174	2015-07-11	15.7762	V	Male	57	33	
1	150734723	2015-07-06	83.2157	S	Male	57	34	
2	150817487	2015-08-02	74.8813	L	Female	71	94	
3	150717776	2015-06-26	15.7762	V	Male	71	35	
4	150721694	2015-07-04	74.8813	L	Female	62	75	

  

	disposition	location	product	age	Year	Month
0	1	9	1267	5	2015	7
1	1	1	1439	36	2015	7
2	1	0	3274	20	2015	8
3	1	0	611	61	2015	6
4	1	0	1893	88	2015	7

[Q3] Can you change the format to DD/MM/YYYY? Show your work.

```
[162]: # write your code here

# combine the year and month column
df['DDMMYY'] = df['treatmentDate'].dt.strftime('%d/%m/%Y')

df['treatmentDate'] = df['DDMMYY']

# drop the year and month column
df = df.drop(columns=['Year', 'Month', 'DDMMYY'])

df.head()
```

```
[162]:
```

	caseNumber	treatmentDate	statWeight	stratum	sex	diagnosis	bodyPart	\
0	150733174	11/07/2015	15.7762	V	Male	57	33	
1	150734723	06/07/2015	83.2157	S	Male	57	34	
2	150817487	02/08/2015	74.8813	L	Female	71	94	
3	150717776	26/06/2015	15.7762	V	Male	71	35	
4	150721694	04/07/2015	74.8813	L	Female	62	75	

  

	disposition	location	product	age
0	1	9	1267	5
1	1	1	1439	36
2	1	0	3274	20
3	1	0	611	61
4	1	0	1893	88

### 3.0.3 Combine Dataframe by .merge and .concat

#### 2.6 Merge

```
[163]: import pandas as pd

superstore_order = pd.read_csv('superstore_order.csv')
superstore_people = pd.read_csv('superstore_people.csv')
superstore_return = pd.read_csv('superstore_return.csv')

[164]: superstore_order.merge(superstore_return[superstore_return["Returned"]=="Yes"],
    on="Order ID",
    how="inner")\
    [["Customer ID", "Returned"]]\
    .drop_duplicates()
```

```
[164]:
```

	Customer ID	Returned
0	ZD-21925	Yes
3	TB-21055	Yes
10	JS-15685	Yes
13	LC-16885	Yes
20	BS-11755	Yes



```

..      ...      ...
688     ED-13885      Yes
689     TS-21205      Yes
696     MF-17665      Yes
702     SH-19975      Yes
705     RB-19435      Yes

```

[222 rows x 2 columns]

[Q2] What does the argument `how="inner"` do?

Ans: Select rows that Share the same condition

[Q3] In your opinion, what information that the result above conveys?

Ans: Result in a table of return from buyer in data order

More merging...

```

[165]: superstore_order.merge(superstore_return,
    on="Order ID" ,
    how="inner")

```

```

[165]:      Row ID      Order ID  Order Date  Ship Date      Ship Mode \
0         19  CA-2014-143336  27/08/2014  01/09/2014    Second Class
1         20  CA-2014-143336  27/08/2014  01/09/2014    Second Class
2         21  CA-2014-143336  27/08/2014  01/09/2014    Second Class
3         56  CA-2016-111682  17/06/2016  18/06/2016    First Class
4         57  CA-2016-111682  17/06/2016  18/06/2016    First Class
..      ...      ...      ...      ...      ...
702      8870  CA-2017-101805  01/12/2017  06/12/2017    Standard Class
703      8871  CA-2017-101805  01/12/2017  06/12/2017    Standard Class
704      8872  CA-2017-101805  01/12/2017  06/12/2017    Standard Class
705      8873  US-2014-105137  10/10/2014  10/10/2014      Same Day
706      8874  US-2014-105137  10/10/2014  10/10/2014      Same Day

      Customer ID      Customer Name      Segment      Country      City \
0         ZD-21925  Zuschuss Donatelli  Consumer  United States  San Francisco
1         ZD-21925  Zuschuss Donatelli  Consumer  United States  San Francisco
2         ZD-21925  Zuschuss Donatelli  Consumer  United States  San Francisco
3         TB-21055      Ted Butterfield  Consumer  United States      Troy
4         TB-21055      Ted Butterfield  Consumer  United States      Troy
..      ...      ...      ...      ...      ...
702      SH-19975      Sally Hughsby  Corporate  United States      Seattle
703      SH-19975      Sally Hughsby  Corporate  United States      Seattle
704      SH-19975      Sally Hughsby  Corporate  United States      Seattle
705      RB-19435      Richard Bierner  Consumer  United States      Columbus
706      RB-19435      Richard Bierner  Consumer  United States      Columbus

```

```

... Region      Product ID      Category Sub-Category \

```

0	...	West	OFF-AR-10003056	Office Supplies	Art
1	...	West	TEC-PH-10001949	Technology	Phones
2	...	West	OFF-BI-10002215	Office Supplies	Binders
3	...	East	OFF-ST-10000604	Office Supplies	Storage
4	...	East	OFF-PA-10001569	Office Supplies	Paper
..	...	...	...	...	...
702	...	West	OFF-BI-10002003	Office Supplies	Binders
703	...	West	FUR-FU-10000023	Furniture	Furnishings
704	...	West	OFF-ST-10002756	Office Supplies	Storage
705	...	East	TEC-MA-10002694	Technology	Machines
706	...	East	OFF-BI-10002429	Office Supplies	Binders

		Product Name	Sales	Quantity \
0		Newell 341	8.560	2
1		Cisco SPA 501G IP Phone	213.480	3
2		Wilson Jones Hanging View Binder White 1	22.720	4
3		Home/Office Personal File Carts	208.560	6
4		Xerox 232	32.400	5
..		...	...	...
702		Ibico Presentation Index for Binding Systems	15.920	5
703		Eldon Wave Desk Accessories	70.680	12
704		Tennsco Stur-D-Stor Boltless Shelving 5 Shelve...	541.240	4
705		Hewlett-Packard Deskjet F4180 All-in-One Color...	101.994	2
706		Premier Elliptical Ring Binder Black	18.264	2

	Discount	Profit	Returned
0	0.0	2.4824	Yes
1	0.2	16.0110	Yes
2	0.2	7.3840	Yes
3	0.0	52.1400	Yes
4	0.0	15.5520	Yes
..	...	...	...
702	0.2	5.3730	Yes
703	0.0	31.0992	Yes
704	0.0	5.4124	Yes
705	0.7	-71.3958	Yes
706	0.7	-13.3936	Yes

[707 rows x 22 columns]

## 2.7) Concatenate

```
[166]: pd.concat([superstore_order, superstore_people], axis=1, join='inner')
```

[166]:	Row ID	Order ID	Order Date	Ship Date	Ship Mode	Customer ID \
	0	1	CA-2016-152156	08/11/2016	11/11/2016	Second Class CG-12520
	1	2	CA-2016-152156	08/11/2016	11/11/2016	Second Class CG-12520
	2	3	CA-2016-138688	12/06/2016	16/06/2016	Second Class DV-13045

```
3      4 US-2015-108966  11/10/2015  18/10/2015  Standard Class  SO-20335
```

	Customer Name	Segment	Country	City	...	\
0	Claire Gute	Consumer	United States	Henderson	...	
1	Claire Gute	Consumer	United States	Henderson	...	
2	Darrin Van Huff	Corporate	United States	Los Angeles	...	
3	Sean ODonnell	Consumer	United States	Fort Lauderdale	...	

	Product ID	Category	Sub-Category	\
0	FUR-BO-10001798	Furniture	Bookcases	
1	FUR-CH-10000454	Furniture	Chairs	
2	OFF-LA-10000240	Office Supplies	Labels	
3	FUR-TA-10000577	Furniture	Tables	

	Product Name	Sales	Quantity	\
0	Bush Somerset Collection Bookcase	261.9600	2	
1	Hon Deluxe Fabric Upholstered Stacking Chairs ...	731.9400	3	
2	Self-Adhesive Address Labels for Typewriters b...	14.6200	2	
3	Bretford CR4500 Series Slim Rectangular Table	957.5775	5	

	Discount	Profit	Person	Region
0	0.00	41.9136	Anna Andreadi	West
1	0.00	219.5820	Chuck Magee	East
2	0.00	6.8714	Kelly Williams	Central
3	0.45	-383.0310	Cassandra Brandow	South

[4 rows x 23 columns]

```
[167]: pd.concat([superstore_order, superstore_people], axis=1, join='outer')
```

```
[167]:
```

	Row ID	Order ID	Order Date	Ship Date	Ship Mode	\
0	1	CA-2016-152156	08/11/2016	11/11/2016	Second Class	
1	2	CA-2016-152156	08/11/2016	11/11/2016	Second Class	
2	3	CA-2016-138688	12/06/2016	16/06/2016	Second Class	
3	4	US-2015-108966	11/10/2015	18/10/2015	Standard Class	
4	5	US-2015-108966	11/10/2015	18/10/2015	Standard Class	
...	...	...	...	...	...	
8875	8876	US-2016-141264	13/08/2016	19/08/2016	Standard Class	
8876	8877	US-2016-141264	13/08/2016	19/08/2016	Standard Class	
8877	8878	CA-2017-126928	17/09/2017	23/09/2017	Standard Class	
8878	8879	CA-2017-126928	17/09/2017	23/09/2017	Standard Class	
8879	8880	US-2015-107944	23/03/2015	25/03/2015	First Class	

	Customer ID	Customer Name	Segment	Country	City	\
0	CG-12520	Claire Gute	Consumer	United States	Henderson	
1	CG-12520	Claire Gute	Consumer	United States	Henderson	
2	DV-13045	Darrin Van Huff	Corporate	United States	Los Angeles	

3	SO-20335	Sean ODonnell	Consumer	United States	Fort Lauderdale
4	SO-20335	Sean ODonnell	Consumer	United States	Fort Lauderdale
...	...	...	...	...	...
8875	CT-11995	Carol Triggs	Consumer	United States	Irving
8876	CT-11995	Carol Triggs	Consumer	United States	Irving
8877	GZ-14470	Gary Zandusky	Consumer	United States	Morristown
8878	GZ-14470	Gary Zandusky	Consumer	United States	Morristown
8879	AM-10360	Alice McCarthy	Corporate	United States	Los Angeles

...	Product ID	Category	Sub-Category	\
0	... FUR-BO-10001798	Furniture	Bookcases	
1	... FUR-CH-10000454	Furniture	Chairs	
2	... OFF-LA-10000240	Office Supplies	Labels	
3	... FUR-TA-10000577	Furniture	Tables	
4	... OFF-ST-10000760	Office Supplies	Storage	
...	...	...	...	...
8875	... OFF-SU-10003505	Office Supplies	Supplies	
8876	... OFF-AP-10002534	Office Supplies	Appliances	
8877	... TEC-MA-10004626	Technology	Machines	
8878	... OFF-ST-10000615	Office Supplies	Storage	
8879	... OFF-PA-10000659	Office Supplies	Paper	

	Product Name	Sales	Quantity	\
0	Bush Somerset Collection Bookcase	261.9600	2	
1	Hon Deluxe Fabric Upholstered Stacking Chairs ...	731.9400	3	
2	Self-Adhesive Address Labels for Typewriters b...	14.6200	2	
3	Bretford CR4500 Series Slim Rectangular Table	957.5775	5	
4	Eldon Fold N Roll Cart System	22.3680	2	
...	...	...	...	...
8875	Premier Electric Letter Opener	185.3760	2	
8876	3.6 Cubic Foot Counter Height Office Refrigerator	58.9240	1	
8877	Lexmark 20R1285 X6650 Wireless All-in-One Printer	480.0000	4	
8878	SimpliFile Personal File Black Granite 15w x 6...	34.0500	3	
8879	TOPS Carbonless Receipt Book Four 2-3/4 x 7-1/...	192.7200	11	

	Discount	Profit	Person	Region
0	0.00	41.9136	Anna Andreadi	West
1	0.00	219.5820	Chuck Magee	East
2	0.00	6.8714	Kelly Williams	Central
3	0.45	-383.0310	Cassandra Brandow	South
4	0.20	2.5164	NaN	NaN
...	...	...	...	...
8875	0.20	-34.7580	NaN	NaN
8876	0.80	-153.2024	NaN	NaN
8877	0.00	225.6000	NaN	NaN
8878	0.00	9.5340	NaN	NaN
8879	0.00	92.5056	NaN	NaN

[8880 rows x 23 columns]

[Q4] What is the difference between `inner` and `outer` on parameter `join` in `pd.concat`?

Ans: `outer`: Includes all rows `inner`: Includes only the common rows

### 3.0.4 Groupby

```
[168]: superstore_order.groupby(['Segment', 'Ship_
      ↪Mode'])[['Sales', 'Quantity', 'Discount', 'Profit']].sum()
```

```
[168]:
```

		Sales	Quantity	Discount	Profit
Segment	Ship Mode				
Consumer	First Class	138594.9328	2455	110.29	18953.7264
	Same Day	53660.6340	1001	43.85	8555.7193
	Second Class	203605.6822	3489	127.29	24701.9148
	Standard Class	627061.3262	10430	443.05	68864.9892
Corporate	First Class	97720.1209	1670	73.07	12660.2526
	Same Day	41716.5550	366	14.50	1120.9222
	Second Class	130759.9288	2027	71.47	15582.1762
	Standard Class	359359.2109	6203	262.82	49832.6780
Home Office	First Class	76743.8674	924	39.82	11829.8821
	Same Day	20968.5170	343	12.50	3909.3442
	Second Class	77175.1080	1148	37.80	12785.8953
	Standard Class	218325.9795	3595	142.14	27298.5786

[Q5] Describe an information that the result above conveys?

Ans: Result in structure that we group first col that use in command is on LHS follow by other, LHS col is bigger chunk

```
[169]: superstore_order["Profit Ratio"] = superstore_order["Profit"] /
      ↪superstore_order["Sales"]
```

```
[170]: superstore_order.groupby(["Category", "Sub-Category"]).agg(mean_profit_ratio =_
      ↪("Profit Ratio", "mean"))
```

```
[170]:
```

		mean_profit_ratio
Category	Sub-Category	
Furniture	Bookcases	-0.127756
	Chairs	0.045028
	Furnishings	0.140782
	Tables	-0.147916
Office Supplies	Appliances	-0.145513
	Art	0.251678
	Binders	-0.191641
	Envelopes	0.421913
	Fasteners	0.301157

Technology	Labels	0.429984
	Paper	0.425586
	Storage	0.092382
	Supplies	0.104970
	Accessories	0.219012
	Copiers	0.317826
	Machines	-0.059535
	Phones	0.118926

[Q6] Describe an information that the result above conveys?

Ans: Result in mean profit in percentage of each sub-category in category

### 3.0.5 Pivot and Melt

Pivot

```
[180]: superstore_order.pivot_table(index="State", columns="Ship Mode", values="Order_ID",
    ↪aggfunc="count").fillna(0).head(10)
```

```
[180]: Ship Mode          First Class  Same Day  Second Class  Standard Class
State
Alabama                9.0         1.0         18.0         30.0
Arizona              42.0        15.0         22.0        123.0
Arkansas             10.0         2.0          8.0         35.0
California          302.0       106.0        346.0       1000.0
Colorado             43.0         5.0         32.0         95.0
Connecticut          19.0         8.0         11.0         39.0
Delaware             16.0         2.0         13.0         55.0
District of Columbia  0.0         0.0          3.0          7.0
Florida              47.0        25.0         57.0        210.0
Georgia              19.0        15.0         31.0        108.0
```

```
[181]: pivot_table_result = superstore_order.pivot_table(index="State", columns="Ship_
    ↪Mode", values="Order ID", aggfunc="count").fillna(0)
print(pivot_table_result)
```

```
Ship Mode          First Class  Same Day  Second Class  Standard Class
State
Alabama                9.0         1.0         18.0         30.0
Arizona              42.0        15.0         22.0        123.0
Arkansas             10.0         2.0          8.0         35.0
California          302.0       106.0        346.0       1000.0
Colorado             43.0         5.0         32.0         95.0
Connecticut          19.0         8.0         11.0         39.0
Delaware             16.0         2.0         13.0         55.0
District of Columbia  0.0         0.0          3.0          7.0
Florida              47.0        25.0         57.0        210.0
Georgia              19.0        15.0         31.0        108.0
```

Idaho	3.0	0.0	2.0	13.0
Illinois	58.0	24.0	96.0	249.0
Indiana	13.0	3.0	30.0	79.0
Iowa	1.0	1.0	4.0	17.0
Kansas	6.0	1.0	2.0	15.0
Kentucky	12.0	5.0	49.0	62.0
Louisiana	7.0	2.0	14.0	15.0
Maine	0.0	0.0	0.0	5.0
Maryland	18.0	7.0	12.0	63.0
Massachusetts	14.0	4.0	35.0	71.0
Michigan	20.0	16.0	43.0	151.0
Minnesota	9.0	4.0	13.0	59.0
Mississippi	3.0	4.0	7.0	36.0
Missouri	7.0	2.0	20.0	24.0
Montana	1.0	1.0	0.0	13.0
Nebraska	6.0	3.0	6.0	20.0
Nevada	4.0	1.0	12.0	17.0
New Hampshire	2.0	0.0	10.0	13.0
New Jersey	5.0	1.0	20.0	87.0
New Mexico	1.0	0.0	9.0	22.0
New York	155.0	57.0	183.0	606.0
North Carolina	36.0	14.0	40.0	139.0
North Dakota	0.0	0.0	5.0	2.0
Ohio	66.0	47.0	84.0	199.0
Oklahoma	5.0	6.0	7.0	44.0
Oregon	20.0	0.0	15.0	81.0
Pennsylvania	103.0	9.0	78.0	341.0
Rhode Island	16.0	0.0	21.0	16.0
South Carolina	3.0	5.0	18.0	16.0
South Dakota	2.0	0.0	0.0	9.0
Tennessee	21.0	2.0	24.0	118.0
Texas	125.0	37.0	161.0	537.0
Utah	4.0	2.0	19.0	28.0
Vermont	0.0	0.0	1.0	2.0
Virginia	39.0	4.0	33.0	115.0
Washington	56.0	34.0	97.0	265.0
West Virginia	0.0	0.0	0.0	3.0
Wisconsin	12.0	3.0	10.0	66.0
Wyoming	0.0	0.0	0.0	1.0

Melt

```
[182]: melted_result = pd.melt(pivot_table_result.reset_index(), id_vars=["State"],
    ↪var_name="Ship Mode", value_name="Order Count")
print(melted_result)
```

	State	Ship Mode	Order Count
0	Alabama	First Class	9.0

1	Arizona	First Class	42.0
2	Arkansas	First Class	10.0
3	California	First Class	302.0
4	Colorado	First Class	43.0
..	...	...	...
191	Virginia	Standard Class	115.0
192	Washington	Standard Class	265.0
193	West Virginia	Standard Class	3.0
194	Wisconsin	Standard Class	66.0
195	Wyoming	Standard Class	1.0

[196 rows x 3 columns]

[Q7] What is the advantage of using melt?

Ans: Simplifying grouping & aggregation process

[Q8] From the superstore\_order, display the ascending order considering values in the 'Profit' column to group the 'Category'.

```
[183]: #enter your code
profit_table = superstore_order.groupby(["Category"]).agg(sum_profit=("Profit",
↪ "sum"))
profit_table.sort_values(by="sum_profit", ascending=True)
```

```
[183]:          sum_profit
Category
Furniture      16858.5619
Office Supplies 105827.0238
Technology      133410.4932
```

[Q9] Create a new column that calculates the total price (sale\*quantity) before discount then group by 'product id' and 'category', then show the mean of the total price

```
[184]: #enter your code here
superstore_order['total_price'] = superstore_order['Quantity'] *
↪ superstore_order['Sales']

superstore_order.groupby(["Category", "Product ID"]).agg(total_price_mean =
↪ ("total_price", "mean"))
```

```
[184]:          total_price_mean
Category Product ID
Furniture FUR-B0-10000112      7426.566000
          FUR-B0-10000330      1258.192000
          FUR-B0-10000362      1726.898000
          FUR-B0-10000468       426.532400
          FUR-B0-10000711      3194.100000
...          ...
```



Technology	TEC-PH-10004912	747.320000
	TEC-PH-10004922	673.249500
	TEC-PH-10004924	57.149333
	TEC-PH-10004959	412.009000
	TEC-PH-10004977	2441.475429

[1846 rows x 1 columns]

[Q10] Complete the function to apply ratio column that calculates from First Class and Standard Class columns on pivot\_table\_result

```
[185]: # [Q10] Complete the function to apply `ratio` column that calculates from
      ↪ `First Class` and `Standard Class` columns on `pivot_table_result`

      # function to transform the ratio
      def get_class_ratio(row):

          # get the first class column
          first_class = row ["First Class"]

          # get the standard class column
          standard_class = row ["Standard Class"]

          # calculate the ratio
          ratio = first_class/standard_class

          return ratio

      pivot_table_result["ratio"] = pivot_table_result.apply(get_class_ratio, axis=1)

      pivot_table_result.head()
```

```
[185]: Ship Mode    First Class    Same Day    Second Class    Standard Class    ratio
State
Alabama           9.0           1.0           18.0           30.0  0.300000
Arizona          42.0          15.0           22.0          123.0  0.341463
Arkansas          10.0           2.0           8.0           35.0  0.285714
California        302.0         106.0          346.0         1000.0  0.302000
Colorado          43.0           5.0           32.0           95.0  0.452632
```

[Q11] After complete Q10, What does the apply function do? Ans: apply a function along a particular axis of a DataFrame

[Q12] Create a new column(short\_ratio) that works the same as Q10 but with lambda function

```
[186]: # [Q12] Create a new column(`short_ratio`) that works the same as Q10 but with
      ↪ `lambda` function
```

```

pivot_table_result["short_ratio"] = pivot_table_result.apply(lambda row:
    ↪row["First Class"] / row["Standard Class"], axis=1)

pivot_table_result.head()

```

```

[186]: Ship Mode    First Class    Same Day    Second Class    Standard Class    ratio \
State
Alabama           9.0           1.0           18.0           30.0    0.300000
Arizona          42.0          15.0           22.0          123.0    0.341463
Arkansas          10.0           2.0           8.0           35.0    0.285714
California        302.0         106.0          346.0         1000.0    0.302000
Colorado          43.0           5.0           32.0           95.0    0.452632

Ship Mode    short_ratio
State
Alabama      0.300000
Arizona      0.341463
Arkansas     0.285714
California   0.302000
Colorado     0.452632

```

[Q13] What is the difference between `apply` and `lambda` function? give 2 examples use case.

Ans: `apply` is method to apply function along particular axis.

`lambda` is work like function with no need to use `def`, small way to create function.

often used together in data manipulation tasks.

Use case 1:

Using `apply` with a `lambda` function to transform a column

```

[187]: temperatures = pd.DataFrame({'Celsius': [10, 20, 30, 40]})
temperatures['Fahrenheit'] = temperatures['Celsius'].apply(lambda x: (x * 9/5)
    ↪+ 32)
print(temperatures)

```

	Celsius	Fahrenheit
0	10	50.0
1	20	68.0
2	30	86.0
3	40	104.0

Use case 2:

Using `apply` with a `lambda` function to create a new column based on conditions

```

[188]: import pandas as pd

```

```
sales = pd.DataFrame({'Sales': [100, 200, 300, 400], 'Profit': [10, -20, 30, -40]})
sales['Performance'] = sales['Profit'].apply(lambda x: 'Good' if x > 0 else 'Bad')
print(sales)
```

	Sales	Profit	Performance
0	100	10	Good
1	200	-20	Bad
2	300	30	Good
3	400	-40	Bad