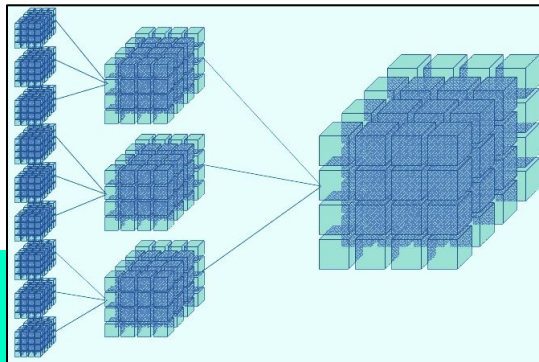


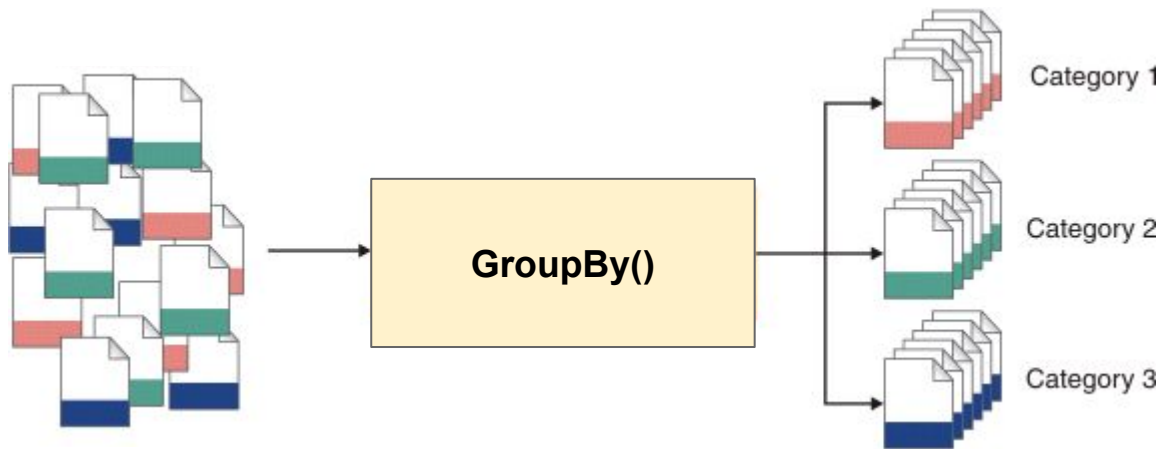
DATA AGGREGATION AND GROUP OPERATIONS.



Chapter 10.

STEP 1: CATEGORIZE DATA

Popular groups include income, age, profession.



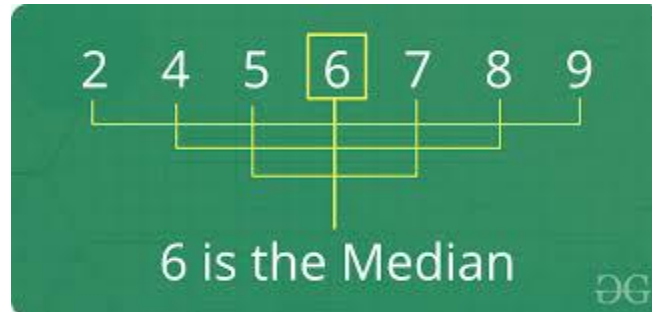
STEP 2: COMPUTE GROUP SUMMARY STATISTICS

count	Number of non-NA values in the group
sum	Sum of non-NA values
mean	Mean of non-NA values
median	Arithmetic median of non-NA values
std, var	Unbiased ($n - 1$ denominator) standard deviation and variance
min, max	Minimum and maximum of non-NA values
prod	Product of non-NA values
first, last	First and last non-NA values

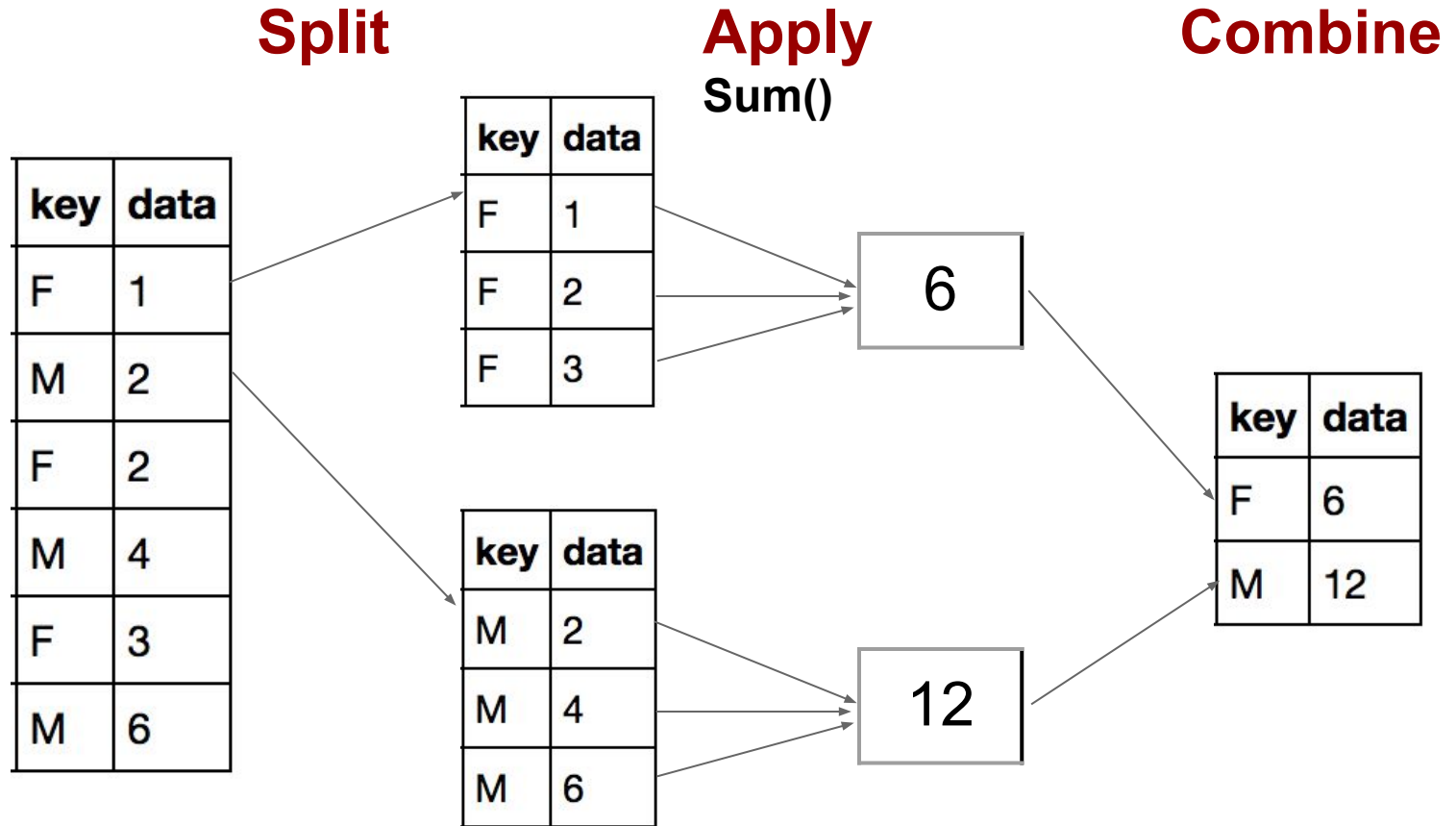
MEAN VS MEDIAN

The mean (average) of a data set is found by adding all numbers in the data set and then dividing by the number of values in the set.

The median is the middle value when a data set is ordered from least to greatest.



SPLIT-APPLY-COMBINE OF A SIMPLE AGGREGATION



CREATE A DATAFRAME

```
df = DataFrame({'key': ['F', 'M', 'F', 'M', 'F', 'M'],  
               'data': [1, 2, 2, 4, 3, 6]},  
               columns=['key', 'data'])  
df
```

	key	data
0	F	1
1	M	2
2	F	2
3	M	4
4	F	3
5	M	6

SPLIT - APPLY - COMBINE

```
grouped = df['data'].groupby(df['key'])  
grouped.sum()
```

key

F 6

M 12

Name: data, dtype: int64

CREATE A DATAFRAME II

```
df = DataFrame({'key1': ['F', 'F', 'F', 'M', 'F', 'M'],  
                'key2': ['20-25', '20-25', '20-25', '25-30', '30-35', '30-35'],  
                'data1': [2, 2, 3, 3, 4, 4]},  
               columns=['key1', 'key2', 'data1'])
```

df

	key1	key2	data1
0	F	20-25	2
1	F	20-25	2
2	F	20-25	3
3	M	25-30	3
4	F	30-35	4
5	M	30-35	4

IF WE GROUP THE DATA BY USING TWO KEYS...

the resulting Series has a **hierarchical index**.

```
grouped = df['data1'].groupby([df['key1'], df['key2']])  
grouped.sum()
```

key1	key2	
F	20-25	7
	30-35	4
M	25-30	3
	30-35	4

Name: data1, dtype: int64

PIVOT TABLE

Month (Multiple Items) 

Sum of Net Sales		Product			
Region	Salesman	FastCar	RapidZoo	SuperGlue	Grand Total
Middle	Joseph	\$ 3,623	\$ 4,782	\$ 7,055	\$ 15,460
	Lawrence	\$ 5,908	\$ 4,642	\$ 4,593	\$ 15,143
	Maria	\$ 6,502	\$ 3,969	\$ 5,408	\$ 15,879
	Matt	\$ 4,170	\$ 6,093	\$ 5,039	\$ 15,302
Middle Total		\$ 20,203	\$ 18,486	\$ 22,095	\$ 60,784
North	Joseph	\$ 3,643	\$ 5,846	\$ 6,574	\$ 16,063
	Lawrence	\$ 4,456	\$ 6,658	\$ 7,685	\$ 18,799
	Maria	\$ 6,235	\$ 4,616	\$ 3,612	\$ 14,463
	Matt	\$ 3,868	\$ 3,526	\$ 3,254	\$ 10,648
North Total		\$ 18,202	\$ 20,646	\$ 20,125	\$ 60,973
West	Joseph	\$ 5,507	\$ 5,186	\$ 4,882	\$ 15,575
	Lawrence	\$ 4,082	\$ 3,272	\$ 6,124	\$ 13,478
	Maria	\$ 5,520	\$ 5,461	\$ 4,872	\$ 15,853
	Matt	\$ 6,737	\$ 4,598	\$ 4,233	\$ 15,568
West Total		\$ 21,846	\$ 20,517	\$ 20,111	\$ 62,474
Grand Total		\$ 60,251	\$ 61,049	\$ 61,331	\$ 184,631

A PIVOT TABLE IS A DATA SUMMARIZATION
TOOL.

IT AGGREGATES A TABLE OF DATA BY ONE OR
MORE KEYS, ARRANGING SOME OF THE GROUP
KEYS ALONG THE ROWS AND SOME ALONG THE
COLUMNS.

PIVOT TABLE EXAMPLE

`sum()` aggregate function

```
df.pivot_table('data1', index='key2', columns='key1', aggfunc='sum')
```

key2	key1	
	F	M
20-25	7.0	NaN
25-30	NaN	3.0
30-35	4.0	4.0