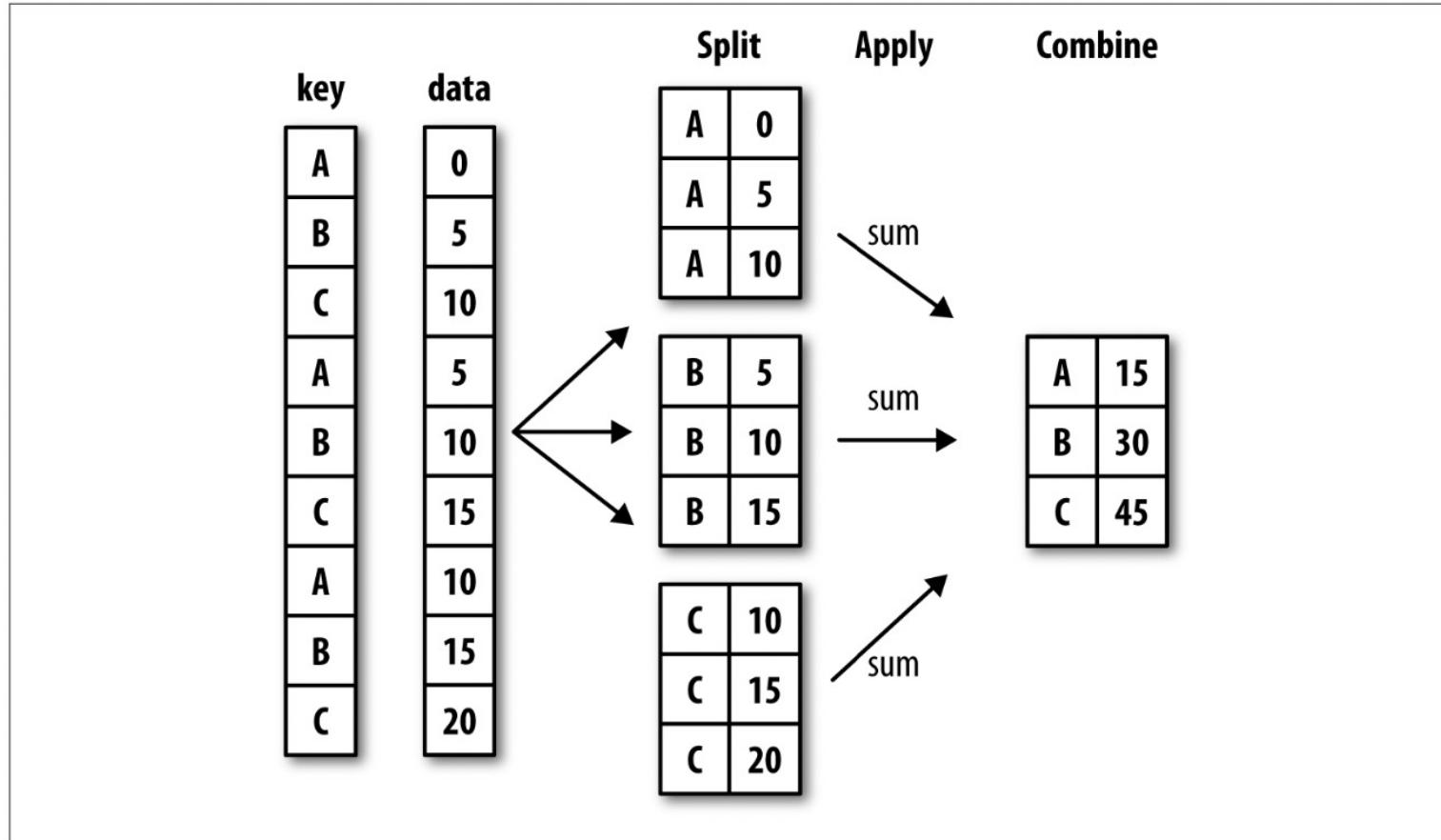


Lecture 9: Data Aggregation and Group Operations

Agenda:

- Split a pandas object into pieces using one or more keys (in the form of functions, arrays, or DataFrame column names)
- Calculate group summary statistics, like count, mean, or standard deviation, or a user-defined function
- Apply within-group transformations or other manipulations, like normalization, linear regression, rank, or subset selection
- Compute pivot tables and cross-tabulations
- Perform quantile analysis and other statistical group analyses

GroupBy Mechanics



- A list or array of values that is the same length as the axis being grouped
- A value indicating a column name in a DataFrame
- A dict or Series giving a correspondence between the values on the axis being grouped and the group names
- A function to be invoked on the axis index or the individual labels in the index

Figure 1. Illustration of a group aggregation

let's get started

```
In [10]: df = pd.DataFrame({'key1' : ['a', 'a', 'b', 'b', 'a'],
....:                         'key2' : ['one', 'two', 'one', 'two', 'one'],
....:                         'data1' : np.random.randn(5),
....:                         'data2' : np.random.randn(5)})
```

```
In [11]: df
Out[11]:
```

	data1	data2	key1	key2
0	-0.204708	1.393406	a	one
1	0.478943	0.092908	a	two
2	-0.519439	0.281746	b	one
3	-0.555730	0.769023	b	two
4	1.965781	1.246435	a	one

To compute group means we can call the GroupBy's mean method:

```
In [14]: grouped.mean()
Out[14]:
key1
a    0.746672
b   -0.537585
Name: data1, dtype: float64
```

Groupby method with the column (a Series) at key1:

```
In [12]: grouped = df['data1'].groupby(df['key1'])
```

```
In [13]: grouped
Out[13]: <pandas.core.groupby.SeriesGroupBy object at 0x7faa31537390>
```

If instead we had passed multiple arrays as a list, we'd get something different:

```
In [15]: means = df['data1'].groupby([df['key1'], df['key2']]).mean()
```

```
In [16]: means
```

```
Out[16]:
```

```
key1  key2
a    one      0.880536
      two      0.478943
b    one     -0.519439
      two     -0.555730
Name: data1, dtype: float64
```

```
In [17]: means.unstack()
```

```
Out[17]:
```

```
key2      one      two
key1
a    0.880536  0.478943
b   -0.519439 -0.555730
```

In this example, the group keys are all Series, though they could be any arrays of the right length:

```
In [18]: states = np.array(['Ohio', 'California', 'California', 'Ohio', 'Ohio'])
```

```
In [19]: years = np.array([2005, 2005, 2006, 2005, 2006])
```

```
In [20]: df['data1'].groupby([states, years]).mean()
```

```
Out[20]:
```

```
California  2005      0.478943
              2006     -0.519439
Ohio        2005     -0.380219
              2006      1.965781
Name: data1, dtype: float64
```

```
In [21]: df.groupby('key1').mean()
```

```
Out[21]:
```

	data1	data2
key1		
a	0.746672	0.910916
b	-0.537585	0.525384

```
In [22]: df.groupby(['key1', 'key2']).mean()
```

```
Out[22]:
```

		data1	data2
key1	key2		
a	one	0.880536	1.319920
	two	0.478943	0.092908
b	one	-0.519439	0.281746
	two	-0.555730	0.769023

```
In [23]: df.groupby(['key1', 'key2']).size()
```

```
Out[23]:
```

key1	key2	
a	one	2
	two	1
b	one	1
	two	1

```
dtype: int64
```

Take note that any missing values in a group key will be excluded from the result.

Iterating Over Groups

```
In [24]: for name, group in df.groupby('key1'):
....:     print(name)
....:     print(group)
....:
```

```
a
   data1      data2 key1 key2
0 -0.204708  1.393406    a  one
1  0.478943  0.092908    a  two
4  1.965781  1.246435    a  one
```

```
b
   data1      data2 key1 key2
2 -0.519439  0.281746    b  one
3 -0.555730  0.769023    b  two
```

In the case of multiple keys, the first element in the tuple will be a tuple of key values:

```
In [25]: for (k1, k2), group in df.groupby(['key1', 'key2']):
....:     print((k1, k2))
....:     print(group)
....:
```

```
('a', 'one')
   data1      data2 key1 key2
0 -0.204708  1.393406    a  one
4  1.965781  1.246435    a  one
```

```
('a', 'two')
   data1      data2 key1 key2
1  0.478943  0.092908    a  two
```

```
('b', 'one')
   data1      data2 key1 key2
2 -0.519439  0.281746    b  one
```

```
('b', 'two')
   data1      data2 key1 key2
3 -0.555730  0.769023    b  two
```

```
In [26]: pieces = dict(list(df.groupby('key1')))
```

We can print out the groups like so:

```
In [27]: pieces['b']
```

```
Out[27]:
```

	data1	data2	key1	key2
2	-0.519439	0.281746	b	one
3	-0.555730	0.769023	b	two

By default groupby groups on axis=0, but you can group on any of the other axes.

```
In [28]: df.dtypes
```

```
Out[28]:
```

data1	float64
data2	float64
key1	object
key2	object
dtype:	object

```
In [29]: grouped = df.groupby(df.dtypes, axis=1)
```

```
In [30]: for dtype, group in grouped:  
.....:     print(dtype)  
.....:     print(group)  
.....:  
float64  
      data1      data2  
0 -0.204708  1.393406  
1  0.478943  0.092908  
2 -0.519439  0.281746  
3 -0.555730  0.769023  
4  1.965781  1.246435  
object  
      key1  key2  
0      a   one  
1      a   two  
2      b   one  
3      b   two  
4      a   one
```

Selecting a Column or Subset of Columns

```
df.groupby('key1')['data1']
df.groupby('key1')[['data2']]
```

```
df['data1'].groupby(df['key1'])
df[['data2']].groupby(df['key1'])
```

```
In [32]: s_grouped = df.groupby(['key1', 'key2'])['data2']
```

```
In [33]: s_grouped
```

```
Out[33]: <pandas.core.groupby.SeriesGroupBy object at 0x7faa30c78da0>
```

```
In [34]: s_grouped.mean()
```

```
Out[34]:
```

```
key1  key2
```

```
a    one    1.319920
     two    0.092908
```

```
b    one    0.281746
     two    0.769023
```

```
Name: data2, dtype: float64
```

```
In [31]: df.groupby(['key1', 'key2'])[['data2']].mean()
```

```
Out[31]:
```

		data2
	key1	key2
a	one	1.319920
	two	0.092908
b	one	0.281746
	two	0.769023

Grouping with Dicts and Series

```
In [35]: people = pd.DataFrame(np.random.randn(5, 5),
.....           columns=['a', 'b', 'c', 'd', 'e'],
.....           index=['Joe', 'Steve', 'Wes', 'Jim', 'Travis'])
```

```
In [36]: people.iloc[2:3, [1, 2]] = np.nan # Add a few NA values
```

```
In [37]: people  
Out[37]:
```

	a	b	c	d	e
Joe	1.007189	-1.296221	0.274992	0.228913	1.352917
Steve	0.886429	-2.001637	-0.371843	1.669025	-0.438570
Wes	-0.539741	NaN	NaN	-1.021228	-0.577087
Jim	0.124121	0.302614	0.523772	0.000940	1.343810
Travis	-0.713544	-0.831154	-2.370232	-1.860761	-0.860757

```
In [38]: mapping = {'a': 'red', 'b': 'red', 'c': 'blue',
.....                 'd': 'blue', 'e': 'red', 'f' : 'orange'}
```

```
In [39]: by_column = people.groupby(mapping, axis=1)
```

```
In [40]: by_column.sum()
```

```
Out[40]:
```

	blue	red
Joe	0.503905	1.063885
Steve	1.297183	-1.553778
Wes	-1.021228	-1.116829
Jim	0.524712	1.770545
Travis	-4.230992	-2.405455

The same functionality holds for Series, which can be viewed as a fixed-size mapping:

```
In [41]: map_series = pd.Series(mapping)
```

```
In [42]: map_series
```

```
Out[42]:
```

```
a      red  
b      red  
c    blue  
d    blue  
e      red  
f  orange
```

```
dtype: object
```

```
In [43]: people.groupby(map_series, axis=1).count()
```

```
Out[43]:
```

	blue	red
Joe	2	3
Steve	2	3
Wes	1	2
Jim	2	3
Travis	2	3

Grouping with Functions

```
In [44]: people.groupby(len).sum()
```

```
Out[44]:
```

	a	b	c	d	e
3	0.591569	-0.993608	0.798764	-0.791374	2.119639
5	0.886429	-2.001637	-0.371843	1.669025	-0.438570
6	-0.713544	-0.831154	-2.370232	-1.860761	-0.860757

Mixing functions with arrays, dicts, or Series is not a problem as everything gets converted to arrays internally:

```
In [45]: key_list = ['one', 'one', 'one', 'two', 'two']
```

```
In [46]: people.groupby([len, key_list]).min()
```

```
Out[46]:
```

	a	b	c	d	e
3 one	-0.539741	-1.296221	0.274992	-1.021228	-0.577087
two	0.124121	0.302614	0.523772	0.000940	1.343810
5 one	0.886429	-2.001637	-0.371843	1.669025	-0.438570
6 two	-0.713544	-0.831154	-2.370232	-1.860761	-0.860757

Grouping by Index Levels

```
In [47]: columns = pd.MultiIndex.from_arrays([[['US', 'US', 'US', 'JP', 'JP'],
.....:                               [1, 3, 5, 1, 3]],
.....:                               names=['cty', 'tenor'])]
```

```
In [48]: hier_df = pd.DataFrame(np.random.randn(4, 5), columns=columns)
```

```
In [49]: hier_df
```

```
Out[49]:
```

	cty	US		JP		
tenor		1	3	5	1	3
0	0.560145	-1.265934	0.119827	-1.063512	0.332883	
1	-2.359419	-0.199543	-1.541996	-0.970736	-1.307030	
2	0.286350	0.377984	-0.753887	0.331286	1.349742	
3	0.069877	0.246674	-0.011862	1.004812	1.327195	

```
In [50]: hier_df.groupby(level='cty', axis=1).count()
```

```
Out[50]:
```

	cty	JP	US
0	2	3	
1	2	3	
2	2	3	
3	2	3	

Data Aggregation

Function name	Description
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

*Optimized **groupby** methods*

```
In [51]: df
Out[51]:
   data1      data2  key1  key2
0 -0.204708  1.393406    a  one
1  0.478943  0.092908    a  two
2 -0.519439  0.281746    b  one
3 -0.555730  0.769023    b  two
4  1.965781  1.246435    a  one
```

```
In [52]: grouped = df.groupby('key1')
```

```
In [54]: def peak_to_peak(arr):
....:     return arr.max() - arr.min()
```

```
In [55]: grouped.agg(peak_to_peak)
```

```
Out[55]:
   data1      data2
key1
a  2.170488  1.300498
b  0.036292  0.487276
```

```
In [53]: grouped['data1'].quantile(0.9)
Out[53]:
key1
a    1.668413
b   -0.523068
Name: data1, dtype: float64
```

```
In [56]: grouped.describe()
```

```
Out[56]:
   data1
   count      mean       std      min      25%      50%      75%
key1
a    3.0  0.746672  1.109736 -0.204708  0.137118  0.478943  1.222362
b    2.0 -0.537585  0.025662 -0.555730 -0.546657 -0.537585 -0.528512
```

```
   data2
   max  count      mean       std      min      25%      50%
key1
a  1.965781  3.0  0.910916  0.712217  0.092908  0.669671  1.246435
b -0.519439  2.0  0.525384  0.344556  0.281746  0.403565  0.525384
```

```
   75%      max
key1
a  1.319920  1.393406
b  0.647203  0.769023
```

Column-Wise and Multiple Function Application

```
In [57]: tips = pd.read_csv('examples/tips.csv')
```

Add tip percentage of total bill

```
In [58]: tips['tip_pct'] = tips['tip'] / tips['total_bill']
```

```
In [59]: tips[:6]
```

```
Out[59]:
```

	total_bill	tip	smoker	day	time	size	tip_pct
0	16.99	1.01	No	Sun	Dinner	2	0.059447
1	10.34	1.66	No	Sun	Dinner	3	0.160542
2	21.01	3.50	No	Sun	Dinner	3	0.166587
3	23.68	3.31	No	Sun	Dinner	2	0.139780
4	24.59	3.61	No	Sun	Dinner	4	0.146808
5	25.29	4.71	No	Sun	Dinner	4	0.186240

```
In [60]: grouped = tips.groupby(['day', 'smoker'])
```

```
In [61]: grouped_pct = grouped['tip_pct']
```

```
In [62]: grouped_pct.agg('mean')
```

```
Out[62]:
```

day	smoker	tip_pct
Fri	No	0.151650
	Yes	0.174783
Sat	No	0.158048
	Yes	0.147906
Sun	No	0.160113
	Yes	0.187250
Thur	No	0.160298
	Yes	0.163863

```
Name: tip_pct, dtype: float64
```

```
In [63]: grouped_pct.agg(['mean', 'std', peak_to_peak])  
Out[63]:
```

		mean	std	peak_to_peak
day	smoker			
Fri	No	0.151650	0.028123	0.067349
	Yes	0.174783	0.051293	0.159925
Sat	No	0.158048	0.039767	0.235193
	Yes	0.147906	0.061375	0.290095
Sun	No	0.160113	0.042347	0.193226
	Yes	0.187250	0.154134	0.644685
Thur	No	0.160298	0.038774	0.193350
	Yes	0.163863	0.039389	0.151240

```
In [64]: grouped_pct.agg([('foo', 'mean'), ('bar', np.std)])  
Out[64]:
```

		foo	bar
day	smoker		
Fri	No	0.151650	0.028123
	Yes	0.174783	0.051293
Sat	No	0.158048	0.039767
	Yes	0.147906	0.061375
Sun	No	0.160113	0.042347
	Yes	0.187250	0.154134
Thur	No	0.160298	0.038774
	Yes	0.163863	0.039389

```
In [65]: functions = ['count', 'mean', 'max']
```

```
In [66]: result = grouped[['tip_pct', 'total_bill']].agg(functions)
```

```
In [67]: result
```

```
Out[67]:
```

		tip_pct		total_bill			
		count	mean	max	count	mean	max
day	smoker						
Fri	No	4	0.151650	0.187735	4	18.420000	22.75
	Yes	15	0.174783	0.263480	15	16.813333	40.17
Sat	No	45	0.158048	0.291990	45	19.661778	48.33
	Yes	42	0.147906	0.325733	42	21.276667	50.81
Sun	No	57	0.160113	0.252672	57	20.506667	48.17
	Yes	19	0.187250	0.710345	19	24.120000	45.35
Thur	No	45	0.160298	0.266312	45	17.113111	41.19
	Yes	17	0.163863	0.241255	17	19.190588	43.11

```
In [68]: result['tip_pct']
```

```
Out[68]:
```

		count	mean	max
day	smoker			
Fri	No	4	0.151650	0.187735
	Yes	15	0.174783	0.263480
Sat	No	45	0.158048	0.291990
	Yes	42	0.147906	0.325733
Sun	No	57	0.160113	0.252672
	Yes	19	0.187250	0.710345
Thur	No	45	0.160298	0.266312
	Yes	17	0.163863	0.241255

```
In [71]: grouped.agg({'tip' : np.max, 'size' : 'sum'})
```

```
Out[71]:
```

		tip	size
day	smoker		
Fri	No	3.50	9
	Yes	4.73	31
Sat	No	9.00	115
	Yes	10.00	104
Sun	No	6.00	167
	Yes	6.50	49
Thur	No	6.70	112
	Yes	5.00	40

```
In [69]: ftuples = [('Durchschnitt', 'mean'), ('Abweichung', np.var)]
```

```
In [70]: grouped['tip_pct', 'total_bill'].agg(ftuples)
```

```
Out[70]:
```

		tip_pct	total_bill		
		Durchschnitt	Abweichung	Durchschnitt	Abweichung
day	smoker				
Fri	No	0.151650	0.000791	18.420000	25.596333
	Yes	0.174783	0.002631	16.813333	82.562438
Sat	No	0.158048	0.001581	19.661778	79.908965
	Yes	0.147906	0.003767	21.276667	101.387535
Sun	No	0.160113	0.001793	20.506667	66.099980
	Yes	0.187250	0.023757	24.120000	109.046044
Thur	No	0.160298	0.001503	17.113111	59.625081
	Yes	0.163863	0.001551	19.190588	69.808518

```
In [72]: grouped.agg({'tip_pct' : ['min', 'max', 'mean', 'std'],
```

```
.....: 'size' : 'sum'})
```

```
Out[72]:
```

		tip_pct		size	
		min	max	mean	std
day	smoker				
Fri	No	0.120385	0.187735	0.151650	0.028123
	Yes	0.103555	0.263480	0.174783	0.051293
Sat	No	0.056797	0.291990	0.158048	0.039767
	Yes	0.035638	0.325733	0.147906	0.061375
Sun	No	0.059447	0.252672	0.160113	0.042347
	Yes	0.065660	0.710345	0.187250	0.154134
Thur	No	0.072961	0.266312	0.160298	0.038774
	Yes	0.090014	0.241255	0.163863	0.039389

Returning Aggregated Data Without Row Indexes

```
In [73]: tips.groupby(['day', 'smoker'], as_index=False).mean()
```

```
Out[73]:
```

	day	smoker	total_bill	tip	size	tip_pct
0	Fri	No	18.420000	2.812500	2.250000	0.151650
1	Fri	Yes	16.813333	2.714000	2.066667	0.174783
2	Sat	No	19.661778	3.102889	2.555556	0.158048
3	Sat	Yes	21.276667	2.875476	2.476190	0.147906
4	Sun	No	20.506667	3.167895	2.929825	0.160113
5	Sun	Yes	24.120000	3.516842	2.578947	0.187250
6	Thur	No	17.113111	2.673778	2.488889	0.160298
7	Thur	Yes	19.190588	3.030000	2.352941	0.163863

Apply: General split-apply-combine

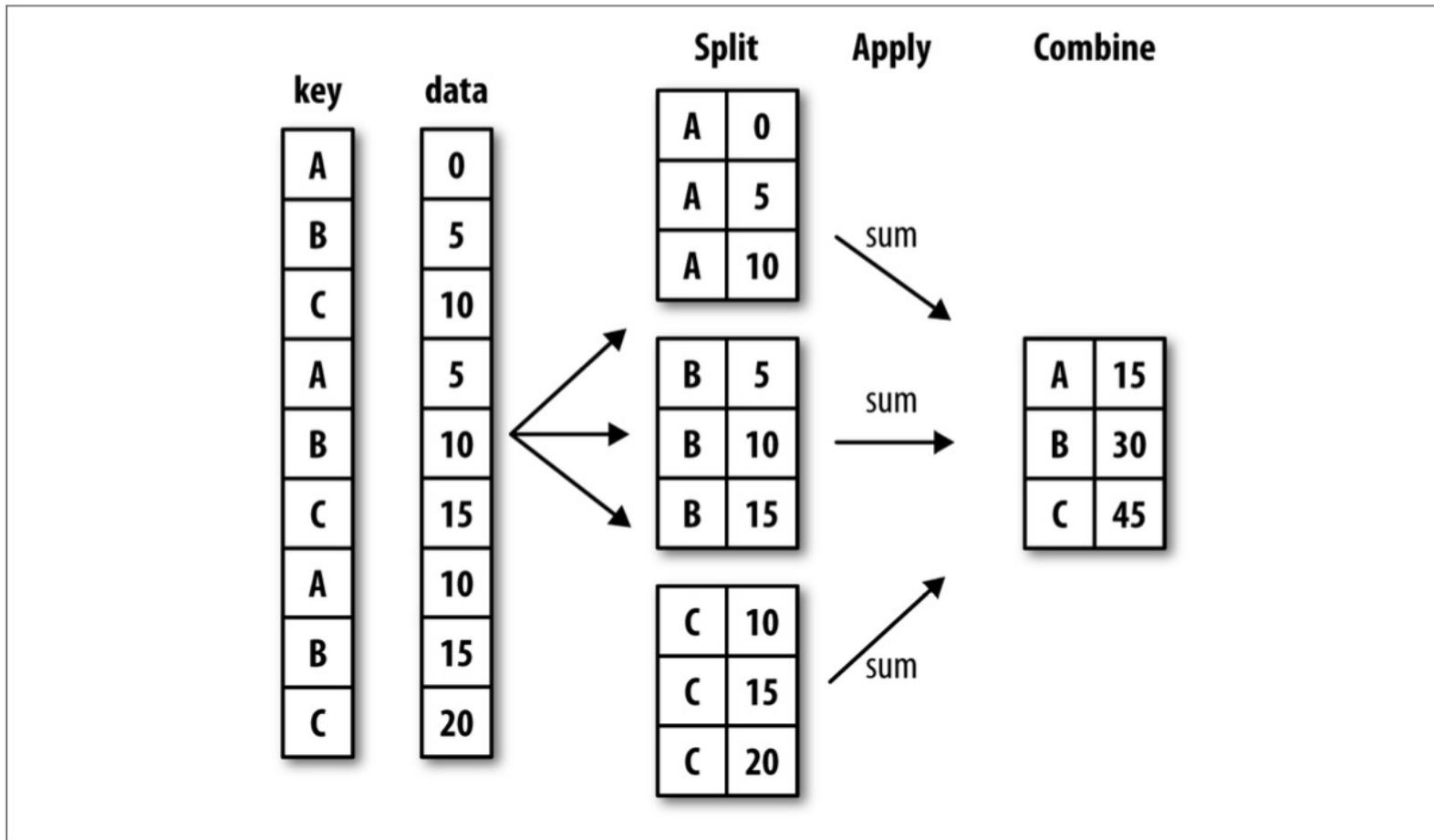


Figure 2. Illustration of a group aggregation

```
In [74]: def top(df, n=5, column='tip_pct'):
....:     return df.sort_values(by=column)[-n:]
```

```
In [75]: top(tips, n=6)
```

```
Out[75]:
   total_bill  tip  smoker  day    time  size  tip_pct
109      14.31  4.00     Yes  Sat  Dinner    2  0.279525
183      23.17  6.50     Yes  Sun  Dinner    4  0.280535
232      11.61  3.39     No   Sat  Dinner    2  0.291990
67       3.07  1.00     Yes  Sat  Dinner    1  0.325733
178      9.60  4.00     Yes  Sun  Dinner    2  0.416667
172      7.25  5.15     Yes  Sun  Dinner    2  0.710345
```

```
In [76]: tips.groupby('smoker').apply(top)
```

```
Out[76]:
```

smoker		total_bill	tip	smoker	day	time	size	tip_pct
No	88	24.71	5.85	No	Thur	Lunch	2	0.236746
	185	20.69	5.00	No	Sun	Dinner	5	0.241663
	51	10.29	2.60	No	Sun	Dinner	2	0.252672
	149	7.51	2.00	No	Thur	Lunch	2	0.266312
	232	11.61	3.39	No	Sat	Dinner	2	0.291990
Yes	109	14.31	4.00	Yes	Sat	Dinner	2	0.279525
	183	23.17	6.50	Yes	Sun	Dinner	4	0.280535
	67	3.07	1.00	Yes	Sat	Dinner	1	0.325733
	178	9.60	4.00	Yes	Sun	Dinner	2	0.416667
	172	7.25	5.15	Yes	Sun	Dinner	2	0.710345

```
In [77]: tips.groupby(['smoker', 'day']).apply(top, n=1, column='total_bill')
```

```
Out[77]:
```

		total_bill	tip	smoker	day	time	size	tip_pct
smoker	day							
No	Fri	94	22.75	3.25	No	Fri	Dinner	2 0.142857
	Sat	212	48.33	9.00	No	Sat	Dinner	4 0.186220
	Sun	156	48.17	5.00	No	Sun	Dinner	6 0.103799
	Thur	142	41.19	5.00	No	Thur	Lunch	5 0.121389
Yes	Fri	95	40.17	4.73	Yes	Fri	Dinner	4 0.117750
	Sat	170	50.81	10.00	Yes	Sat	Dinner	3 0.196812
	Sun	182	45.35	3.50	Yes	Sun	Dinner	3 0.077178
	Thur	197	43.11	5.00	Yes	Thur	Lunch	4 0.115982

```
In [78]: result = tips.groupby('smoker')['tip_pct'].describe()
```

```
In [79]: result
```

```
Out[79]:
```

```
count      mean       std      min      25%      50%      75%  \
smoker
No        151.0   0.159328  0.039910  0.056797  0.136906  0.155625  0.185014
Yes       93.0    0.163196  0.085119  0.035638  0.106771  0.153846  0.195059
max
```

```
smoker
No     0.291990
Yes    0.710345
```

```
f = lambda x: x.describe()
grouped.apply(f)
```

```
In [80]: result.unstack('smoker')
```

```
Out[80]:
```

```
smoker
count      No      151.000000
               Yes     93.000000
mean      No      0.159328
               Yes     0.163196
std       No      0.039910
               Yes     0.085119
min      No      0.056797
               Yes     0.035638
25%      No      0.136906
               Yes     0.106771
50%      No      0.155625
               Yes     0.153846
75%      No      0.185014
               Yes     0.195059
max      No      0.291990
               Yes     0.710345
dtype: float64
```

Suppressing the Group Keys

```
In [81]: tips.groupby('smoker', group_keys=False).apply(top)
```

```
Out[81]:
```

	total_bill	tip	smoker	day	time	size	tip_pct
88	24.71	5.85	No	Thur	Lunch	2	0.236746
185	20.69	5.00	No	Sun	Dinner	5	0.241663
51	10.29	2.60	No	Sun	Dinner	2	0.252672
149	7.51	2.00	No	Thur	Lunch	2	0.266312
232	11.61	3.39	No	Sat	Dinner	2	0.291990
109	14.31	4.00	Yes	Sat	Dinner	2	0.279525
183	23.17	6.50	Yes	Sun	Dinner	4	0.280535
67	3.07	1.00	Yes	Sat	Dinner	1	0.325733
178	9.60	4.00	Yes	Sun	Dinner	2	0.416667
172	7.25	5.15	Yes	Sun	Dinner	2	0.710345

Quantile and Bucket Analysis

```
In [82]: frame = pd.DataFrame({'data1': np.random.randn(1000),
.....: 'data2': np.random.randn(1000)})
```

```
In [83]: quartiles = pd.cut(frame.data1, 4)
```

```
In [84]: quartiles[:10]
```

```
Out[84]:
```

```
0    (-1.23, 0.489]
1    (-2.956, -1.23]
2    (-1.23, 0.489]
3    (0.489, 2.208]
4    (-1.23, 0.489]
5    (0.489, 2.208]
6    (-1.23, 0.489]
7    (-1.23, 0.489]
8    (0.489, 2.208]
9    (0.489, 2.208]
```

```
Name: data1, dtype: category
```

```
Categories (4, interval[float64]): [(-2.956, -1.23] < (-1.23, 0.489] < (0.489, 2.
208] < (2.208, 3.928)]
```

```
In [85]: def get_stats(group):
```

```
....:     return {'min': group.min(), 'max': group.max(),
....:             'count': group.count(), 'mean': group.mean()}
```

```
In [86]: grouped = frame.data2.groupby(quartiles)
```

```
In [87]: grouped.apply(get_stats).unstack()
```

```
Out[87]:
```

		count	max	mean	min
data1					
	(-2.956, -1.23]	95.0	1.670835	-0.039521	-3.399312
	(-1.23, 0.489]	598.0	3.260383	-0.002051	-2.989741
	(0.489, 2.208]	297.0	2.954439	0.081822	-3.745356
	(2.208, 3.928]	10.0	1.765640	0.024750	-1.929776

```
# Return quantile numbers
```

```
In [88]: grouping = pd.qcut(frame.data1, 10, labels=False)
```

```
In [89]: grouped = frame.data2.groupby(grouping)
```

```
In [90]: grouped.apply(get_stats).unstack()
```

```
Out[90]:
```

	count	max	mean	min
data1				
0	100.0	1.670835	-0.049902	-3.399312
1	100.0	2.628441	0.030989	-1.950098
2	100.0	2.527939	-0.067179	-2.925113
3	100.0	3.260383	0.065713	-2.315555
4	100.0	2.074345	-0.111653	-2.047939
5	100.0	2.184810	0.052130	-2.989741
6	100.0	2.458842	-0.021489	-2.223506
7	100.0	2.954439	-0.026459	-3.056990
8	100.0	2.735527	0.103406	-3.745356
9	100.0	2.377020	0.220122	-2.064111

Pivot Tables and Cross-Tabulation

Suppose you wanted to compute a table of group means (the default pivot_table aggregation type) arranged by day and smoker on the rows:

```
In [130]: tips.pivot_table(index=['day', 'smoker'])
```

```
Out[130]:
```

		size	tip	tip_pct	total_bill
day	smoker				
Fri	No	2.250000	2.812500	0.151650	18.420000
	Yes	2.066667	2.714000	0.174783	16.813333
Sat	No	2.555556	3.102889	0.158048	19.661778
	Yes	2.476190	2.875476	0.147906	21.276667
Sun	No	2.929825	3.167895	0.160113	20.506667
	Yes	2.578947	3.516842	0.187250	24.120000
Thur	No	2.488889	2.673778	0.160298	17.113111
	Yes	2.352941	3.030000	0.163863	19.190588

```
In [131]: tips.pivot_table(['tip_pct', 'size'], index=['time', 'day'], columns='smoker')
```

```
Out[131]:
```

		size		tip_pct	
smoker		No	Yes	No	Yes
Dinner	time	day			
	Fri				
Lunch	Dinner	Fri	2.000000	2.222222	0.139622
		Sat	2.555556	2.476190	0.158048
	Lunch	Sun	2.929825	2.578947	0.160113
		Thur	2.000000	NaN	0.159744
Dinner	Lunch	Fri	3.000000	1.833333	0.187735
		Thur	2.500000	2.352941	0.160311
Dinner	Dinner				0.188937
					0.163863

```
In [132]: tips.pivot_table(['tip_pct', 'size'], index=['time', 'day'],
....:                           columns='smoker', margins=True)
```

```
Out[132]:
```

smoker		size			tip_pct		
		No	Yes	All	No	Yes	All
time	day						
Dinner	Fri	2.000000	2.222222	2.166667	0.139622	0.165347	0.158916
	Sat	2.555556	2.476190	2.517241	0.158048	0.147906	0.153152
	Sun	2.929825	2.578947	2.842105	0.160113	0.187250	0.166897
	Thur	2.000000	NaN	2.000000	0.159744	NaN	0.159744
Lunch	Fri	3.000000	1.833333	2.000000	0.187735	0.188937	0.188765
	Thur	2.500000	2.352941	2.459016	0.160311	0.163863	0.161301
All		2.668874	2.408602	2.569672	0.159328	0.163196	0.160803

To use a different aggregation function,
pass it to aggfunc.

```
In [133]: tips.pivot_table('tip_pct', index=['time', 'smoker'], columns='day',
....:                           aggfunc=len, margins=True)
```

```
Out[133]:
```

		day	Fri	Sat	Sun	Thur	All
time	smoker						
Dinner	No		3.0	45.0	57.0	1.0	106.0
	Yes		9.0	42.0	19.0	NaN	70.0
Lunch	No		1.0	NaN	NaN	44.0	45.0
	Yes		6.0	NaN	NaN	17.0	23.0
All			19.0	87.0	76.0	62.0	244.0

```
In [134]: tips.pivot_table('tip_pct', index=['time', 'size', 'smoker'],
.....:                               columns='day', aggfunc='mean', fill_value=0)
```

```
Out[134]:
```

day		Fri	Sat	Sun	Thur
time	size	smoker			
Dinner	1	No	0.000000	0.137931	0.000000
		Yes	0.000000	0.325733	0.000000
	2	No	0.139622	0.162705	0.168859
		Yes	0.171297	0.148668	0.207893
	3	No	0.000000	0.154661	0.152663
		Yes	0.000000	0.144995	0.152660
	4	No	0.000000	0.150096	0.148143
		Yes	0.117750	0.124515	0.193370
	5	No	0.000000	0.000000	0.206928
		Yes	0.000000	0.106572	0.065660
...					
Lunch	1	No	0.000000	0.000000	0.000000
		Yes	0.223776	0.000000	0.000000
	2	No	0.000000	0.000000	0.000000
		Yes	0.181969	0.000000	0.000000
	3	No	0.187735	0.000000	0.000000
		Yes	0.000000	0.000000	0.084246
	4	No	0.000000	0.000000	0.000000
		Yes	0.000000	0.000000	0.204952
	5	No	0.000000	0.000000	0.000000
		Yes	0.000000	0.000000	0.138919
	6	No	0.000000	0.000000	0.000000
		Yes	0.000000	0.000000	0.155410

```
[21 rows x 4 columns]
```

Table 2 for a summary of pivot_table methods.

Function name	Description
values	Column name or names to aggregate; by default aggregates all numeric columns
index	Column names or other group keys to group on the rows of the resulting pivot table
columns	Column names or other group keys to group on the columns of the resulting pivot table
aggfunc	Aggregation function or list of functions ('mean' by default); can be any function valid in a groupby context
fill_value	Replace missing values in result table
dropna	If True, do not include columns whose entries are all NA
margins	Add row/column subtotals and grand total (False by default)

Cross-Tabulations: Crosstab

A cross-tabulation (or *crosstab* for short) is a special case of a pivot table that computes group frequencies.
Here is an example:

```
In [138]: data
Out[138]:
   Sample Nationality   Handedness
0      1        USA  Right-handed
1      2       Japan  Left-handed
2      3        USA  Right-handed
3      4       Japan  Right-handed
4      5       Japan  Left-handed
5      6       Japan  Right-handed
6      7        USA  Right-handed
7      8       USA  Left-handed
8      9       Japan  Right-handed
9     10        USA  Right-handed
```

```
In [139]: pd.crosstab(data.Nationality, data.Handedness, margins=True)
Out[139]:
   Handedness  Left-handed  Right-handed  All
   Nationality
   Japan           2           3          5
   USA             1           4          5
   All             3           7         10

In [140]: pd.crosstab([tips.time, tips.day], tips.smoker, margins=True)
Out[140]:
   smoker      No  Yes  All
   time    day
   Dinner  Fri    3    9   12
                  Sat   45   42   87
                  Sun   57   19   76
                  Thur   1    0    1
   Lunch   Fri    1    6    7
                  Thur   44   17   61
   All        151  93  244
```