

Pandas Basics

for Data Analysis

By Armin Khayati

Pandas Basics

- Pandas contains data structures and data manipulation tools designed to make data cleaning and analysis fast and easy in Python.
- `pandas` is often used in tandem with numerical computing tools like NumPy and SciPy, analytical libraries like `statsmodels` and `scikit-learn`, and data visualization libraries like `matplotlib`.
- While `pandas` adopts many coding idioms from NumPy, the biggest difference is that `pandas` is designed for working with **tabular** or **heterogeneous** data.
- NumPy, by contrast, is best suited for working with **homogeneous** numerical array data.

Pandas Basics

- Throughout the rest of the slides, we use the following import convention for pandas:

```
In [1]: import pandas as pd
```

- Thus, whenever you see `pd.` in code, it's referring to pandas.
- You may also find it easier to import `Series` and `DataFrame` into the local namespace since they are so frequently used

```
In [2]: from pandas import Series, DataFrame
```

Pandas Basics

- To get started with pandas, you will need to get comfortable with its two workhorse data structures: `Series` and `DataFrame`.
- While they are not a universal solution for every problem, they provide a solid, easy-to-use basis for most applications.

Series

- A `Series` is a one-dimensional array-like object containing a sequence of values and an associated array of data labels, called its `index`.
- The simplest `Series` is formed from only an array of data:

```
In [11]: obj = pd.Series([4, 7, -5, 3])
```

```
In [12]: obj
```

```
Out[12]:
```

```
0      4
```

```
1      7
```

```
2     -5
```

```
3      3
```

```
dtype: int64
```

Series

- The string representation of a Series displayed interactively shows the index on the left and the values on the right.
- Since we did not specify an index for the data, a default one consisting of the integers 0 through $N - 1$ (where N is the length of the data) is created.

```
In [13]: obj.values
```

```
Out[13]: array([ 4,  7, -5,  3])
```

```
In [14]: obj.index # like range(4)
```

```
Out[14]: RangeIndex(start=0, stop=4, step=1)
```

Series

- Often it will be desirable to create a Series with an index identifying each data point with a label:

```
In [15]: obj2 = pd.Series([4, 7, -5, 3], index=['d', 'b', 'a', 'c'])
```

```
In [16]: obj2
```

```
Out[16]:
```

```
d      4
```

```
b      7
```

```
a     -5
```

```
c      3
```

```
dtype: int64
```

```
In [17]: obj2.index
```

```
Out[17]: Index(['d', 'b', 'a', 'c'], dtype='object')
```

Series

- Compared with NumPy arrays, you can use labels in the index when selecting single values or a set of values:

```
In [18]: obj2['a']
```

```
Out[18]: -5
```

```
In [19]: obj2['d'] = 6
```

```
In [20]: obj2[['c', 'a', 'd']]
```

```
Out[20]:
```

```
c      3
```

```
a     -5
```

```
d      6
```

```
dtype: int64
```

- Here `['c', 'a', 'd']` is interpreted as a list of indices, even though it contains strings instead of integers.

Series

- Using NumPy functions or NumPy-like operations, such as filtering with a boolean array, scalar multiplication, or applying math functions, will preserve the index-value link:

```
In [21]: obj2[obj2 > 0]
```

```
Out[21]:
```

```
d    6
b    7
c    3
dtype: int64
```

```
In [22]: obj2 * 2
```

```
Out[22]:
```

```
d    12
b    14
a   -10
c     6
dtype: int64
```

```
In [23]: np.exp(obj2)
```

```
Out[23]:
```

```
d    403.428793
b   1096.633158
a     0.006738
c    20.085537
dtype: float64
```

Series

- Another way to think about a Series is as a fixed-length, ordered dict, as it is a mapping of index values to data values. It can be used in many contexts where you might use a dict:

```
In [24]: 'b' in obj2
```

```
Out[24]: True
```

```
In [25]: 'e' in obj2
```

```
Out[25]: False
```

Series

- You can create a Series from a Python dictionary by passing the `dict`:

```
In [26]: sdata = {'Ohio': 35000, 'Texas': 71000, 'Oregon': 16000, 'Utah': 5000}
```

```
In [27]: obj3 = pd.Series(sdata)
```

```
In [28]: obj3
```

```
Out[28]:
```

```
Ohio      35000
```

```
Oregon     16000
```

```
Texas      71000
```

```
Utah        5000
```

```
dtype: int64
```

Series

- When you are only passing a dict, the index in the resulting Series will have the dict's keys in sorted order.
- You can override this by passing the dict keys in the order you want them to appear in the resulting Series:

```
In [29]: states = ['California', 'Ohio', 'Oregon', 'Texas']
```

```
In [30]: obj4 = pd.Series(sdata, index=states)
```

```
In [31]: obj4
```

```
Out[31]:
```

California	NaN
Ohio	35000.0
Oregon	16000.0
Texas	71000.0

dtype: float64

Series

- Here, three values found in `sdata` were placed in the appropriate locations, but since no value for 'California' was found, it appears as `NaN` (not a number), which is considered in pandas to mark missing or *NA* values.
- Since 'Utah' was not included in `states` , it is excluded from the resulting object.

Series

- The `isnull` and `notnull` functions in pandas should be used to detect missing data:

```
In [32]: pd.isnull(obj4)
```

```
Out[32]:
```

California	True
Ohio	False
Oregon	False
Texas	False

```
dtype: bool
```

```
In [33]: pd.notnull(obj4)
```

```
Out[33]:
```

California	False
Ohio	True
Oregon	True
Texas	True

```
dtype: bool
```

Series

- Series also has these as instance methods:

```
In [34]: obj4.isnull()
Out[34]:
California    True
Ohio          False
Oregon        False
Texas         False
dtype: bool
```

Series

- A useful Series feature for many applications is that it automatically aligns by index label in arithmetic operations:

```
In [35]: obj3
Out[35]:
Ohio      35000
Oregon    16000
Texas     71000
Utah       5000
dtype: int64
```

```
In [36]: obj4
Out[36]:
California      NaN
Ohio           35000.0
Oregon          16000.0
Texas           71000.0
dtype: float64
```

```
In [37]: obj3 + obj4
Out[37]:
California      NaN
Ohio           70000.0
Oregon          32000.0
Texas          142000.0
Utah            NaN
dtype: float64
```


Series

- If you have experience with databases, you can think about Data alignment as being similar to a join operation.
- Both the Series object itself and its index have a name attribute, which integrates with other key areas of pandas functionality:

```
In [38]: obj4.name = 'population'
```

```
In [39]: obj4.index.name = 'state'
```

```
In [40]: obj4
```

```
Out[40]:
```

```
state
```

```
California      NaN
```

```
Ohio            35000.0
```

```
Oregon          16000.0
```

```
Texas           71000.0
```

```
Name: population, dtype: float64
```

Series

- A Series's index can be altered in-place by assignment:

```
In [41]: obj
```

```
Out[41]:
```

```
0    4
```

```
1    7
```

```
2   -5
```

```
3    3
```

```
dtype: int64
```

```
In [42]: obj.index = ['Bob', 'Steve', 'Jeff', 'Ryan']
```

```
In [43]: obj
```

```
Out[43]:
```

```
Bob    4
```

```
Steve  7
```

```
Jeff  -5
```

```
Ryan   3
```

```
dtype: int64
```

DataFrame

- A DataFrame represents a rectangular table of data and contains an ordered collection of columns, each of which can be a different value type (numeric, string, boolean, etc.).
- The DataFrame has both a row and column index; it can be thought of as a dict of Series all sharing the same index.
- Under the hood, the data is stored as one or more two-dimensional blocks rather than a list, dict, or some other collection of one-dimensional arrays.

DataFrame

- There are many ways to construct a DataFrame, though one of the most common is from a dict of equal-length lists or NumPy arrays:

```
data = {'state': ['Ohio', 'Ohio', 'Ohio', 'Nevada', 'Nevada', 'Nevada'],  
        'year': [2000, 2001, 2002, 2001, 2002, 2003],  
        'pop': [1.5, 1.7, 3.6, 2.4, 2.9, 3.2]}  
frame = pd.DataFrame(data)
```

- The resulting DataFrame will have its index assigned automatically as with Series, and the columns are placed in sorted order:

```
In [45]: frame
```

```
Out[45]:
```

	pop	state	year
0	1.5	Ohio	2000
1	1.7	Ohio	2001
2	3.6	Ohio	2002
3	2.4	Nevada	2001
4	2.9	Nevada	2002
5	3.2	Nevada	2003

- .

DataFrame

- For large DataFrames, the `head` method selects only the first five rows:

```
In [46]: frame.head()
```

```
Out[46]:
```

	pop	state	year
0	1.5	Ohio	2000
1	1.7	Ohio	2001
2	3.6	Ohio	2002
3	2.4	Nevada	2001
4	2.9	Nevada	2002

- If you specify a sequence of columns, the DataFrame's columns will be arranged in that order:

```
In [47]: pd.DataFrame(data, columns=['year', 'state', 'pop'])
```

```
Out[47]:
```

	year	state	pop
0	2000	Ohio	1.5
1	2001	Ohio	1.7
2	2002	Ohio	3.6
3	2001	Nevada	2.4
4	2002	Nevada	2.9
5	2003	Nevada	3.2

DataFrame

- If you pass a column that isn't contained in the dict, it will appear with missing values in the result:

```
In [48]: frame2 = pd.DataFrame(data, columns=['year', 'state', 'pop', 'debt'],  
.....:                           index=['one', 'two', 'three', 'four',  
.....:                           'five', 'six'])
```

```
In [49]: frame2
```

```
Out[49]:
```

	year	state	pop	debt
one	2000	Ohio	1.5	NaN
two	2001	Ohio	1.7	NaN
three	2002	Ohio	3.6	NaN
four	2001	Nevada	2.4	NaN
five	2002	Nevada	2.9	NaN
six	2003	Nevada	3.2	NaN

```
In [50]: frame2.columns
```

```
Out[50]: Index(['year', 'state', 'pop', 'debt'], dtype='object')
```

DataFrame

- A column in a DataFrame can be retrieved as a Series either by dict-like notation or by attribute:

```
In [51]: frame2['state']  
Out[51]:  
one      Ohio  
two      Ohio  
three    Ohio  
four     Nevada  
five     Nevada  
six      Nevada  
Name: state, dtype: object
```

```
In [52]: frame2.year  
Out[52]:  
one      2000  
two      2001  
three    2002  
four     2001  
five     2002  
six      2003  
Name: year, dtype: int64
```

DataFrame

- Note that the returned Series have the same index as the DataFrame, and their `name` attribute has been appropriately set.
- Rows can also be retrieved by position or name with the special `loc` attribute (much more on this later).

```
In [53]: frame2.loc['three']
```

```
Out[53]:
```

```
year      2002
```

```
state     Ohio
```

```
pop       3.6
```

```
debt      NaN
```

```
Name: three, dtype: object
```


DataFrame

- Columns can be modified by assignment. For example, the empty 'debt' column could be assigned a scalar value or an array of values:

```
In [54]: frame2['debt'] = 16.5
```

```
In [55]: frame2
```

```
Out[55]:
```

	year	state	pop	debt
one	2000	Ohio	1.5	16.5
two	2001	Ohio	1.7	16.5
three	2002	Ohio	3.6	16.5
four	2001	Nevada	2.4	16.5
five	2002	Nevada	2.9	16.5
six	2003	Nevada	3.2	16.5

```
In [56]: frame2['debt'] = np.arange(6.)
```

```
In [57]: frame2
```

```
Out[57]:
```

	year	state	pop	debt
one	2000	Ohio	1.5	0.0
two	2001	Ohio	1.7	1.0
three	2002	Ohio	3.6	2.0
four	2001	Nevada	2.4	3.0
five	2002	Nevada	2.9	4.0
six	2003	Nevada	3.2	5.0

DataFrame

- When you are assigning lists or arrays to a column, the value's length must match the length of the DataFrame.
- If you assign a Series, its labels will be realigned exactly to the DataFrame's index, inserting missing values in any holes:

```
In [58]: val = pd.Series([-1.2, -1.5, -1.7], index=['two', 'four', 'five'])
```

```
In [59]: frame2['debt'] = val
```

```
In [60]: frame2
```

```
Out[60]:
```

	year	state	pop	debt
one	2000	Ohio	1.5	NaN
two	2001	Ohio	1.7	-1.2
three	2002	Ohio	3.6	NaN
four	2001	Nevada	2.4	-1.5
five	2002	Nevada	2.9	-1.7
six	2003	Nevada	3.2	NaN

DataFrame

- Assigning a column that doesn't exist will create a new column.
- The `del` keyword will delete columns as with a dict.
- As an example of `del` , I first add a new column of boolean values where the `state` column equals 'Ohio' :

```
In [61]: frame2['eastern'] = frame2.state == 'Ohio'
```

```
In [62]: frame2
```

```
Out[62]:
```

	year	state	pop	debt	eastern
one	2000	Ohio	1.5	NaN	True
two	2001	Ohio	1.7	-1.2	True
three	2002	Ohio	3.6	NaN	True
four	2001	Nevada	2.4	-1.5	False
five	2002	Nevada	2.9	-1.7	False
six	2003	Nevada	3.2	NaN	False

DataFrame

- New columns cannot be created with the `frame2.eastern` syntax.
- The `del` method can then be used to remove this column:

```
In [63]: del frame2['eastern']
```

```
In [64]: frame2.columns
```

```
Out[64]: Index(['year', 'state', 'pop', 'debt'], dtype='object')
```

DataFrame

- Another common form of data is a nested dict of dicts:

```
In [65]: pop = {'Nevada': {2001: 2.4, 2002: 2.9},  
.....:         'Ohio': {2000: 1.5, 2001: 1.7, 2002: 3.6}}
```

- If the nested dict is passed to the DataFrame, pandas will interpret the outer dict keys as the columns and the inner keys as the row indices:

```
In [66]: frame3 = pd.DataFrame(pop)
```

```
In [67]: frame3
```

```
Out[67]:
```

```
• .
```

	Nevada	Ohio
2000	NaN	1.5
2001	2.4	1.7
2002	2.9	3.6

DataFrame

- You can transpose the DataFrame (swap rows and columns) with similar syntax to a NumPy array:

```
In [68]: frame3.T
```

```
Out[68]:
```

	2000	2001	2002
Nevada	NaN	2.4	2.9
Ohio	1.5	1.7	3.6

- The keys in the inner dicts are combined and sorted to form the index in the result.

DataFrame

- If a DataFrame's `index` and `columns` have their name attributes set, these will also be displayed:

```
In [72]: frame3.index.name = 'year'; frame3.columns.name = 'state'
```

```
In [73]: frame3
```

```
Out[73]:
```

```
state  Nevada  Ohio
```

```
year
```

```
2000      NaN    1.5
```

```
2001      2.4    1.7
```

```
2002      2.9    3.6
```

DataFrame

- As with Series, the values attribute returns the data contained in the DataFrame as a two-dimensional ndarray:

```
In [74]: frame3.values
```

```
Out[74]:
```

```
array([[ nan,  1.5],  
       [ 2.4,  1.7],  
       [ 2.9,  3.6]])
```

```
In [75]: frame2.values
```

```
Out[75]:
```

```
array([[2000, 'Ohio', 1.5, nan],  
       [2001, 'Ohio', 1.7, -1.2],  
       [2002, 'Ohio', 3.6, nan],  
       [2001, 'Nevada', 2.4, -1.5],  
       [2002, 'Nevada', 2.9, -1.7],  
       [2003, 'Nevada', 3.2, nan]], dtype=object)
```


Possible data inputs to DataFrame constructor

Type	Notes
2D ndarray	A matrix of data, passing optional row and column labels
dict of arrays, lists, or tuples	Each sequence becomes a column in the DataFrame; all sequences must be the same length
NumPy structured/record array	Treated as the “dict of arrays” case
dict of Series	Each value becomes a column; indexes from each Series are unioned together to form the result’s row index if no explicit index is passed
dict of dicts	Each inner dict becomes a column; keys are unioned to form the row index as in the “dict of Series” case
List of dicts or Series	Each item becomes a row in the DataFrame; union of dict keys or Series indexes become the DataFrame’s column labels
List of lists or tuples	Treated as the “2D ndarray” case
Another DataFrame	The DataFrame’s indexes are used unless different ones are passed
NumPy MaskedArray	Like the “2D ndarray” case except masked values become NA/missing in the DataFrame result

Index Objects

- pandas's Index objects are responsible for holding the axis labels and other metadata (like the axis name or names).
- Any array or other sequence of labels you use when constructing a Series or DataFrame is internally converted to an Index:

```
In [76]: obj = pd.Series(range(3), index=['a', 'b', 'c'])
```

```
In [77]: index = obj.index
```

```
In [78]: index
```

- Out[78]: Index(['a', 'b', 'c'], dtype='object')

```
In [79]: index[1:]
```

```
Out[79]: Index(['b', 'c'], dtype='object')
```

Index Objects

- Index objects are immutable and thus can't be modified by the user:

```
index[1] = 'd' # TypeError
```

Index Objects

- Immutability makes it safer to share Index objects among data structures:

```
In [80]: labels = pd.Index(np.arange(3))
```

```
In [81]: labels
```

```
Out[81]: Int64Index([0, 1, 2], dtype='int64')
```

```
In [82]: obj2 = pd.Series([1.5, -2.5, 0], index=labels)
```

```
In [83]: obj2
```

```
Out[83]:
```

```
0      1.5
```

```
1     -2.5
```

```
2       0.0
```

```
dtype: float64
```

```
In [84]: obj2.index is labels
```

```
Out[84]: True
```

Index Objects

- In addition to being array-like, an Index also behaves like a fixed-size set:

```
In [85]: frame3
Out[85]:
```

state	Nevada	Ohio
year		
2000	NaN	1.5
2001	2.4	1.7
2002	2.9	3.6

```
In [86]: frame3.columns
Out[86]: Index(['Nevada', 'Ohio'], dtype='object', name='state')
```

```
In [87]: 'Ohio' in frame3.columns
Out[87]: True
```

```
In [88]: 2003 in frame3.index
Out[88]: False
```

Index Objects

- Unlike Python sets, a pandas Index can contain duplicate labels:

```
In [89]: dup_labels = pd.Index(['foo', 'foo', 'bar', 'bar'])
```

```
In [90]: dup_labels
```

```
Out[90]: Index(['foo', 'foo', 'bar', 'bar'], dtype='object')
```

Index Objects

- Selections with duplicate labels will select all occurrences of that label.
- Each Index has a number of methods and properties for set logic, which answer other common questions about the data it contains.

Some Index methods and properties

Method	Description
append	Concatenate with additional Index objects, producing a new Index
difference	Compute set difference as an Index
intersection	Compute set intersection
union	Compute set union
isin	Compute boolean array indicating whether each value is contained in the passed collection
delete	Compute new Index with element at index <i>i</i> deleted
drop	Compute new Index by deleting passed values
insert	Compute new Index by inserting element at index <i>i</i>
is_monotonic	Returns True if each element is greater than or equal to the previous element
is_unique	Returns True if the Index has no duplicate values
unique	Compute the array of unique values in the Index

Essential Functionality

- This section will walk you through the fundamental mechanics of interacting with the data contained in a Series or DataFrame.

Reindexing

- An important method on pandas objects is `reindex`, which means to create a new object with the data conformed to a new index.
- Consider an example:

```
In [91]: obj = pd.Series([4.5, 7.2, -5.3, 3.6], index=['d', 'b', 'a', 'c'])
```

```
In [92]: obj
```

```
Out[92]:
```

```
d      4.5
```

```
b      7.2
```

```
a     -5.3
```

```
c      3.6
```

```
dtype: float64
```

Reindexing

- Calling `reindex` on this Series rearranges the data according to the new index, introducing missing values if any index values were not already present:

```
In [93]: obj2 = obj.reindex(['a', 'b', 'c', 'd', 'e'])
```

```
In [94]: obj2
```

```
Out[94]:
```

```
a    -5.3
```

```
b     7.2
```

```
c     3.6
```

```
d     4.5
```

```
e     NaN
```

```
dtype: float64
```

Reindexing

- For ordered data like time series, it may be desirable to do some interpolation or filling of values when reindexing.

```
In [95]: obj3 = pd.Series(['blue', 'purple', 'yellow'], index=[0, 2, 4])
```

```
In [96]: obj3
```

```
Out[96]:
```

```
0      blue
```

```
2    purple
```

```
4    yellow
```

```
dtype: object
```

Reindexing

- The method option allows us to do this, using a method such as `ffill` , which forward-fills the values:

```
In [97]: obj3.reindex(range(6), method='ffill')
```

```
Out[97]:
```

```
0      blue
```

```
1      blue
```

```
2    purple
```

```
3    purple
```

```
4    yellow
```

```
5    yellow
```

```
dtype: object
```

Reindexing

- With DataFrame, reindex can alter either the (row) index, columns, or both.

```
In [98]: frame = pd.DataFrame(np.arange(9).reshape((3, 3)),  
.....:                        index=['a', 'c', 'd'],  
.....:                        columns=['Ohio', 'Texas', 'California'])
```

```
In [99]: frame
```

```
Out[99]:
```

	Ohio	Texas	California
a	0	1	2
c	3	4	5
d	6	7	8

Reindexing

- When passed only a sequence, it reindexes the rows in the result:

```
In [100]: frame2 = frame.reindex(['a', 'b', 'c', 'd'])
```

```
In [101]: frame2
```

```
Out[101]:
```

	Ohio	Texas	California
a	0.0	1.0	2.0
b	NaN	NaN	NaN
c	3.0	4.0	5.0
d	6.0	7.0	8.0

Reindexing

- The columns can be reindexed with the columns keyword:

```
In [102]: states = ['Texas', 'Utah', 'California']
```

```
In [103]: frame.reindex(columns=states)
```

```
Out[103]:
```

	Texas	Utah	California
a	1	NaN	2
c	4	NaN	5
d	7	NaN	8

Reindexing

- As we'll explore in more detail, you can reindex more succinctly by label-indexing with `loc`, and many users prefer to use it exclusively:

```
In [104]: frame.loc[['a', 'b', 'c', 'd'], states]
```

```
Out[104]:
```

	Texas	Utah	California
a	1.0	NaN	2.0
b	NaN	NaN	NaN
c	4.0	NaN	5.0
d	7.0	NaN	8.0

Reindex function arguments

Argument	Description
<code>index</code>	New sequence to use as index. Can be Index instance or any other sequence-like Python data structure. An Index will be used exactly as is without any copying.
<code>method</code>	Interpolation (fill) method; 'ffill' fills forward, while 'bfill' fills backward.
<code>fill_value</code>	Substitute value to use when introducing missing data by reindexing.
<code>limit</code>	When forward- or backfilling, maximum size gap (in number of elements) to fill.
<code>tolerance</code>	When forward- or backfilling, maximum size gap (in absolute numeric distance) to fill for inexact matches.
<code>level</code>	Match simple Index on level of MultiIndex; otherwise select subset of.
<code>copy</code>	If <code>True</code> , always copy underlying data even if new index is equivalent to old index; if <code>False</code> , do not copy the data when the indexes are equivalent.

Dropping Entries from an Axis

- Dropping one or more entries from an axis is easy if you already have an index array or list without those entries.
- As that can require a bit of munging and set logic, the `drop` method will return a new object with the indicated value or values deleted from an axis:

```
In [105]: obj = pd.Series(np.arange(5.), index=['a', 'b', 'c', 'd', 'e'])
```

```
In [106]: obj
```

```
Out[106]:
```

```
a    0.0
```

```
b    1.0
```

```
c    2.0
```

```
d    3.0
```

```
e    4.0
```

```
dtype: float64
```

Dropping Entries from an Axis

```
In [107]: new_obj = obj.drop('c')
```

```
In [108]: new_obj
```

```
Out[108]:
```

```
a    0.0
```

```
b    1.0
```

```
d    3.0
```

```
e    4.0
```

```
dtype: float64
```

```
In [109]: obj.drop(['d', 'c'])
```

```
Out[109]:
```

```
a    0.0
```

```
b    1.0
```

```
e    4.0
```

```
dtype: float64
```

Dropping Entries from an Axis

- With DataFrame, index values can be deleted from either axis.
- To illustrate this, we first create an example DataFrame:

```
In [110]: data = pd.DataFrame(np.arange(16).reshape((4, 4)),  
.....:                        index=['Ohio', 'Colorado', 'Utah', 'New York'],  
.....:                        columns=['one', 'two', 'three', 'four'])
```

```
In [111]: data
```

```
Out[111]:
```

	one	two	three	four
Ohio	0	1	2	3
Colorado	4	5	6	7
Utah	8	9	10	11
New York	12	13	14	15

Dropping Entries from an Axis

- Calling `drop` with a sequence of labels will drop values from the row labels (axis 0):

```
In [112]: data.drop(['Colorado', 'Ohio'])
```

```
Out[112]:
```

	one	two	three	four
Utah	8	9	10	11
New York	12	13	14	15

Dropping Entries from an Axis

- You can drop values from the columns by passing `axis=1` or `axis='columns'` :

```
In [113]: data.drop('two', axis=1)
```

```
Out[113]:
```

	one	three	four
Ohio	0	2	3
Colorado	4	6	7
Utah	8	10	11
New York	12	14	15

```
In [114]: data.drop(['two', 'four'], axis='columns')
```

```
Out[114]:
```

	one	three
Ohio	0	2
Colorado	4	6
Utah	8	10
New York	12	14

Dropping Entries from an Axis

- Many functions, like `drop`, which modify the size or shape of a Series or DataFrame, can manipulate an object in-place without returning a new object:

```
In [115]: obj.drop('c', inplace=True)
```

```
In [116]: obj
```

```
Out[116]:
```

```
a    0.0
```

```
b    1.0
```

```
d    3.0
```

```
e    4.0
```

```
dtype: float64
```


Indexing, Selection, and Filtering

- Series indexing (`obj[...]`) works analogously to NumPy array indexing, except you can use the Series's index values instead of only integers.

```
In [117]: obj = pd.Series(np.arange(4.), index=['a', 'b', 'c', 'd'])
```

```
In [118]: obj
```

```
Out[118]:
```

```
a    0.0
```

```
b    1.0
```

```
c    2.0
```

```
d    3.0
```

```
dtype: float64
```

```
In [119]: obj['b']
```

```
Out[119]: 1.0
```

Indexing, Selection, and Filtering

```
In [120]: obj[1]
```

```
Out[120]: 1.0
```

```
In [121]: obj[2:4]
```

```
Out[121]:
```

```
c      2.0
```

```
d      3.0
```

```
dtype: float64
```

```
In [122]: obj[['b', 'a', 'd']]
```

```
Out[122]:
```

```
b      1.0
```

```
a      0.0
```

```
d      3.0
```

```
dtype: float64
```

```
In [123]: obj[[1, 3]]
```

```
Out[123]:
```

```
b      1.0
```

```
d      3.0
```

```
dtype: float64
```

```
In [124]: obj[obj < 2]
```

```
Out[124]:
```

```
a      0.0
```

```
b      1.0
```

```
dtype: float64
```

Indexing, Selection, and Filtering

- Slicing with labels behaves differently than normal Python slicing in that the end-point is inclusive:

```
In [125]: obj['b':'c']
```

```
Out[125]:
```

```
b      1.0
```

```
c      2.0
```

```
dtype: float64
```

```
In [126]: obj['b':'c'] = 5
```

```
In [127]: obj
```

```
Out[127]:
```

```
a      0.0
```

```
b      5.0
```

```
c      5.0
```

```
d      3.0
```

```
dtype: float64
```

Indexing, Selection, and Filtering

- Indexing into a DataFrame is for retrieving one or more columns either with a single value or sequence:

```
In [128]: data = pd.DataFrame(np.arange(16).reshape((4, 4)),  
.....:                        index=['Ohio', 'Colorado', 'Utah', 'New York'],  
.....:                        columns=['one', 'two', 'three', 'four'])
```

```
In [129]: data
```

```
Out[129]:
```

	one	two	three	four
Ohio	0	1	2	3
Colorado	4	5	6	7
Utah	8	9	10	11
New York	12	13	14	15

Indexing, Selection, and Filtering

- Another use case is in indexing with a boolean DataFrame, such as one produced by a scalar comparison:

```
In [134]: data < 5
```

```
Out[134]:
```

	one	two	three	four
Ohio	True	True	True	True
Colorado	True	False	False	False
Utah	False	False	False	False
New York	False	False	False	False

```
In [135]: data[data < 5] = 0
```

```
In [136]: data
```

```
Out[136]:
```

	one	two	three	four
Ohio	0	0	0	0
Colorado	0	5	6	7
Utah	8	9	10	11
New York	12	13	14	15

Indexing, Selection, and Filtering

```
In [130]: data['two']
```

```
Out[130]:
```

```
Ohio      1
```

```
Colorado  5
```

```
Utah      9
```

```
New York 13
```

```
Name: two, dtype: int64
```

```
In [131]: data[['three', 'one']]
```

```
Out[131]:
```

```
      three  one
```

```
Ohio      2    0
```

```
Colorado  6    4
```

```
Utah     10    8
```

```
New York 14   12
```

Indexing, Selection, and Filtering

- Indexing like this has a few special cases. First, slicing or selecting data with a boolean array:

```
In [132]: data[:2]
```

```
Out[132]:
```

	one	two	three	four
Ohio	0	1	2	3
Colorado	4	5	6	7

```
In [133]: data[data['three'] > 5]
```

```
Out[133]:
```

	one	two	three	four
Colorado	4	5	6	7
Utah	8	9	10	11
New York	12	13	14	15

Indexing, Selection, and Filtering

- Another use case is in indexing with a boolean DataFrame, such as one produced by a scalar comparison:

```
In [134]: data < 5
```

```
Out[134]:
```

	one	two	three	four
Ohio	True	True	True	True
Colorado	True	False	False	False
Utah	False	False	False	False
New York	False	False	False	False

```
In [135]: data[data < 5] = 0
```

```
In [136]: data
```

```
Out[136]:
```

	one	two	three	four
Ohio	0	0	0	0
Colorado	0	5	6	7
Utah	8	9	10	11
New York	12	13	14	15

Selection with loc and iloc

- For DataFrame label-indexing on the rows, we introduce the special indexing operators `loc` and `iloc`.
- They enable you to select a subset of the rows and columns from a DataFrame with NumPy-like notation using either axis labels (`loc`) or integers (`iloc`).

Selection with loc and iloc

- As a preliminary example, let's select a single row and multiple columns by label:

```
In [137]: data.loc['Colorado', ['two', 'three']]  
Out[137]:  
two      5  
three    6  
Name: Colorado, dtype: int64
```

Selection with loc and iloc

- We'll then perform some similar selections with integers using `iloc`:

```
In [138]: data.iloc[2, [3, 0, 1]]
```

```
Out[138]:
```

```
four      11
```

```
one        8
```

```
two        9
```

```
Name: Utah, dtype: int64
```

```
In [139]: data.iloc[2]
```

```
Out[139]:
```

```
one        8
```

```
two        9
```

```
three     10
```

```
four     11
```

```
Name: Utah, dtype: int64
```

```
In [140]: data.iloc[[1, 2], [3, 0, 1]]
```

```
Out[140]:
```

	four	one	two
Colorado	7	0	5
Utah	11	8	9

Selection with loc and iloc

- Both indexing functions work with slices in addition to single labels or lists of labels:

```
In [141]: data.loc[:, 'Utah', 'two']
```

```
Out[141]:
```

```
Ohio      0
```

```
Colorado  5
```

```
Utah      9
```

```
Name: two, dtype: int64
```

```
In [142]: data.iloc[:, :3][data.three > 5]
```

```
Out[142]:
```

	one	two	three
Colorado	0	5	6
Utah	8	9	10
New York	12	13	14

Indexing options with DataFrame

Type	Notes
<code>df[val]</code>	Select single column or sequence of columns from the DataFrame; special case conveniences: boolean array (filter rows), slice (slice rows), or boolean DataFrame (set values based on some criterion)
<code>df.loc[val]</code>	Selects single row or subset of rows from the DataFrame by label
<code>df.loc[:, val]</code>	Selects single column or subset of columns by label
<code>df.loc[val1, val2]</code>	Select both rows and columns by label
<code>df.iloc[where]</code>	Selects single row or subset of rows from the DataFrame by integer position
<code>df.iloc[:, where]</code>	Selects single column or subset of columns by integer position
<code>df.iloc[where_i, where_j]</code>	Select both rows and columns by integer position
<code>df.at[label_i, label_j]</code>	Select a single scalar value by row and column label
<code>df.iat[i, j]</code>	Select a single scalar value by row and column position (integers)
reindex method	Select either rows or columns by labels
get_value, set_value methods	Select single value by row and column label

Arithmetic and Data Alignment

- When you are adding together objects, if any index pairs are not the same, the respective index in the result will be the union of the index pairs.
- For users with database experience, this is similar to an automatic outer join on the index labels.

Arithmetic and Data Alignment

- Let's look at an example:

```
In [151]: s2 = pd.Series([-2.1, 3.6, -1.5, 4, 3.1],  
.....:                  index=['a', 'c', 'e', 'f', 'g'])
```

```
In [152]: s1
```

```
Out[152]:
```

```
a    7.3
```

```
c   -2.5
```

```
d    3.4
```

```
e    1.5
```

```
dtype: float64
```

```
In [153]: s2
```

```
Out[153]:
```

```
a   -2.1
```

```
c    3.6
```

```
e   -1.5
```

```
f    4.0
```

```
g    3.1
```

```
dtype: float64
```

Arithmetic and Data Alignment

- Adding these together yields:

```
In [154]: s1 + s2
Out[154]:
a      5.2
c      1.1
d      NaN
e      0.0
f      NaN
g      NaN
dtype: float64
```


Arithmetic and Data Alignment

- In the case of DataFrame, alignment is performed on both the rows and the columns:

```
In [155]: df1 = pd.DataFrame(np.arange(9.).reshape((3, 3)), columns=list('bcd'),  
.....:                      index=['Ohio', 'Texas', 'Colorado'])
```

```
In [156]: df2 = pd.DataFrame(np.arange(12.).reshape((4, 3)), columns=list('bde'),  
.....:                      index=['Utah', 'Ohio', 'Texas', 'Oregon'])
```

```
In [157]: df1
```

```
Out[157]:
```

	b	c	d
Ohio	0.0	1.0	2.0
Texas	3.0	4.0	5.0
Colorado	6.0	7.0	8.0

```
In [158]: df2
```

```
Out[158]:
```

	b	d	e
Utah	0.0	1.0	2.0
Ohio	3.0	4.0	5.0
Texas	6.0	7.0	8.0
Oregon	9.0	10.0	11.0

Arithmetic and Data Alignment

- Adding these together returns a DataFrame whose index and columns are the unions of the ones in each DataFrame:

```
In [159]: df1 + df2
```

```
Out[159]:
```

	b	c	d	e
Colorado	NaN	NaN	NaN	NaN
Ohio	3.0	NaN	6.0	NaN
Oregon	NaN	NaN	NaN	NaN
Texas	9.0	NaN	12.0	NaN
Utah	NaN	NaN	NaN	NaN

Arithmetic and Data Alignment

- Since the 'c' and 'e' columns are not found in both DataFrame objects, they appear as all missing in the result.
- The same holds for the rows whose labels are not common to both objects.

Arithmetic and Data Alignment

- If you add DataFrame objects with no column or row labels in common, the result will contain all nulls:

```
In [161]: df2 = pd.DataFrame({'B': [3, 4]})
```

```
In [162]: df1
```

```
Out[162]:
```

```
   A
0  1
1  2
```

```
In [163]: df2
```

```
Out[163]:
```

```
   B
0  3
1  4
```

```
In [164]: df1 - df2
```

```
Out[164]:
```

```
   A  B
0 NaN NaN
1 NaN NaN
```

Arithmetic methods with fill values

- In arithmetic operations between differently indexed objects, you might want to fill with a special value, like 0, when an axis label is found in one object but not the other:

```
In [165]: df1 = pd.DataFrame(np.arange(12.).reshape((3, 4)),  
.....:                      columns=list('abcd'))
```

```
In [166]: df2 = pd.DataFrame(np.arange(20.).reshape((4, 5)),  
.....:                      columns=list('abcde'))
```

```
In [167]: df2.loc[1, 'b'] = np.nan
```

```
In [168]: df1
```

```
Out[168]:
```

	a	b	c	d
0	0.0	1.0	2.0	3.0
1	4.0	5.0	6.0	7.0
2	8.0	9.0	10.0	11.0

```
In [169]: df2
```

```
Out[169]:
```

	a	b	c	d	e
0	0.0	1.0	2.0	3.0	4.0
1	5.0	NaN	7.0	8.0	9.0
2	10.0	11.0	12.0	13.0	14.0
3	15.0	16.0	17.0	18.0	19.0

Arithmetic methods with fill values

- Adding these together results in NA values in the locations that don't overlap:

```
In [170]: df1 + df2
```

```
Out[170]:
```

	a	b	c	d	e
0	0.0	2.0	4.0	6.0	NaN
1	9.0	NaN	13.0	15.0	NaN
2	18.0	20.0	22.0	24.0	NaN
3	NaN	NaN	NaN	NaN	NaN

Arithmetic methods with fill values

- Using the add method on df1 , I pass df2 and an argument to fill_value :

```
In [171]: df1.add(df2, fill_value=0)
```

```
Out[171]:
```

	a	b	c	d	e
0	0.0	2.0	4.0	6.0	4.0
1	9.0	5.0	13.0	15.0	9.0
2	18.0	20.0	22.0	24.0	14.0
3	15.0	16.0	17.0	18.0	19.0

Flexible arithmetic methods

Method	Description
<code>add, radd</code>	Methods for addition (+)
<code>sub, rsub</code>	Methods for subtraction (-)
<code>div, rdiv</code>	Methods for division (/)
<code>floordiv, rfloordiv</code>	Methods for floor division (//)
<code>mul, rmul</code>	Methods for multiplication (*)
<code>pow, rpow</code>	Methods for exponentiation (**)

Operations between DataFrame and Series

- As with NumPy arrays of different dimensions, arithmetic between DataFrame and Series is also defined.

Operations between DataFrame and Series

- First, as a motivating example, consider the difference between a two-dimensional array and one of its rows:

```
In [175]: arr = np.arange(12.).reshape((3, 4))
```

```
In [176]: arr
```

```
Out[176]:
```

```
array([[ 0.,  1.,  2.,  3.],
       [ 4.,  5.,  6.,  7.],
       [ 8.,  9., 10., 11.]])
```

```
In [177]: arr[0]
```

```
Out[177]: array([ 0.,  1.,  2.,  3.])
```

```
In [178]: arr - arr[0]
```

```
Out[178]:
```

```
array([[ 0.,  0.,  0.,  0.],
       [ 4.,  4.,  4.,  4.],
       [ 8.,  8.,  8.,  8.]])
```

Operations between DataFrame and Series

- When we subtract `arr[0]` from `arr`, the subtraction is performed once for each row.
- This is referred to as broadcasting.

Operations between DataFrame and Series

- Operations between a DataFrame and a Series are similar:

```
In [179]: frame = pd.DataFrame(np.arange(12.).reshape((4, 3)),  
.....:                        columns=list('bde'),  
.....:                        index=['Utah', 'Ohio', 'Texas', 'Oregon'])
```

```
In [180]: series = frame.iloc[0]
```

```
In [181]: frame
```

```
Out[181]:
```

	b	d	e
Utah	0.0	1.0	2.0
Ohio	3.0	4.0	5.0
Texas	6.0	7.0	8.0
Oregon	9.0	10.0	11.0

```
In [182]: series
```

```
Out[182]:
```

```
b    0.0
```

```
d    1.0
```

```
e    2.0
```

```
Name: Utah, dtype: float64
```

Operations between DataFrame and Series

- By default, arithmetic between DataFrame and Series matches the index of the Series on the DataFrame's columns, broadcasting down the rows:

```
In [183]: frame - series
```

```
Out[183]:
```

	b	d	e
Utah	0.0	0.0	0.0
Ohio	3.0	3.0	3.0
Texas	6.0	6.0	6.0
Oregon	9.0	9.0	9.0

Operations between DataFrame and Series

- If an index value is not found in either the DataFrame's columns or the Series's index, the objects will be reindexed to form the union:

```
In [184]: series2 = pd.Series(range(3), index=['b', 'e', 'f'])
```

```
In [185]: frame + series2
```

```
Out[185]:
```

	b	d	e	f
Utah	0.0	NaN	3.0	NaN
Ohio	3.0	NaN	6.0	NaN
Texas	6.0	NaN	9.0	NaN
Oregon	9.0	NaN	12.0	NaN

Operations between DataFrame and Series

- If you want to instead broadcast over the columns, matching on the rows, you have to use one of the arithmetic methods. For example:

```
In [186]: series3 = frame['d']
```

```
In [187]: frame
```

```
Out[187]:
```

	b	d	e
Utah	0.0	1.0	2.0
Ohio	3.0	4.0	5.0
Texas	6.0	7.0	8.0
Oregon	9.0	10.0	11.0

```
In [188]: series3
```

```
Out[188]:
```

Utah	1.0
Ohio	4.0
Texas	7.0
Oregon	10.0

Name: d, dtype: float64

```
In [189]: frame.sub(series3, axis='index')
```

```
Out[189]:
```

	b	d	e
Utah	-1.0	0.0	1.0
Ohio	-1.0	0.0	1.0
Texas	-1.0	0.0	1.0
Oregon	-1.0	0.0	1.0

Operations between DataFrame and Series

- The axis number that you pass is the axis to match on.
- In this case we mean to match on the DataFrame's row index (`axis='index'` or `axis=0`) and broadcast across.

Function Application and Mapping

- NumPy ufuncs (element-wise array methods) also work with pandas objects:

```
In [190]: frame = pd.DataFrame(np.random.randn(4, 3), columns=list('bde'),  
.....:                        index=['Utah', 'Ohio', 'Texas', 'Oregon'])
```

```
In [191]: frame
```

```
Out[191]:
```

	b	d	e
Utah	-0.204708	0.478943	-0.519439
Ohio	-0.555730	1.965781	1.393406
Texas	0.092908	0.281746	0.769023
Oregon	1.246435	1.007189	-1.296221

```
In [192]: np.abs(frame)
```

```
Out[192]:
```

	b	d	e
Utah	0.204708	0.478943	0.519439
Ohio	0.555730	1.965781	1.393406
Texas	0.092908	0.281746	0.769023
Oregon	1.246435	1.007189	1.296221

Function Application and Mapping

- Another frequent operation is applying a function on one-dimensional arrays to each column or row.
- DataFrame's `apply` method does exactly this:

```
In [193]: f = lambda x: x.max() - x.min()
```

```
In [194]: frame.apply(f)
```

```
Out[194]:
```

```
b      1.802165
```

```
d      1.684034
```

```
e      2.689627
```

```
dtype: float64
```

Function Application and Mapping

- In Last Example the function f , which computes the difference between the maximum and minimum of a Series, is invoked once on each column in frame .
- If you pass `axis='columns'` to `apply` , the function will be invoked once per row instead:

```
In [195]: frame.apply(f, axis='columns')
Out[195]:
Utah      0.998382
Ohio      2.521511
Texas     0.676115
Oregon    2.542656
dtype: float64
```

Function Application and Mapping

- Many of the most common array statistics (like sum and mean) are DataFrame methods, so using apply is not necessary.
- The function passed to apply need not return a scalar value; it can also return a Series with multiple values:

```
In [196]: def f(x):  
         ....:     return pd.Series([x.min(), x.max()], index=['min', 'max'])
```

```
In [197]: frame.apply(f)
```

```
Out[197]:
```

	b	d	e
min	-0.555730	0.281746	-1.296221
max	1.246435	1.965781	1.393406

Function Application and Mapping

- Element-wise Python functions can be used, too.
- Suppose you wanted to compute a formatted string from each floating-point value in frame .

Function Application and Mapping

- You can do this with `apply` `map`

```
In [198]: format = lambda x: '%.2f' % x
```

```
In [199]: frame.applymap(format)
```

```
Out[199]:
```

	b	d	e
Utah	-0.20	0.48	-0.52
Ohio	-0.56	1.97	1.39
Texas	0.09	0.28	0.77
Oregon	1.25	1.01	-1.30

Function Application and Mapping

- The reason for the name `applymap` is that `Series` has a `map` method for applying an element-wise function:

```
In [200]: frame['e'].map(format)
Out[200]:
Utah      -0.52
Ohio       1.39
Texas      0.77
Oregon    -1.30
Name: e, dtype: object
```

Sorting and Ranking

- Sorting a dataset by some criterion is another important built-in operation.
- To sort lexicographically by row or column index, use the `sort_index` method, which returns a new, sorted object:

```
In [201]: obj = pd.Series(range(4), index=['d', 'a', 'b', 'c'])
```

```
In [202]: obj.sort_index()
```

```
Out[202]:
```

```
a    1
```

```
b    2
```

```
c    3
```

```
d    0
```

```
dtype: int64
```


Sorting and Ranking

- With a DataFrame, you can sort by index on either axis:

```
In [203]: frame = pd.DataFrame(np.arange(8).reshape((2, 4)),  
.....:                        index=['three', 'one'],  
.....:                        columns=['d', 'a', 'b', 'c'])
```

```
In [204]: frame.sort_index()
```

```
Out[204]:
```

	d	a	b	c
one	4	5	6	7
three	0	1	2	3

```
In [205]: frame.sort_index(axis=1)
```

```
Out[205]:
```

	a	b	c	d
three	1	2	3	0
one	5	6	7	4

Sorting and Ranking

- The data is sorted in ascending order by default, but can be sorted in descending order, too:

```
In [206]: frame.sort_index(axis=1, ascending=False)
```

```
Out[206]:
```

	d	c	b	a
three	0	3	2	1
one	4	7	6	5

- To sort a Series by its values, use its `sort_values` method:

```
In [207]: obj = pd.Series([4, 7, -3, 2])
```

```
In [208]: obj.sort_values()
```

```
Out[208]:
```

```
2    -3
```

```
3     2
```

```
0     4
```

```
1     7
```

```
dtype: int64
```

Sorting and Ranking

- Any missing values are sorted to the end of the Series by default:

```
In [209]: obj = pd.Series([4, np.nan, 7, np.nan, -3, 2])
```

```
In [210]: obj.sort_values()
```

```
Out[210]:
```

```
4    -3.0
```

```
5     2.0
```

```
0     4.0
```

```
2     7.0
```

```
1     NaN
```

```
3     NaN
```

```
dtype: float64
```

Sorting and Ranking

- When sorting a DataFrame, you can use the data in one or more columns as the sort keys. To do so, pass one or more column names to the `by` option of `sort_values`:

```
In [211]: frame = pd.DataFrame({'b': [4, 7, -3, 2], 'a': [0, 1, 0, 1]})
```

```
In [212]: frame
```

```
Out[212]:
```

	a	b
0	0	4
1	1	7
2	0	-3
3	1	2

```
In [213]: frame.sort_values(by='b')
```

```
Out[213]:
```

	a	b
2	0	-3
3	1	2
0	0	4
1	1	7

Sorting and Ranking

- To sort by multiple columns, pass a list of names:

```
In [214]: frame.sort_values(by=['a', 'b'])
```

```
Out[214]:
```

	a	b
2	0	-3
0	0	4
3	1	2
1	1	7

- Ranking assigns ranks from one through the number of valid data points in an array.

Sorting and Ranking

- by default rank breaks ties by assigning each group the mean rank:

```
In [215]: obj = pd.Series([7, -5, 7, 4, 2, 0, 4])
```

```
In [216]: obj.rank()
```

```
Out[216]:
```

```
0    6.5
```

```
1    1.0
```

```
2    6.5
```

```
3    4.5
```

```
4    3.0
```

```
5    2.0
```

```
6    4.5
```

```
dtype: float64
```

Sorting and Ranking

- Ranks can also be assigned according to the order in which they're observed in the data:

```
In [217]: obj.rank(method='first')
```

```
Out[217]:
```

```
0      6.0
```

```
1      1.0
```

```
2      7.0
```

```
3      4.0
```

```
4      3.0
```

```
5      2.0
```

```
6      5.0
```

```
dtype: float64
```

- Here, instead of using the average rank 6.5 for the entries 0 and 2, they instead have been set to 6 and 7 because label 0 precedes label 2 in the data.

Sorting and Ranking

- You can rank in descending order, too:

```
# Assign tie values the maximum rank in the group  
In [218]: obj.rank(ascending=False, method='max')  
Out[218]:  
0      2.0  
1      7.0  
2      2.0  
3      4.0  
4      5.0  
5      6.0  
6      4.0  
dtype: float64
```


Sorting and Ranking

- DataFrame can compute ranks over the rows or the columns:

```
In [219]: frame = pd.DataFrame({'b': [4.3, 7, -3, 2], 'a': [0, 1, 0, 1],  
.....:                        'c': [-2, 5, 8, -2.5]})
```

```
In [220]: frame
```

```
Out[220]:
```

	a	b	c
0	0	4.3	-2.0
1	1	7.0	5.0
2	0	-3.0	8.0
3	1	2.0	-2.5

```
In [221]: frame.rank(axis='columns')
```

```
Out[221]:
```

	a	b	c
0	2.0	3.0	1.0
1	1.0	3.0	2.0
2	2.0	1.0	3.0
3	2.0	3.0	1.0

Tie-breaking methods with rank

Method	Description
'average'	Default: assign the average rank to each entry in the equal group
'min'	Use the minimum rank for the whole group
'max'	Use the maximum rank for the whole group
'first'	Assign ranks in the order the values appear in the data
'dense'	Like <code>method='min'</code> , but ranks always increase by 1 in between groups rather than the number of equal elements in a group

Axis Indexes with Duplicate Labels

- Up until now all of the examples we've looked at have had unique axis labels (index values).
- While many pandas functions (like `reindex`) require that the labels be unique, it's not mandatory.

Axis Indexes with Duplicate Labels

- Let's consider a small Series with duplicate indices:

```
In [222]: obj = pd.Series(range(5), index=['a', 'a', 'b', 'b', 'c'])
```

```
In [223]: obj
```

```
Out[223]:
```

```
a    0
```

```
a    1
```

```
b    2
```

```
b    3
```

```
c    4
```

```
dtype: int64
```

Axis Indexes with Duplicate Labels

- The index's `is_unique` property can tell you whether its labels are unique or not:

```
In [224]: obj.index.is_unique  
Out[224]: False
```

- Data selection is one of the main things that behaves differently with duplicates.
- Indexing a label with multiple entries returns a Series, while single entries return a scalar value:

```
In [225]: obj['a']  
Out[225]:  
a      0  
a      1  
dtype: int64
```

- .

```
In [226]: obj['c']  
Out[226]: 4
```

Axis Indexes with Duplicate Labels

- The same logic extends to indexing rows in a DataFrame:

```
In [227]: df = pd.DataFrame(np.random.randn(4, 3), index=['a', 'a', 'b', 'b'])
```

```
In [228]: df
```

```
Out[228]:
```

	0	1	2
a	0.274992	0.228913	1.352917
a	0.886429	-2.001637	-0.371843
b	1.669025	-0.438570	-0.539741
b	0.476985	3.248944	-1.021228

```
In [229]: df.loc['b']
```

```
Out[229]:
```

	0	1	2
b	1.669025	-0.438570	-0.539741
b	0.476985	3.248944	-1.021228