

Lecture 13, 26 September 2023

Pandas (Python and data analysis)

- Built on top of numpy

Series and data frames

- Numpy defines homogeneous n-dimensional arrays
- Data science works with tables: 2-dimensional arrays
- Pandas has two fundamental data structures
 - Series : A column of data
 - Data Frame : A table of data

Key difference

- Numpy indices are always $[0 . . n-1]$ in each dimension
- Pandas allows more flexible "named" indices for rows and columns
 - Dictionary vs list

Load pandas

- Don't need to import numpy unless one is separately using numpy arrays

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

Create a series

- Convert a sequence into a series (column)

```
In [2]: h = ('AA', '2012-02-01', 100, 10.2)
s = pd.Series(h)
type(s)
```

```
Out[2]: pandas.core.series.Series
```

```
In [3]: s
```

```
Out[3]: 0      AA
1  2012-02-01
2      100
3      10.2
dtype: object
```

Convert a dictionary to a series

- Keys become "row indices"

```
In [4]: d = {'name' : 'IBM', 'date' : '2010-09-08', 'shares' : 100, 'price' : 10.2}
ds = pd.Series(d)
type(ds)
```

```
Out[4]: pandas.core.series.Series
```

```
In [5]: ds
```

```
Out[5]: name      IBM
date    2010-09-08
shares      100
price      10.2
dtype: object
```

Creating an index

```
In [6]: f = ['FB', '2001-08-02', 90, 3.2]
fs = pd.Series(f, index = ['name', 'date', 'shares', 'price'])
```

```
In [7]: fs
```

```
Out[7]: name          FB
date      2001-08-02
shares      90
price       3.2
dtype: object
```

Accessing elements

- Use named index, or position
- Use slices, sublists

```
In [8]: fs['shares']
```

```
Out[8]: 90
```

```
In [9]: fs[0], fs['name']
```

```
Out[9]: ('FB', 'FB')
```

```
In [10]: fs[0:2]
```

```
Out[10]: name          FB
date      2001-08-02
dtype: object
```

```
In [11]: fs[[0,2]]
```

```
Out[11]: name          FB
shares      90
dtype: object
```

```
In [12]: fs['name':'price']
```

```
Out[12]: name          FB
date      2001-08-02
shares      90
price       3.2
dtype: object
```

```
In [13]: fs['price':'name']
```

```
Out[13]: Series([], dtype: object)
```

```
In [14]: fs[['price', 'name']]
```

```
Out[14]: price      3.2
name          FB
dtype: object
```

```
In [15]: fs[0:3]
```

```
Out[15]: name          FB
date      2001-08-02
shares      90
dtype: object
```

Data frames

- A table is a sequence of columns
- A data frame is a sequence of series

```
In [16]: data2 = {'name' : ['AA', 'IBM', 'GOOG'],
                'date' : ['2001-12-01', '2012-02-10', '2010-04-09'],
                'shares' : [100, 30, 90],
                'price' : [12.3, 10.3, 32.2]
                }
df2 = pd.DataFrame(data2)
```

```
In [17]: df2
```

```
Out[17]:
```

	name	date	shares	price
0	AA	2001-12-01	100	12.3
1	IBM	2012-02-10	30	10.3
2	GOOG	2010-04-09	90	32.2

```
In [18]: data3 = (['AA', 'IBM', 'GOOG'],
                  ['2001-12-01', '2012-02-10', '2010-04-09'],
                  [100, 30, 90],
                  [12.3, 10.3, 32.2])
df3 = pd.DataFrame(data3)
```

```
In [19]: df3
```

```
Out[19]:
```

	0	1	2
0	AA	IBM	GOOG
1	2001-12-01	2012-02-10	2010-04-09
2	100	30	90
3	12.3	10.3	32.2

Add a column

```
In [20]: df2['owner'] = 'Unknown'
          # df2['owner'] = ['a', 'b', 'c']
df2
```

```
Out[20]:
```

	name	date	shares	price	owner
0	AA	2001-12-01	100	12.3	Unknown
1	IBM	2012-02-10	30	10.3	Unknown
2	GOOG	2010-04-09	90	32.2	Unknown

```
In [21]: #df2['owner'] = 'Unknown'
          df2['owner2'] = ['a', 'b', 'c']
df2
```

```
Out[21]:
```

	name	date	shares	price	owner	owner2
0	AA	2001-12-01	100	12.3	Unknown	a
1	IBM	2012-02-10	30	10.3	Unknown	b
2	GOOG	2010-04-09	90	32.2	Unknown	c

Add row indices

```
In [22]: df2.index = ['one', 'two', 'three']
```

```
In [23]: df2
```

```
Out[23]:
```

	name	date	shares	price	owner	owner2
one	AA	2001-12-01	100	12.3	Unknown	a
two	IBM	2012-02-10	30	10.3	Unknown	b
three	GOOG	2010-04-09	90	32.2	Unknown	c

Convert one of the columns into an index

```
In [24]: df2 = df2.set_index(['name'])
```

```
In [25]: df2
```

Out[25]:

	date	shares	price	owner	owner2
name					
AA	2001-12-01	100	12.3	Unknown	a
IBM	2012-02-10	30	10.3	Unknown	b
GOOG	2010-04-09	90	32.2	Unknown	c

Replace an index

```
In [26]: df2 = df2.set_index(['price'])
```

```
In [27]: df2
```

Out[27]:

	date	shares	owner	owner2
price				
12.3	2001-12-01	100	Unknown	a
10.3	2012-02-10	30	Unknown	b
32.2	2010-04-09	90	Unknown	c

Use multiple columns for indexing

```
In [28]: df2 = pd.DataFrame(data2)
df2['owner'] = 'Unknown'
df2 = df2.set_index(['name', 'price'])
```

```
In [29]: df2
```

Out[29]:

		date	shares	owner
name	price			
AA	12.3	2001-12-01	100	Unknown
IBM	10.3	2012-02-10	30	Unknown
GOOG	32.2	2010-04-09	90	Unknown

Accessing values in a dataframe

By column index

- Similar to projection in relational algebra

```
In [30]: df2[['shares', 'date']]
```

Out[30]:

		shares	date
name	price		
AA	12.3	100	2001-12-01
IBM	10.3	30	2012-02-10
GOOG	32.2	90	2010-04-09

By row index

```
In [31]: df2.loc['AA']
```

```
Out[31]:
```

	date	shares	price	owner
AA	2001-12-01	100	12.3	Unknown

Individual element by position

```
In [32]: df2.loc['AA', 'shares']
```

```
Out[32]: price
12.3    100
Name: shares, dtype: int64
```

Slices, etc

```
In [33]: df2.loc[:, 'shares']
```

```
Out[33]: name price
AA      12.3    100
IBM     10.3     30
GOOG    32.2     90
Name: shares, dtype: int64
```

```
In [34]: df2 = pd.DataFrame(data2)
df2['owner'] = 'Unknown'
df2 = df2.set_index(['name'])
df2
```

```
Out[34]:
```

	date	shares	price	owner
AA	2001-12-01	100	12.3	Unknown
IBM	2012-02-10	30	10.3	Unknown
GOOG	2010-04-09	90	32.2	Unknown

```
In [35]: df2.loc['AA':'IBM', 'shares':'owner']
```

```
Out[35]:
```

	shares	price	owner
AA	100	12.3	Unknown
IBM	30	10.3	Unknown

- Unlike series, cannot use position indices if "real" index exists

```
In [36]: df2 = pd.DataFrame(data2)
df2['owner'] = 'Unknown'
df2
```

```
Out[36]:
```

	name	date	shares	price	owner
0	AA	2001-12-01	100	12.3	Unknown
1	IBM	2012-02-10	30	10.3	Unknown
2	GOOG	2010-04-09	90	32.2	Unknown

```
In [37]: df2.loc[0:1]
```

```
Out[37]:
```

	name	date	shares	price	owner
0	AA	2001-12-01	100	12.3	Unknown
1	IBM	2012-02-10	30	10.3	Unknown

```
In [38]: df2 = df2.set_index(['name'])  
df2
```

Out[38]:

	date	shares	price	owner
name				
AA	2001-12-01	100	12.3	Unknown
IBM	2012-02-10	30	10.3	Unknown
GOOG	2010-04-09	90	32.2	Unknown

```
In [39]: df2.loc[0:2]
```

```
-----
TypeError                                Traceback (most recent call last)
Cell In[39], line 1
----> 1 df2.loc[0:2]

File ~/miniconda3/lib/python3.11/site-packages/pandas/core/indexing.py:1103, in _LocationIndexer._getitem_em__(self, key)
    1100 axis = self.axis or 0
    1102 maybe_callable = com.apply_if_callable(key, self.obj)
-> 1103 return self._getitem_axis(maybe_callable, axis=axis)

File ~/miniconda3/lib/python3.11/site-packages/pandas/core/indexing.py:1323, in _LocIndexer._getitem_axis(self, key, axis)
    1321 if isinstance(key, slice):
    1322     self._validate_key(key, axis)
-> 1323 return self._get_slice_axis(key, axis=axis)
    1324 elif com.is_bool_indexer(key):
    1325     return self._getbool_axis(key, axis=axis)

File ~/miniconda3/lib/python3.11/site-packages/pandas/core/indexing.py:1355, in _LocIndexer._get_slice_axis(self, slice_obj, axis)
    1352 return obj.copy(deep=False)
    1354 labels = obj._get_axis(axis)
-> 1355 indexer = labels.slice_indexer(slice_obj.start, slice_obj.stop, slice_obj.step)
    1357 if isinstance(indexer, slice):
    1358     return self.obj._slice(indexer, axis=axis)

File ~/miniconda3/lib/python3.11/site-packages/pandas/core/indexes/base.py:6344, in Index.slice_indexer(self, start, end, step)
    6300 def slice_indexer(
    6301     self,
    6302     start: Hashable | None = None,
    6303     end: Hashable | None = None,
    6304     step: int | None = None,
    6305 ) -> slice:
    6306     """
    6307     Compute the slice indexer for input labels and step.
    6308     (...)
    6342     slice(1, 3, None)
    6343     """
-> 6344 start_slice, end_slice = self.slice_locs(start, end, step=step)
    6346 # return a slice
    6347 if not is_scalar(start_slice):

File ~/miniconda3/lib/python3.11/site-packages/pandas/core/indexes/base.py:6537, in Index.slice_locs(self, start, end, step)
    6535 start_slice = None
    6536 if start is not None:
-> 6537     start_slice = self.get_slice_bound(start, "left")
    6538 if start_slice is None:
    6539     start_slice = 0

File ~/miniconda3/lib/python3.11/site-packages/pandas/core/indexes/base.py:6452, in Index.get_slice_bound(self, label, side)
    6448 original_label = label
    6450 # For datetime indices label may be a string that has to be converted
    6451 # to datetime boundary according to its resolution.
-> 6452 label = self._maybe_cast_slice_bound(label, side)
    6454 # we need to look up the label
    6455 try:

File ~/miniconda3/lib/python3.11/site-packages/pandas/core/indexes/base.py:6406, in Index._maybe_cast_slice_bound(self, label, side)
    6404 # reject them, if index does not contain label
    6405 if (is_float(label) or is_integer(label)) and label not in self:
-> 6406     self._raise_invalid_indexer("slice", label)
    6408 return label

File ~/miniconda3/lib/python3.11/site-packages/pandas/core/indexes/base.py:4152, in Index._raise_invalid_indexer(self, form, key, reraise)
    4150 if reraise is not lib.no_default:
    4151     raise TypeError(msg) from reraise
-> 4152 raise TypeError(msg)

TypeError: cannot do slice indexing on Index with these indexers [0] of type int
```

Reading csv files

```
In [40]: casts = pd.read_csv('cast.csv', index_col=None)
titles = pd.read_csv('titles.csv', index_col=None)
```

```
In [41]: casts.head()
```

Out[41]:

	title	year	name	type	character	n
0	Closet Monster	2015	Buffy #1	actor	Buffy 4	31.0
1	Suuri illusioni	1985	Homo \$	actor	Guests	22.0
2	Battle of the Sexes	2017	\$hutter	actor	Bobby Riggs Fan	10.0
3	Secret in Their Eyes	2015	\$hutter	actor	2002 Dodger Fan	NaN
4	Steve Jobs	2015	\$hutter	actor	1988 Opera House Patron	NaN

```
In [42]: casts.head(7)
```

Out[42]:

	title	year	name	type	character	n
0	Closet Monster	2015	Buffy #1	actor	Buffy 4	31.0
1	Suuri illusioni	1985	Homo \$	actor	Guests	22.0
2	Battle of the Sexes	2017	\$hutter	actor	Bobby Riggs Fan	10.0
3	Secret in Their Eyes	2015	\$hutter	actor	2002 Dodger Fan	NaN
4	Steve Jobs	2015	\$hutter	actor	1988 Opera House Patron	NaN
5	Straight Outta Compton	2015	\$hutter	actor	Club Patron	NaN
6	Straight Outta Compton	2015	\$hutter	actor	Dopeman	NaN

```
In [43]: titles.tail()
```

Out[43]:

	title	year
49995	Rebel	1970
49996	Suzanne	1996
49997	Bomba	2013
49998	Aao Jao Ghar Tumhara	1984
49999	Mrs. Munck	1995

Filtering data

- Movies after 1985
- Like select in relational algebra

```
In [44]: after85 = titles[titles['year'] > 1985]
after85
```

Out[44]:

	title	year
0	The Rising Son	1990
2	Crucea de piatra	1993
3	Country	2000
4	Gaiking II	2011
5	Medusa (IV)	2015
...
49990	Junebug	2005
49993	Corruption.Gov	2010
49996	Suzanne	1996
49997	Bomba	2013
49999	Mrs. Munck	1995

29814 rows × 2 columns

- Movies in years 1990 - 1999


```
In [45]: t = titles
movies90 = t[(t['year'] >= 1990) &(t['year'] < 2000)]
movies90
```

Out[45]:

	title	year
0	The Rising Son	1990
2	Crucea de piatra	1993
12	Poka Makorer Ghar Bosoti	1996
19	Maa Durga Shakti	1999
24	Conflict of Interest	1993
...
49969	Chi mei wang liang	1998
49979	Gagay: Prinsesa ng brownout	1993
49987	I Won't Dance	1992
49996	Suzanne	1996
49999	Mrs. Munck	1995

4803 rows × 2 columns

Sorting

- All movies named 'Macbeth'
- Sort by year

```
In [46]: macbeth = t[t['title'] == 'Macbeth']
```

```
In [47]: macbeth
```

Out[47]:

	title	year
4226	Macbeth	1913
9322	Macbeth	2006
11722	Macbeth	2013
17166	Macbeth	1997
25847	Macbeth	1998

```
In [48]: macbeth = macbeth.sort_values('year')
```

```
In [49]: macbeth
```

Out[49]:

	title	year
4226	Macbeth	1913
17166	Macbeth	1997
25847	Macbeth	1998
9322	Macbeth	2006
11722	Macbeth	2013

```
In [50]: macbeth = macbeth.sort_index()
```

```
In [51]: macbeth
```

Out[51]:

	title	year
4226	Macbeth	1913
9322	Macbeth	2006
11722	Macbeth	2013
17166	Macbeth	1997
25847	Macbeth	1998

Summaries and descriptive statistics

In [52]: titles.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 50000 entries, 0 to 49999
Data columns (total 2 columns):
#   Column  Non-Null Count  Dtype
---  ---      -
0   title   50000 non-null   object
1   year    50000 non-null   int64
dtypes: int64(1), object(1)
memory usage: 781.4+ KB
```

In [53]: casts.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 75001 entries, 0 to 75000
Data columns (total 6 columns):
#   Column  Non-Null Count  Dtype
---  ---      -
0   title   75000 non-null   object
1   year    75001 non-null   int64
2   name    75001 non-null   object
3   type    75001 non-null   object
4   character 75001 non-null   object
5   n       46035 non-null   float64
dtypes: float64(1), int64(1), object(4)
memory usage: 3.4+ MB
```

In [54]: titles.describe()

Out[54]:

	year
count	50000.000000
mean	1986.106120
std	29.293942
min	1900.000000
25%	1967.000000
50%	1996.000000
75%	2011.000000
max	2024.000000

In [55]: casts.describe()

Out[55]:

	year	n
count	75001.000000	46035.000000
mean	1990.536473	16.814359
std	26.748233	24.695616
min	1912.000000	1.000000
25%	1974.000000	4.000000
50%	2002.000000	10.000000
75%	2012.000000	21.000000
max	2023.000000	701.000000

Descriptive statistics for categorical data

In [56]: casts['name'].describe()

Out[56]:

```
count          75001
unique         29319
top            Ernie Adams
freq           431
Name: name, dtype: object
```

In [57]: housing = pd.read_csv('housing.csv', index_col=None)

In [58]: housing.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20640 entries, 0 to 20639
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  ---                ---
0   longitude              20640 non-null  float64
1   latitude               20640 non-null  float64
2   housing_median_age     20640 non-null  float64
3   total_rooms            20640 non-null  float64
4   total_bedrooms         20433 non-null  float64
5   population              20640 non-null  float64
6   households              20640 non-null  float64
7   median_income          20640 non-null  float64
8   median_house_value     20640 non-null  float64
9   ocean_proximity        20640 non-null  object
dtypes: float64(9), object(1)
memory usage: 1.6+ MB
```

In [59]: housing.describe()

Out[59]:

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_
count	20640.000000	20640.000000	20640.000000	20640.000000	20433.000000	20640.000000	20640.000000	20640.000000	;
mean	-119.569704	35.631861	28.639486	2635.763081	537.870553	1425.476744	499.539680	3.870671	20
std	2.003532	2.135952	12.585558	2181.615252	421.385070	1132.462122	382.329753	1.899822	1:
min	-124.350000	32.540000	1.000000	2.000000	1.000000	3.000000	1.000000	0.499900	:
25%	-121.800000	33.930000	18.000000	1447.750000	296.000000	787.000000	280.000000	2.563400	1:
50%	-118.490000	34.260000	29.000000	2127.000000	435.000000	1166.000000	409.000000	3.534800	1:
75%	-118.010000	37.710000	37.000000	3148.000000	647.000000	1725.000000	605.000000	4.743250	20
max	-114.310000	41.950000	52.000000	39320.000000	6445.000000	35682.000000	6082.000000	15.000100	50