

# Lecture 17, 17 October 2024

## 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)
s
```

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

```
In [3]: type(s)
```

```
Out[3]: pandas.core.series.Series
```

## 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)
ds
```

```
Out[4]: name      IBM
date    2010-09-08
shares      100
price      10.2
dtype: object
```

```
In [5]: type(ds)
```

```
Out[5]: pandas.core.series.Series
```

## Creating an index

- Provide a separate sequence of index headers, same length as values

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

```
fs
```

```
Out[6]: name          FB
       date      2001-08-02
       shares      90
       price       3.2
       dtype: object
```

## Accessing elements

- Using named index

```
In [7]: fs['shares']
```

```
Out[7]: 90
```

- Using position
  - Note: `fs[0]` also works but is *deprecated* -- may be disallowed in future versions

```
In [8]: fs.iloc[0], fs['name']
```

```
Out[8]: ('FB', 'FB')
```

- Using a slice of positions
  - Here the indices behave like positions in a list
  - Slice `[i:j]` runs from `i` to `j-1`

```
In [9]: fs.iloc[0:2]
```

```
Out[9]: name          FB
       date      2001-08-02
       dtype: object
```

- Using a sequence of positions

```
In [10]: fs.iloc[[0,2]]
```

```
Out[10]: name      FB
        shares     90
        dtype: object
```

- Can do the same with named indices
  - Slice uses position of indices
  - However, includes last index in the slice, unlike positional slice

```
In [11]: fs['name':'price']
```

```
Out[11]: name          FB
       date      2001-08-02
       shares      90
       price       3.2
       dtype: object
```

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

```
Out[12]: Series([], dtype: object)
```

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

```
Out[13]: price      3.2
        name        FB
        dtype: object
```

## Data frames

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

```
In [14]: 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)
df2
```

```
Out[14]:
```

	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 [15]: type(df2)
```

```
Out[15]: pandas.core.frame.DataFrame
```

- We can create a data frame from an anonymous sequence of sequences
  - In this case, each inner sequence is interpreted as a *row*, not a *column*
- Both rows and columns are indexed by position

```
In [16]: 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)
df3
```

```
Out[16]:
```

	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

- We can add a column
  - Provide a default value for all rows, or
  - Provide a sequence of values

```
In [17]: df2['owner'] = 'Unknown'
# df2['owner'] = ['a', 'b', 'c']
df2
```

```
Out[17]:
```

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

```
Out[18]:
```

	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

- If we provide a sequence of values, it must cover all rows

```
In [19]: #df2['owner'] = 'Unknown'
df2['owner3'] = ['a', 'b']
df2
```

```

-----
ValueError                                Traceback (most recent call last)
Cell In[19], line 2
      1 #df2['owner'] = 'Unknown'
----> 2 df2['owner3'] = ['a', 'b']
      3 df2

File ~/python-venv/venv/lib/python3.11/site-packages/pandas/core/frame.py:4311, in DataFrame._setitem__(self, key, value)
    4308     self._setitem_array([[key], value])
    4309 else:
    4310     # set column
-> 4311     self._set_item(key, value)

File ~/python-venv/venv/lib/python3.11/site-packages/pandas/core/frame.py:4524, in DataFrame._set_item(self, key, value)
    4514 def _set_item(self, key, value) -> None:
    4515     """
    4516     Add series to DataFrame in specified column.
    4517     (...)
    4522     ensure homogeneity.
    4523     """
-> 4524     value, refs = self._sanitize_column(value)
    4526     if (
    4527         key in self.columns
    4528         and value.ndim == 1
    4529         and not isinstance(value.dtype, ExtensionDtype)
    4530     ):
    4531         # broadcast across multiple columns if necessary
    4532         if not self.columns.is_unique or isinstance(self.columns, MultiIndex):

File ~/python-venv/venv/lib/python3.11/site-packages/pandas/core/frame.py:5266, in DataFrame._sanitize_column(self, value)
    5263     return _reindex_for_setitem(value, self.index)
    5265 if is_list_like(value):
-> 5266     com.require_length_match(value, self.index)
    5267 arr = sanitize_array(value, self.index, copy=True, allow_2d=True)
    5268 if (
    5269     isinstance(value, Index)
    5270     and value.dtype == "object"
    (...)
    5273     # TODO: Remove kludge in sanitize_array for string mode when enforcing
    5274     # this deprecation

File ~/python-venv/venv/lib/python3.11/site-packages/pandas/core/common.py:573, in require_length_match(data, index)
    569     """
    570     Check the length of data matches the length of the index.
    571     """
    572     if len(data) != len(index):
-> 573         raise ValueError(
    574             "Length of values "
    575             f"({len(data)}) "
    576             "does not match length of index "
    577             f"({len(index)})"
    578         )

ValueError: Length of values (2) does not match length of index (3)

```

## Add row indices

```
In [20]: df2.index = ['one', 'two', 'three']
df2
```

```
Out[20]:
```

	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

- Note that we lose the previous indices (one, two, three)

```
In [21]: df2 = df2.set_index(['name'])
df2
```

```
Out[21]:
```

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

## Replace an index

- Again, we lose the previous index, `name`

```
In [22]: df2 = df2.set_index(['price'])
df2
```

```
Out[22]:
```

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

## Use multiple columns for indexing

```
In [23]: df2 = pd.DataFrame(data2) # Reset data frame to original
df2['owner'] = 'Unknown'
df2 = df2.set_index(['name', 'price'])
df2
```

```
Out[23]:
```

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

## Accessing values in a dataframe

### By column index

- Similar to projection in relational algebra

```
In [24]: df2[['shares', 'date']]
```

```
Out[24]:
```

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

### By row index

```
In [25]: df2.loc['AA']
```

```
Out[25]:
```

	date	shares	owner
	<b>price</b>		
<b>12.3</b>	2001-12-01	100	Unknown

## Individual element - specify row and column index

```
In [26]: df2.loc['AA', 'shares']
```

```
Out[26]: price
12.3 100
Name: shares, dtype: int64
```

## Slices, etc

```
In [27]: df2.loc[:, 'shares'] # All rows, column 'shares'
```

```
Out[27]: name price
AA 12.3 100
IBM 10.3 30
GOOG 32.2 90
Name: shares, dtype: int64
```

```
In [28]: df2 = pd.DataFrame(data2) # Reset data frame to original
df2['owner'] = 'Unknown'
df2 = df2.set_index(['name'])
df2
```

```
Out[28]:
```

	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

- An arbitrary subtable of rows and columns

```
In [29]: df2.loc['AA':'IBM', 'shares':'owner']
```

```
Out[29]:
```

	shares	price	owner
name			
AA	100	12.3	Unknown
IBM	30	10.3	Unknown

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

```
In [30]: df2 = pd.DataFrame(data2) # Reset data frame to original
df2['owner'] = 'Unknown'
df2
```

```
Out[30]:
```

	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

- We can slice the rows
  - Note that for data frames, 0, 1, ... are treated as labels, so slice works like for named indices

```
In [31]: df2.loc[0:1]
```

```
Out[31]:
```

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

- Now create a row index

```
In [32]: df2 = df2.set_index(['name'])
df2
```

Out [32]:

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

- Can no longer slice rows by position
  - This is because the row numbers are really labels, not positions
  - The behaviour for `Series` is inconsistent with this interpretation for `DataFrame`

In [33]: `df2.loc[0:2]`

```

-----
TypeError                                 Traceback (most recent call last)
Cell In[33], line 1
----> 1 df2.loc[0:2]

File ~/python-venv/venv/lib/python3.11/site-packages/pandas/core/indexing.py:1191, in _LocationIndexer._getitem
__(self, key)
    1189 maybe_callable = com.apply_if_callable(key, self.obj)
    1190 maybe_callable = self._check_deprecated_callable_usage(key, maybe_callable)
-> 1191 return self._getitem_axis(maybe_callable, axis=axis)

File ~/python-venv/venv/lib/python3.11/site-packages/pandas/core/indexing.py:1411, in _iLocIndexer._getitem_axis
(self, key, axis)
    1409 if isinstance(key, slice):
    1410     self._validate_key(key, axis)
-> 1411     return self._get_slice_axis(key, axis=axis)
    1412 elif com.is_bool_indexer(key):
    1413     return self._getbool_axis(key, axis=axis)

File ~/python-venv/venv/lib/python3.11/site-packages/pandas/core/indexing.py:1443, in _iLocIndexer._get_slice_axi
s(self, slice_obj, axis)
    1440 return obj.copy(deep=False)
    1442 labels = obj._get_axis(axis)
-> 1443 indexer = labels.slice_indexer(slice_obj.start, slice_obj.stop, slice_obj.step)
    1445 if isinstance(indexer, slice):
    1446     return self.obj._slice(indexer, axis=axis)

File ~/python-venv/venv/lib/python3.11/site-packages/pandas/core/indexes/base.py:6662, in Index.slice_indexer(se
lf, start, end, step)
    6618 def slice_indexer(
    6619     self,
    6620     start: Hashable | None = None,
    6621     end: Hashable | None = None,
    6622     step: int | None = None,
    6623 ) -> slice:
    6624     """
    6625     Compute the slice indexer for input labels and step.
    6626     (...)
    6660     slice(1, 3, None)
    6661     """
-> 6662 start_slice, end_slice = self.slice_locs(start, end, step=step)
    6664 # return a slice
    6665 if not is_scalar(start_slice):

File ~/python-venv/venv/lib/python3.11/site-packages/pandas/core/indexes/base.py:6879, in Index.slice_locs(self,
start, end, step)
    6877 start_slice = None
    6878 if start is not None:
-> 6879     start_slice = self.get_slice_bound(start, "left")
    6880 if start_slice is None:
    6881     start_slice = 0

File ~/python-venv/venv/lib/python3.11/site-packages/pandas/core/indexes/base.py:6794, in Index.get_slice_bound
(self, label, side)
    6790 original_label = label
    6792 # For datetime indices label may be a string that has to be converted
    6793 # to datetime boundary according to its resolution.
-> 6794 label = self._maybe_cast_slice_bound(label, side)
    6796 # we need to look up the label
    6797 try:

File ~/python-venv/venv/lib/python3.11/site-packages/pandas/core/indexes/base.py:6727, in Index._maybe_cast_slic
e_bound(self, label, side)
    6725 # reject them, if index does not contain label
    6726 if (is_float(label) or is_integer(label)) and label not in self:
-> 6727     self._raise_invalid_indexer("slice", label)
    6729 return label

File ~/python-venv/venv/lib/python3.11/site-packages/pandas/core/indexes/base.py:4301, in Index._raise_invalid_i
ndexer(self, form, key, reraise)
    4299 if reraise is not lib.no_default:
    4300     raise TypeError(msg) from reraise
-> 4301 raise TypeError(msg)

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

```

## Reading csv files

- By convention, the first line of a csv file is interpreted to be column names
- `index_col=None` says do not create a row index from any of the given columns



```
In [34]: casts = pd.read_csv('cast.csv', index_col=None)
titles = pd.read_csv('titles.csv', index_col=None)
```

- Examine the first few rows using `head()`
  - Default is 5 rows
  - Can ask for `n` rows

```
In [35]: casts.head()
```

```
Out[35]:
```

	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 [36]: casts.head(7)
```

```
Out[36]:
```

	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

- Likewise, `tail()` gives last few lines

```
In [37]: titles.tail()
```

```
Out[37]:
```

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

```
In [38]: titles.tail(8)
```

```
Out[38]:
```

	title	year
49992	Legend of Horror	1972
49993	Corruption.Gov	2010
49994	Lille Fridolf blir morfar	1957
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 [39]: after85 = titles[titles['year'] > 1985]
# select * from titles where year > 1985
after85
```

```
Out[39]:
```

	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

- Project the output of select on column `title`

```
In [40]: after85titles = titles[titles['year'] > 1985]['title']
after85titles
```

```
Out[40]:
```

0	The Rising Son
2	Crucea de piatra
3	Country
4	Gaiking II
5	Medusa (IV)
...	...
49990	Junebug
49993	Corruption.Gov
49996	Suzanne
49997	Bomba
49999	Mrs. Munck

Name: title, Length: 29814, dtype: object

- Boolean combinations of conditions
  - `&` for and, `|` for or, `~` for not
- Movies in years 1990 - 1999

```
In [41]: t = titles
movies90 = t[(t['year'] >= 1990) & (t['year'] < 2000)]
movies90
```

```
Out[41]:
```

	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

- Complement of the previous condition, using negation

```
In [42]: t = titles
notmovies90 = t[~((t['year'] >= 1990) & (t['year'] < 2000))]
notmovies90
```

```
Out[42]:
```

	title	year
1	The Thousand Plane Raid	1969
3	Country	2000
4	Gaiking II	2011
5	Medusa (IV)	2015
6	The Fresh Air Will Do You Good	2008
...	...	...
49993	Corruption.Gov	2010
49994	Lille Fridolf blir morfar	1957
49995	Rebel	1970
49997	Bomba	2013
49998	Aao Jao Ghar Tumhara	1984

45197 rows × 2 columns

- Complement of the previous condition, using or

```
In [43]: t = titles
notmovies90or = t[(t['year'] < 1990) | (t['year'] >= 2000)]
notmovies90or
```

```
Out[43]:
```

	title	year
1	The Thousand Plane Raid	1969
3	Country	2000
4	Gaiking II	2011
5	Medusa (IV)	2015
6	The Fresh Air Will Do You Good	2008
...	...	...
49993	Corruption.Gov	2010
49994	Lille Fridolf blir morfar	1957
49995	Rebel	1970
49997	Bomba	2013
49998	Aao Jao Ghar Tumhara	1984

45197 rows × 2 columns

## Sorting

- All movies named 'Macbeth'

```
In [44]: macbeth = t[t['title'] == 'Macbeth']
macbeth
```

```
Out[44]:
```

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

- Sort by year
  - Note that sort is in-place

```
In [45]: macbeth = macbeth.sort_values('year')
macbeth
```

```
Out[45]:
```

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

- To restore original order, sort by index

```
In [46]: macbeth = macbeth.sort_index()
macbeth
```

```
Out[46]:
```

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

## Summaries and descriptive statistics

- `info()` gives overall summary of a data frame

```
In [47]: 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 [48]: 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
```

- `describe()` gives statistical summary of numeric columns

```
In [49]: titles.describe()
```

```
Out[49]:
```

	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 [50]: casts.describe()
```

```
Out[50]:
```

	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

- Can also get summary for a non-numeric column

```
In [51]: casts['name'].describe()
```

```
Out[51]:
```

count	75001
unique	29319
top	Ernie Adams
freq	431

Name: name, dtype: object

```
In [52]: casts['character'].describe()
```

```
Out[52]:
```

count	75001
unique	50299
top	Himself
freq	405

Name: character, dtype: object

- Another example, housing data by locality in California

```
In [53]: housing = pd.read_csv('housing.csv', index_col=None)
```

```
In [54]: housing.head()
```

```
Out[54]:
```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_hc
0	-122.23	37.88	41.0	880.0	129.0	322.0	126.0	8.3252	
1	-122.22	37.86	21.0	7099.0	1106.0	2401.0	1138.0	8.3014	
2	-122.24	37.85	52.0	1467.0	190.0	496.0	177.0	7.2574	
3	-122.25	37.85	52.0	1274.0	235.0	558.0	219.0	5.6431	
4	-122.25	37.85	52.0	1627.0	280.0	565.0	259.0	3.8462	

```
In [55]: 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 [56]: housing.describe()
```

```
Out[56]:
```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_ir
<b>count</b>	20640.000000	20640.000000	20640.000000	20640.000000	20433.000000	20640.000000	20640.000000	20640.0
<b>mean</b>	-119.569704	35.631861	28.639486	2635.763081	537.870553	1425.476744	499.539680	3.8
<b>std</b>	2.003532	2.135952	12.585558	2181.615252	421.385070	1132.462122	382.329753	1.8
<b>min</b>	-124.350000	32.540000	1.000000	2.000000	1.000000	3.000000	1.000000	0.4
<b>25%</b>	-121.800000	33.930000	18.000000	1447.750000	296.000000	787.000000	280.000000	2.5
<b>50%</b>	-118.490000	34.260000	29.000000	2127.000000	435.000000	1166.000000	409.000000	3.5
<b>75%</b>	-118.010000	37.710000	37.000000	3148.000000	647.000000	1725.000000	605.000000	4.7
<b>max</b>	-114.310000	41.950000	52.000000	39320.000000	6445.000000	35682.000000	6082.000000	15.0

- Only one non-numeric column, `ocean_proximity`

```
In [57]: housing['ocean_proximity'].describe()
```

```
Out[57]: count      20640
unique         5
top    <1H OCEAN
freq         9136
Name: ocean_proximity, dtype: object
```