

# Python Lab Exercise #2

## Objectives:

- Load .csv files into pandas DataFrames
- Describe and manipulate data in Series and DataFrames
- Visualize data using DataFrame methods and matplotlib



In [103...]

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from scipy import stats
from sklearn.datasets import load_iris
```

## What is Pandas?

Pandas, as [the Anaconda docs](#) tell us, offers us "High-performance, easy-to-use data structures and data analysis tools." It's something like "Excel for Python", but it's quite a bit more powerful.

Let's read in the heart dataset.

Pandas has many methods for reading different types of files. Note that here we have a .csv file.

Read about this dataset [here](#).

In [65]:

```
heart_df = pd.read_csv('heart.csv')
```

The output of the `.read_csv()` function is a pandas *DataFrame*, which has a familiar tabular structure of rows and columns.

In [66]:

```
type(heart_df)
```

Out[66]: pandas.core.frame.DataFrame

In [67]: heart\_df.head()

Out[67]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
<b>0</b>	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
<b>1</b>	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
<b>2</b>	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
<b>3</b>	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
<b>4</b>	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1

In [68]: heart\_df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 303 entries, 0 to 302
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype  
 --- 
  0   age         303 non-null    int64  
  1   sex         303 non-null    int64  
  2   cp          303 non-null    int64  
  3   trestbps    303 non-null    int64  
  4   chol        303 non-null    int64  
  5   fbs         303 non-null    int64  
  6   restecg     303 non-null    int64  
  7   thalach     303 non-null    int64  
  8   exang       303 non-null    int64  
  9   oldpeak     303 non-null    float64 
  10  slope        303 non-null    int64  
  11  ca          303 non-null    int64  
  12  thal        303 non-null    int64  
  13  target       303 non-null    int64  
dtypes: float64(1), int64(13)
memory usage: 33.3 KB
```

In [69]: heart\_df.describe().T

Out[69]:

	count	mean	std	min	25%	50%	75%	max
<b>age</b>	303.0	54.366337	9.082101	29.0	47.5	55.0	61.0	77.0
<b>sex</b>	303.0	0.683168	0.466011	0.0	0.0	1.0	1.0	1.0
<b>cp</b>	303.0	0.966997	1.032052	0.0	0.0	1.0	2.0	3.0
<b>trestbps</b>	303.0	131.623762	17.538143	94.0	120.0	130.0	140.0	200.0
<b>chol</b>	303.0	246.264026	51.830751	126.0	211.0	240.0	274.5	564.0
<b>fbs</b>	303.0	0.148515	0.356198	0.0	0.0	0.0	0.0	1.0
<b>restecg</b>	303.0	0.528053	0.525860	0.0	0.0	1.0	1.0	2.0
<b>thalach</b>	303.0	149.646865	22.905161	71.0	133.5	153.0	166.0	202.0

	count	mean	std	min	25%	50%	75%	max
<b>exang</b>	303.0	0.326733	0.469794	0.0	0.0	0.0	1.0	1.0
<b>oldpeak</b>	303.0	1.039604	1.161075	0.0	0.0	0.8	1.6	6.2
<b>slope</b>	303.0	1.399340	0.616226	0.0	1.0	1.0	2.0	2.0
<b>ca</b>	303.0	0.729373	1.022606	0.0	0.0	0.0	1.0	4.0
<b>thal</b>	303.0	2.313531	0.612277	0.0	2.0	2.0	3.0	3.0
<b>target</b>	303.0	0.544554	0.498835	0.0	0.0	1.0	1.0	1.0

In [70]:

heart\_df

Out[70]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
<b>0</b>	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
<b>1</b>	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
<b>2</b>	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
<b>3</b>	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
<b>4</b>	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
<b>298</b>	57	0	0	140	241	0	1	123	1	0.2	1	0	3	0
<b>299</b>	45	1	3	110	264	0	1	132	0	1.2	1	0	3	0
<b>300</b>	68	1	0	144	193	1	1	141	0	3.4	1	2	3	0
<b>301</b>	57	1	0	130	131	0	1	115	1	1.2	1	1	3	0
<b>302</b>	57	0	1	130	236	0	0	174	0	0.0	1	1	2	0

303 rows × 14 columns

## DataFrames and Series

Two main types of pandas objects are the DataFrame and the Series, the latter being in effect a single column of the former:

In [71]:

```
age_series = heart_df['age']
type(age_series)
```

Out[71]:

pandas.core.series.Series

Notice how we can isolate a column of our DataFrame simply by using square brackets together with the name of the column.

Both Series and DataFrames have an *index* as well:

In [72]:

heart\_df.index

```
Out[72]: RangeIndex(start=0, stop=303, step=1)
```

```
In [73]: age_series.index
```

```
Out[73]: RangeIndex(start=0, stop=303, step=1)
```

Pandas is built on top of NumPy, and we can always access the NumPy array underlying a DataFrame using `.values`.

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

```
Out[74]: array([[63.,  1.,  3., ...,  0.,  1.,  1.],
   [37.,  1.,  2., ...,  0.,  2.,  1.],
   [41.,  0.,  1., ...,  0.,  2.,  1.],
   ...,
   [68.,  1.,  0., ...,  2.,  3.,  0.],
   [57.,  1.,  0., ...,  1.,  3.,  0.],
   [57.,  0.,  1., ...,  1.,  2.,  0.]])
```

## Basic DataFrame Attributes and Methods

### `.head()`

```
In [75]: heart_df.head()
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
<b>0</b>	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
<b>1</b>	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
<b>2</b>	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
<b>3</b>	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
<b>4</b>	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1

### `.tail()`

```
In [76]: heart_df.tail()
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
<b>298</b>	57	0	0	140	241	0	1	123	1	0.2	1	0	3	0
<b>299</b>	45	1	3	110	264	0	1	132	0	1.2	1	0	3	0
<b>300</b>	68	1	0	144	193	1	1	141	0	3.4	1	2	3	0
<b>301</b>	57	1	0	130	131	0	1	115	1	1.2	1	1	3	0
<b>302</b>	57	0	1	130	236	0	0	174	0	0.0	1	1	2	0

## .info()

In [77]:

```
heart_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 303 entries, 0 to 302
Data columns (total 14 columns):
 #   Column   Non-Null Count  Dtype  
 --- 
  0   age      303 non-null    int64  
  1   sex      303 non-null    int64  
  2   cp       303 non-null    int64  
  3   trestbps 303 non-null    int64  
  4   chol     303 non-null    int64  
  5   fbs      303 non-null    int64  
  6   restecg  303 non-null    int64  
  7   thalach  303 non-null    int64  
  8   exang    303 non-null    int64  
  9   oldpeak  303 non-null    float64 
  10  slope    303 non-null    int64  
  11  ca       303 non-null    int64  
  12  thal     303 non-null    int64  
  13  target   303 non-null    int64  
dtypes: float64(1), int64(13)
memory usage: 33.3 KB
```

## .describe()

In [78]:

```
heart_df.describe()
```

Out[78]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach
<b>count</b>	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000
<b>mean</b>	54.366337	0.683168	0.966997	131.623762	246.264026	0.148515	0.528053	149.646865
<b>std</b>	9.082101	0.466011	1.032052	17.538143	51.830751	0.356198	0.525860	22.905161
<b>min</b>	29.000000	0.000000	0.000000	94.000000	126.000000	0.000000	0.000000	71.000000
<b>25%</b>	47.500000	0.000000	0.000000	120.000000	211.000000	0.000000	0.000000	133.500000
<b>50%</b>	55.000000	1.000000	1.000000	130.000000	240.000000	0.000000	1.000000	153.000000
<b>75%</b>	61.000000	1.000000	2.000000	140.000000	274.500000	0.000000	1.000000	166.000000
<b>max</b>	77.000000	1.000000	3.000000	200.000000	564.000000	1.000000	2.000000	202.000000

## .dtypes

In [79]:

```
heart_df.dtypes
```

Out[79]:

age	int64
sex	int64

```

cp           int64
trestbps    int64
chol          int64
fbs           int64
restecg      int64
thalach       int64
exang          int64
oldpeak     float64
slope          int64
ca             int64
thal           int64
target         int64
dtype: object

```

## .shape

In [80]: `heart_df.shape`

Out[80]: (303, 14)

## Exploratory Plots

Let's make ourselves a histogram of ages:

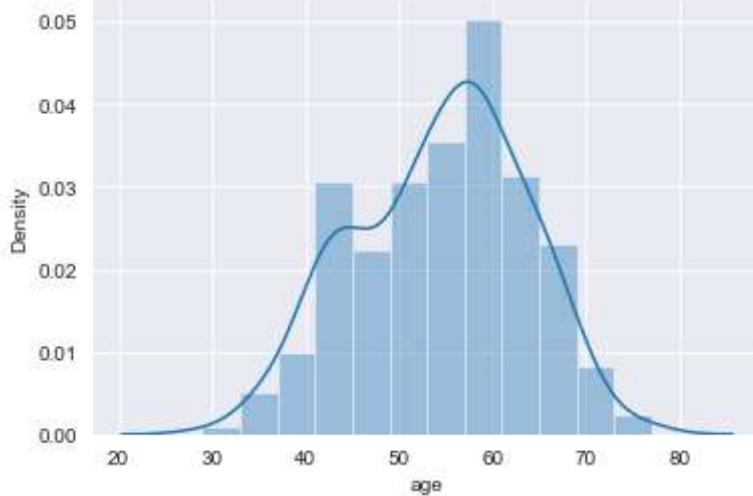
```

In [81]: sns.set_style('darkgrid')
sns.distplot(a=heart_df['age']);
# For more recent versions of seaborn:
# sns.histplot(data=heart_df['age'], kde=True);

```

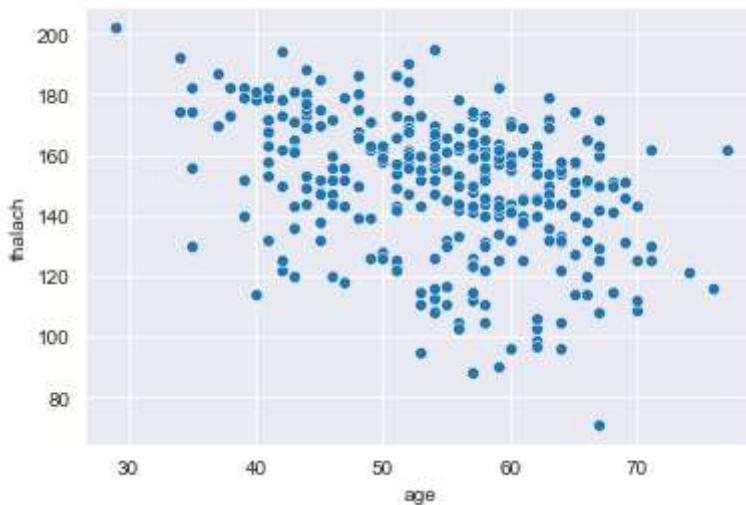
C:\Users\Aungkyaw\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

`warnings.warn(msg, FutureWarning)`



And while we're at it let's do a scatter plot of maximum heart rate vs. age:

```
In [82]: sns.scatterplot(x=heart_df['age'], y=heart_df['thalach']);
```



## Adding to a DataFrame

### Adding Rows

Here are two rows that our engineer accidentally left out of the .csv file, expressed as a Python dictionary:

```
In [83]: extra_rows = {'age': [40, 30], 'sex': [1, 0], 'cp': [0, 0], 'trestbps': [120, 130],
                  'chol': [240, 200],
                  'fbs': [0, 0], 'restecg': [1, 0], 'thalach': [120, 122], 'exang': [0, 1],
                  'oldpeak': [0.1, 1.0], 'slope': [1, 1], 'ca': [0, 1], 'thal': [2, 3],
                  'target': [0, 0]}
extra_rows
```

```
Out[83]: {'age': [40, 30],
           'sex': [1, 0],
           'cp': [0, 0],
           'trestbps': [120, 130],
           'chol': [240, 200],
           'fbs': [0, 0],
           'restecg': [1, 0],
           'thalach': [120, 122],
           'exang': [0, 1],
           'oldpeak': [0.1, 1.0],
           'slope': [1, 1],
           'ca': [0, 1],
           'thal': [2, 3],
           'target': [0, 0]}
```

How can we add this to the bottom of our dataset?

```
In [84]: # Let's first turn this into a DataFrame.
          # We can use the .from_dict() method.

missing = pd.DataFrame(extra_rows)
missing
```

```
Out[84]:   age  sex  cp  trestbps  chol  fbs  restecg  thalach  exang  oldpeak  slope  ca  thal  target
```

	age	sex	cp	trestbps	chol	fb	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	40	1	0	120	240	0	1	120	0	0.1	1	0	2	0
1	30	0	0	130	200	0	0	122	1	1.0	1	1	3	0

In [85]:

```
# Now we just need to concatenate the two DataFrames together.
# Note the `ignore_index` parameter! We'll set that to True.
```

```
heart_augmented = pd.concat([heart_df, missing],
                             ignore_index=True)
```

In [86]:

```
# Let's check the end to make sure we were successful!
```

```
heart_augmented.tail()
```

Out[86]:

	age	sex	cp	trestbps	chol	fb	restecg	thalach	exang	oldpeak	slope	ca	thal	target
300	68	1	0	144	193	1	1	141	0	3.4	1	2	3	0
301	57	1	0	130	131	0	1	115	1	1.2	1	1	3	0
302	57	0	1	130	236	0	0	174	0	0.0	1	1	2	0
303	40	1	0	120	240	0	1	120	0	0.1	1	0	2	0
304	30	0	0	130	200	0	0	122	1	1.0	1	1	3	0

## Adding Columns

Adding a column is very easy in `pandas`. Let's add a new column to our dataset called "test", and set all of its values to 0.

In [87]:

```
heart_augmented['test'] = 0
```

In [88]:

```
heart_augmented.head()
```

Out[88]:

	age	sex	cp	trestbps	chol	fb	restecg	thalach	exang	oldpeak	slope	ca	thal	target	test
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1	0
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1	0
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1	0
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1	0
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1	0

I can also add columns whose values are functions of existing columns.

Suppose I want to add the cholesterol column ("chol") to the resting systolic blood pressure column ("trestbps"):

In [89]: `heart_augmented['chol+trestbps'] = heart_augmented['chol'] + heart_augmented['trestbps']`

In [90]: `heart_augmented.head()`

Out[90]:

	age	sex	cp	trestbps	chol	fb	restecg	thalach	exang	oldpeak	slope	ca	thal	target	test
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1	0
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1	0
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1	0
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1	0
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1	0

## Filtering

We can use filtering techniques to see only certain rows of our data. If we wanted to see only the rows for patients 70 years of age or older, we can simply type:

In [91]: `heart_augmented['age'] >= 70`

Out[91]:

0	False
1	False
2	False
3	False
4	False
...	
300	False
301	False
302	False
303	False
304	False

Name: age, Length: 305, dtype: bool

In [92]: `heart_augmented[heart_augmented['age'] >= 70]`

Out[92]:

	age	sex	cp	trestbps	chol	fb	restecg	thalach	exang	oldpeak	slope	ca	thal	target	test
25	71	0	1	160	302	0	1	162	0	0.4	2	2	2	1	C
60	71	0	2	110	265	1	0	130	0	0.0	2	1	2	1	C
129	74	0	1	120	269	0	0	121	1	0.2	2	1	2	1	C
144	76	0	2	140	197	0	2	116	0	1.1	1	0	2	1	C
145	70	1	1	156	245	0	0	143	0	0.0	2	0	2	1	C
151	71	0	0	112	149	0	1	125	0	1.6	1	0	2	1	C
225	70	1	0	145	174	0	1	125	1	2.6	0	0	3	0	C
234	70	1	0	130	322	0	0	109	0	2.4	1	3	2	0	C

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target	test
238	77	1	0	125	304	0	0	162	1	0.0	2	3	2	0	C
240	70	1	2	160	269	0	1	112	1	2.9	1	1	3	0	C

Use '&' for "and" and '|' for "or".

## Exercise

Display the patients who are 70 or over as well as the patients whose trestbps score is greater than 170.

In [95]:

```
heart_augmented.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 305 entries, 0 to 304
Data columns (total 16 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   age              305 non-null    int64  
 1   sex              305 non-null    int64  
 2   cp               305 non-null    int64  
 3   trestbps         305 non-null    int64  
 4   chol             305 non-null    int64  
 5   fbs              305 non-null    int64  
 6   restecg          305 non-null    int64  
 7   thalach          305 non-null    int64  
 8   exang            305 non-null    int64  
 9   oldpeak          305 non-null    float64 
 10  slope            305 non-null    int64  
 11  ca               305 non-null    int64  
 12  thal              305 non-null    int64  
 13  target            305 non-null    int64  
 14  test              305 non-null    int64  
 15  chol+trestbps    305 non-null    int64  
dtypes: float64(1), int64(15)
memory usage: 38.2 KB
```

In [102...]:

```
df_filtered = heart_augmented[(heart_augmented['age'] > 70) & (heart_augmented['trestbps'] > 170)]
df_filtered
```

Out[102...]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target	test
25	71	0	1	160	302	0	1	162	0	0.4	2	2	2	1	C

In [98]:

```
heart_augmented[heart_augmented['age'] > 70]
```

Out[98]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target	test
25	71	0	1	160	302	0	1	162	0	0.4	2	2	2	1	C
60	71	0	2	110	265	1	0	130	0	0.0	2	1	2	1	C

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target	test
129	74	0	1	120	269	0	0	121	1	0.2	2	1	2	1	C
144	76	0	2	140	197	0	2	116	0	1.1	1	0	2	1	C
151	71	0	0	112	149	0	1	125	0	1.6	1	0	2	1	C
238	77	1	0	125	304	0	0	162	1	0.0	2	3	2	0	C

## Exploratory Plot

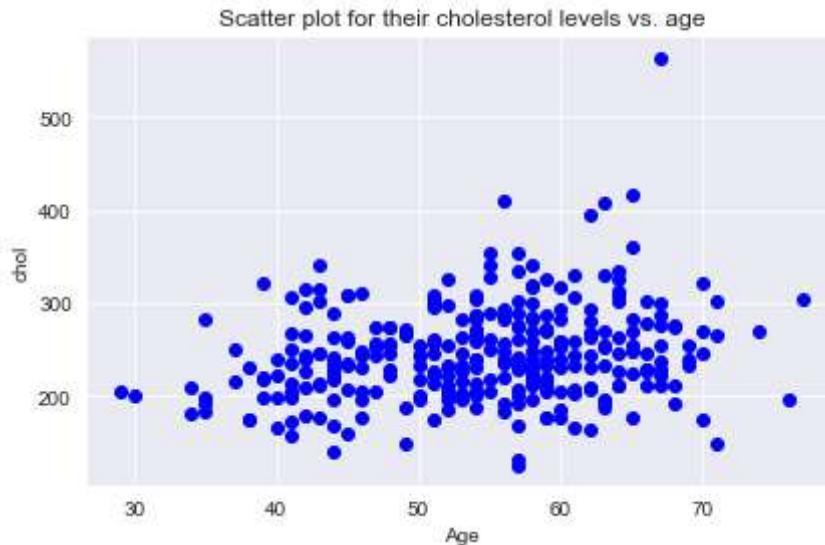
Using the subframe we just made, let's make a scatter plot of their cholesterol levels vs. age and color by sex:

In [111...]

```
fig, ax = plt.subplots()

ax.scatter(heart_augmented['age'], heart_augmented['chol'], c='blue')
ax.set_xlabel('Age')
ax.set_ylabel('chol')
ax.set_title('Scatter plot for their cholesterol levels vs. age')

plt.tight_layout()
```



## .loc and .iloc

We can use `.loc` to get, say, the first ten values of the age and resting blood pressure ("trestbps") columns:

In [112...]

```
heart_augmented.loc
```

Out[112...]

```
<pandas.core.indexing._LocIndexer at 0x24f450195e0>
```

In [113...]

```
heart_augmented.loc[:9, ['age', 'trestbps']]
```

Out[113...]

	age	trestbps
0	63	145
1	37	130
2	41	130
3	56	120
4	57	120
5	57	140
6	56	140
7	44	120
8	52	172
9	57	150

.iloc is used for selecting locations in the DataFrame **by number**:

In [ ]:

heart\_augmented.iloc

In [114...]

heart\_augmented.iloc[3, 0]

Out[114...]

56

In [115...]

heart\_augmented.head()

Out[115...]

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target	test
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1	0
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1	0
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1	0
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1	0
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1	0

## Exercise

How would we get the same slice as just above by using .iloc() instead of .loc()?

In [133...]

heart\_augmented.iloc[0:10:, [0, 3]]

Out[133...]

	age	trestbps
0	63	145
1	37	130

	age	trestbps
<b>2</b>	41	130
<b>3</b>	56	120
<b>4</b>	57	120
<b>5</b>	57	140
<b>6</b>	56	140
<b>7</b>	44	120
<b>8</b>	52	172
<b>9</b>	57	150

## Statistics

### .mean()

In [134...]

```
heart_augmented.mean()
```

Out[134...]

age	54.239344
sex	0.681967
cp	0.960656
trestbps	131.580328
chol	246.091803
fbs	0.147541
restecg	0.527869
thalach	149.459016
exang	0.327869
oldpeak	1.036393
slope	1.396721
ca	0.727869
thal	2.314754
target	0.540984
test	0.000000
chol+trestbps	377.672131
dtype:	float64

Be careful! Some of these will be not straightforwardly interpretable. What does an average "sex" of 0.682 mean?

### .min()

In [135...]

```
heart_augmented.min()
```

Out[135...]

age	29.0
sex	0.0
cp	0.0
trestbps	94.0
chol	126.0
fbs	0.0
restecg	0.0

```
thalach      71.0
exang        0.0
oldpeak      0.0
slope         0.0
ca            0.0
thal          0.0
target        0.0
test          0.0
chol+trestbps 249.0
dtype: float64
```

### .max()

```
In [136...]: heart_augmented.max()
```

```
Out[136...]: age        77.0
sex         1.0
cp          3.0
trestbps   200.0
chol       564.0
fbs         1.0
restecg    2.0
thalach    202.0
exang      1.0
oldpeak    6.2
slope       2.0
ca          4.0
thal        3.0
target     1.0
test        0.0
chol+trestbps 679.0
dtype: float64
```

## Series Methods

### .value\_counts()

How many different values does slope have? What about sex? And target?

```
In [137...]: heart_augmented['slope'].value_counts()
```

```
Out[137...]: 2    142
1    142
0     21
Name: slope, dtype: int64
```

```
In [138...]: heart_augmented['sex'].value_counts()
```

```
Out[138...]: 1    208
0     97
Name: sex, dtype: int64
```

### .sort\_values()

In [139...]: `heart_augmented['age'].sort_values()`

Out[139...]:

72	29
304	30
58	34
125	34
65	35
	..
25	71
60	71
129	74
144	76
238	77

Name: age, Length: 305, dtype: int64

## pandas -Native Plotting

The `.plot()` and `.hist()` methods available for DataFrames use a wrapper around `matplotlib`:

In [ ]: `heart_augmented.plot(x='age', y='trestbps', kind='scatter');`

In [ ]: `heart_augmented.hist(column='chol');`

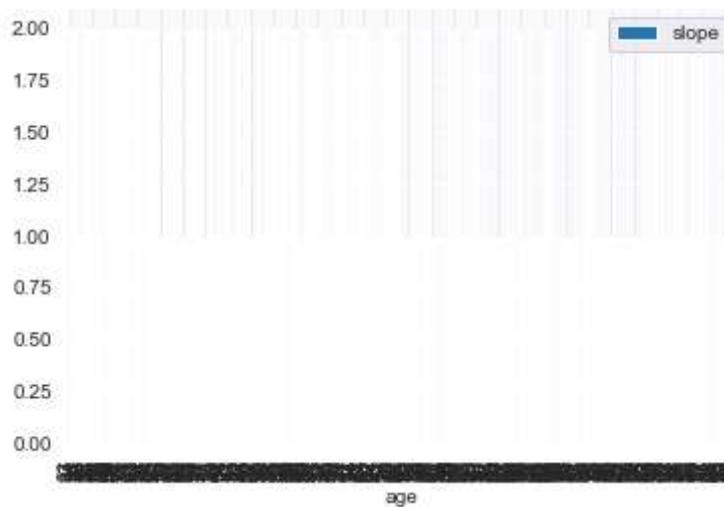
## Exercises

1. Make a bar plot of "age" vs. "slope" for the `heart_augmented` DataFrame.

In [144...]:

```
# Enter your code s
heart_augmented.plot.bar(x='age',y='slope')
```

Out[144...]:

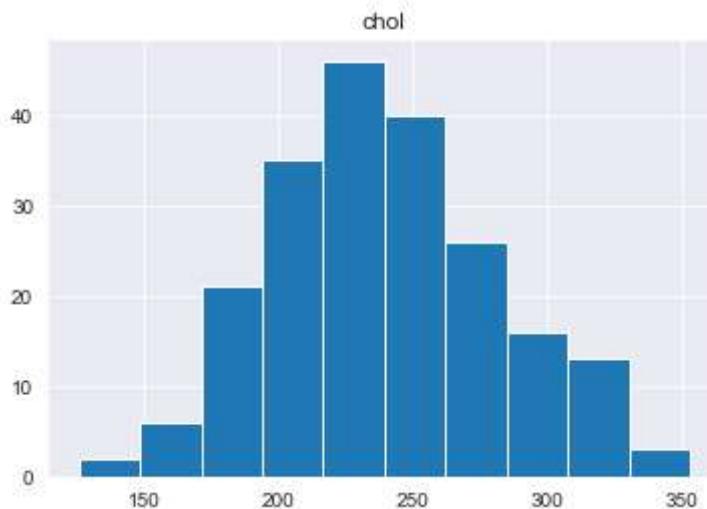


1. Make a histogram of ages for **just the men** in `heart_augmented` (`heart_augmented['sex']=1`).

In [149...]:

```
# Enter your code here
men=heart_augmented[heart_augmented['sex']==1]
men
men.hist(column='chol')
```

Out[149... array([ [`<AxesSubplot:title={'center':'chol'}`] ]], dtype=object)



1. Make separate scatter plots of cholesterol vs. resting systolic blood pressure for the target=0 and the target=1 groups. Put both plots on the same figure and give each an appropriate title.

In [169...]

```
heart_augmented['target'].value_counts()
target0=heart_augmented[heart_augmented['target']==0]
target1=heart_augmented[heart_augmented['target']==1]
heart_augmented.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 305 entries, 0 to 304
Data columns (total 16 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   age         305 non-null    int64  
 1   sex         305 non-null    int64  
 2   cp          305 non-null    int64  
 3   trestbps   305 non-null    int64  
 4   chol        305 non-null    int64  
 5   fbs         305 non-null    int64  
 6   restecg    305 non-null    int64  
 7   thalach    305 non-null    int64  
 8   exang       305 non-null    int64  
 9   oldpeak    305 non-null    float64 
 10  slope       305 non-null    int64  
 11  ca          305 non-null    int64  
 12  thal        305 non-null    int64  
 13  target      305 non-null    int64  
 14  test        305 non-null    int64  
 15  chol+trestbps 305 non-null  int64  
dtypes: float64(1), int64(15)
memory usage: 38.2 KB
```

In [168...]

```
fig, ax = plt.subplots()
```

```
ax.scatter(target0['chol'], target0['restecg'], c='blue')
ax.set_xlabel('Chol')
ax.set_ylabel('resting blood pressure')
ax.set_title('Scatter plot for their cholesterol levels vs. blood pressure')

plt.tight_layout()
```



In [154...]

```
fig, ax = plt.subplots()

ax.scatter(target1['chol'], target1['restecg'], c='blue')
ax.set_xlabel('Chol')
ax.set_ylabel('resting blood pressure')
ax.set_title('Scatter plot for their cholesterol levels vs. blood pressure')

plt.tight_layout()
```



In [ ]:

```
## In R
```

In [174...]

```
%load_ext rpy2.ipython
```

The rpy2.ipython extension is already loaded. To reload it, use:  
%reload\_ext rpy2.ipython

In [193...]

```
%%R -i heart_augmented -w 15 -h 10 --units cm
library("ggplot2")
ggplot(heart_augmented,aes(x=chol, y=restecg)) +
  geom_point() +
  facet_wrap(~target) +
  labs(title="Scatter plot for their cholesterol levels vs. blood pressure",x='Chol',
```

