

Unit 4 - Noisy data

1. Histograms
2. Standard deviation
3. Boxplots

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

1. Histograms

Gym example

Taken from: <https://data36.com/plot-histogram-python-pandas/>

```
In [2]: mu = 168 #mean
sigma = 5 #stddev
sample = 250
np.random.seed(0)
height_f = np.random.normal(mu, sigma, sample).astype(int)
```

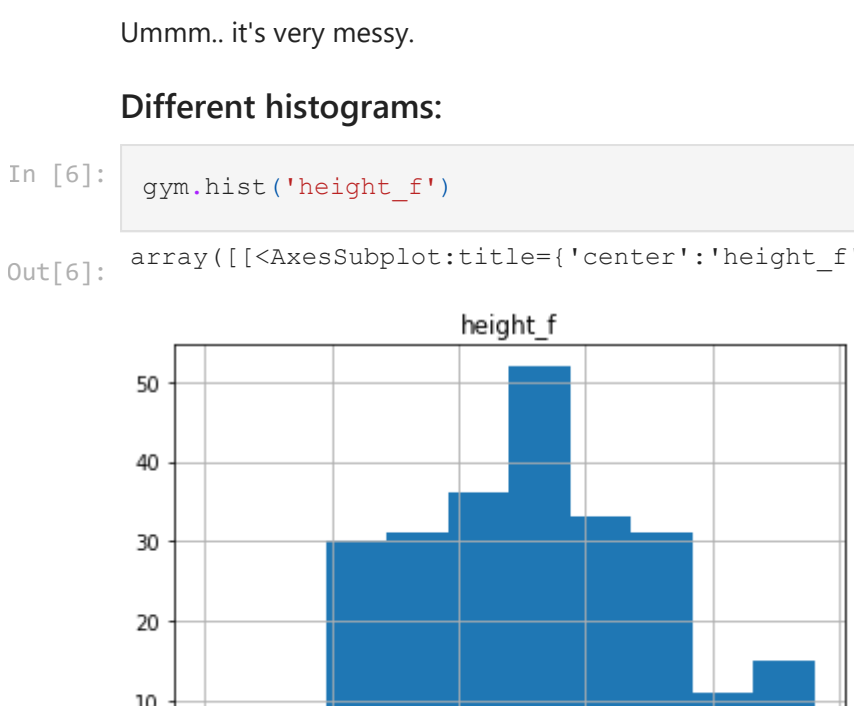
```
In [3]: mu = 176 #mean
sigma = 6 #stddev
sample = 250
np.random.seed(1)
height_m = np.random.normal(mu, sigma, sample).astype(int)
```

```
In [4]: gym = pd.DataFrame({'height_f': height_f, 'height_m': height_m})
gym.head(7)
```

```
Out[4]:   height_f  height_m
0      176      185
1      170      172
2      172      172
3      179      169
4      177      181
5      163      162
6      172      186
```

To begin with - why a histogram?

This is what happens if we just plot lines:

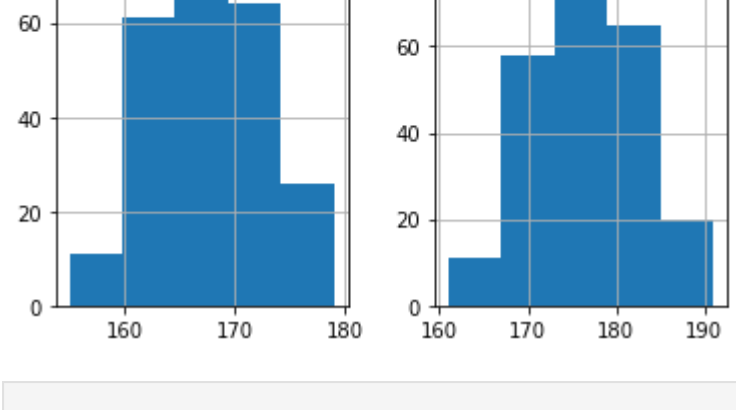


Ummm... it's very messy.

Different histograms:

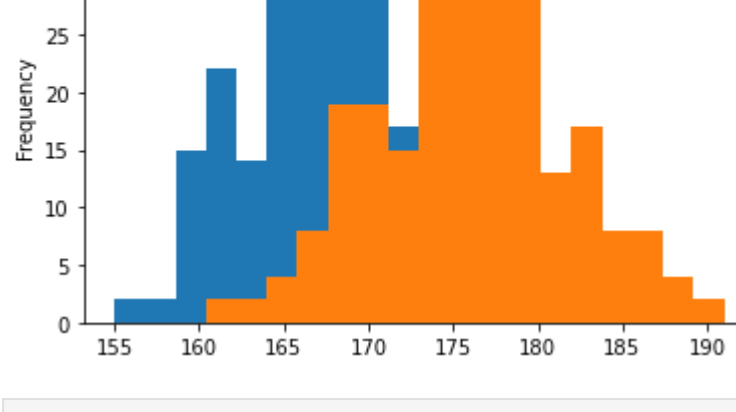
```
In [6]: gym.hist('height_f')
```

Out[6]: array([[<AxesSubplot:title='center':'height_f'>]], dtype=object)



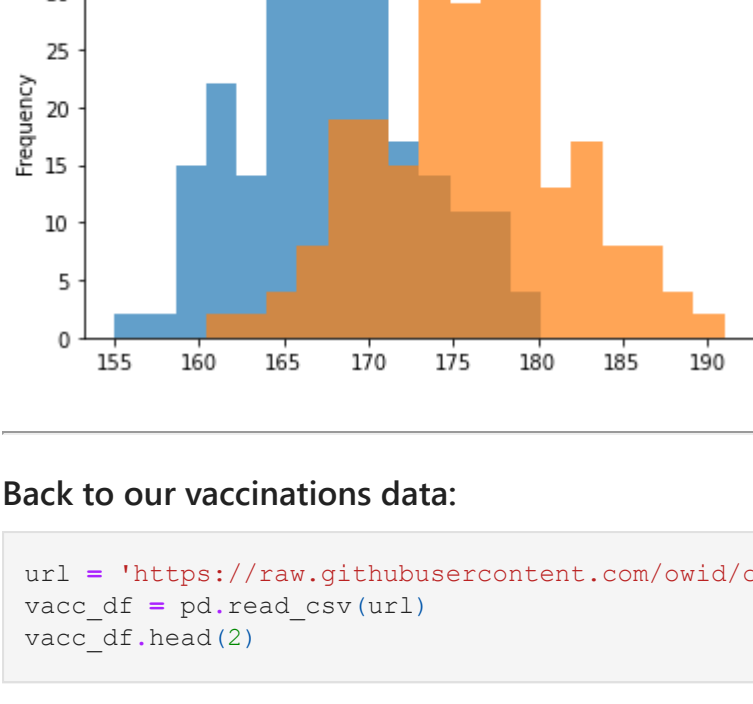
```
In [7]: gym.hist(bins=5)
```

Out[7]: array([[<AxesSubplot:title='center':'height_f'>],
[<AxesSubplot:title='center':'height_m'>]], dtype=object)



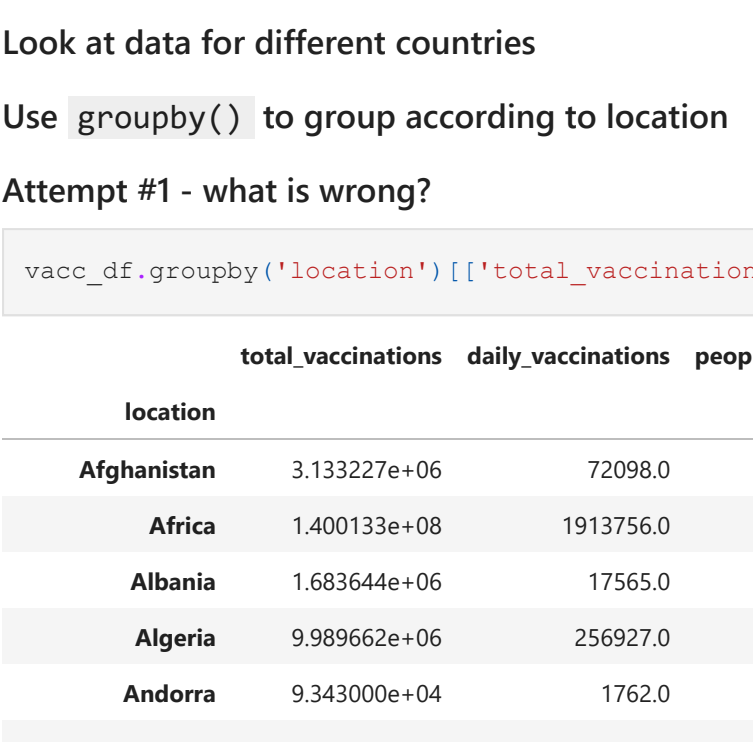
```
In [8]: gym.plot.hist(bins=20)
```

Out[8]: <AxesSubplot:ylabel='Frequency'>



```
In [9]: gym.plot.hist(bins=20, alpha=0.7)
```

Out[9]: <AxesSubplot:ylabel='Frequency'>



Back to our vaccinations data:

```
In [10]: url = 'https://raw.githubusercontent.com/owid/covid-19-data/master/public/data/vaccinations/vaccinations.csv'
vacc_df = pd.read_csv(url)
vacc_df.head(2)
```

```
Out[10]:   location  iso_code  date  total_vaccinations  people_vaccinated  people_fully_vaccinated  total_boosters  daily_vaccinations_raw  daily_vaccinations_smoothed
0  Afghanistan    AFG  2021-02-22              0.0                0.0                NaN                NaN                NaN                NaN
1  Afghanistan    AFG  2021-02-23              NaN                NaN                NaN                NaN                NaN                NaN
```

Look at data for different countries

Use `groupby()` to group according to location

Attempt #1 - what is wrong?

```
In [11]: vacc_df.groupby('location')[['total_vaccinations', 'daily_vaccinations', 'people_fully_vaccinated', 'people_fully_vaccinated_per_hundred']]
```

```
Out[11]:   total_vaccinations  daily_vaccinations  people_fully_vaccinated  people_fully_vaccinated_per_hundred
location
Afghanistan      3.133227e+06          72098.0          4.307440e+05                1.08
Africa           1.400133e+08          1913756.0          5.702818e+07                4.15
Albania           1.683644e+06           17565.0           7.540200e+05                26.25
Algeria           9.989662e+06          256927.0          4.174623e+06                9.36
Andorra           9.343000e+04           1762.0           4.183100e+04                54.08
...
Wallis and Futuna  1.023100e+04            343.0           5.016000e+03                45.21
World             6.113763e+09          43375360.0          2.563762e+09                32.56
Yemen             3.229340e+05           10240.0           1.490900e+04                0.05
Zambia            6.703610e+05           13814.0           2.919470e+05                1.54
Zimbabwe          5.218672e+06           75000.0           2.188630e+06                14.50
```

234 rows × 4 columns

Attempt #2

Is this better?

```
In [12]: vacc_df.fillna(0).groupby('location')[['total_vaccinations', 'daily_vaccinations', 'people_fully_vaccinated', 'people_fully_vaccinated_per_hundred']]
```

```
Out[12]:   total_vaccinations  daily_vaccinations  people_fully_vaccinated  people_fully_vaccinated_per_hundred
location
Afghanistan      NaN                NaN                NaN                NaN
Africa            NaN                NaN                NaN                NaN
Albania           NaN                NaN                NaN                NaN
Algeria           NaN                NaN                NaN                NaN
Andorra           NaN                NaN                NaN                NaN
...
Wallis and Futuna NaN                NaN                NaN                NaN
World             6.113763e+09          43375360.0          2.563762e+09                32.56
Yemen            NaN                NaN                NaN                NaN
Zambia           NaN                NaN                NaN                NaN
Zimbabwe         NaN                NaN                NaN                NaN
```

234 rows × 4 columns

Attempt #3 - change the 0 from string to int, and finally, it works :-)

```
In [13]: vacc_df.fillna(0).groupby('location')[['total_vaccinations', 'daily_vaccinations', 'people_fully_vaccinated', 'people_fully_vaccinated_per_hundred']]
```

```
Out[13]:   total_vaccinations  daily_vaccinations  people_fully_vaccinated  people_fully_vaccinated_per_hundred
location
Afghanistan      3.133227e+06          72098.0          4.307440e+05                1.08
Africa           1.400133e+08          1913756.0          5.702818e+07                4.15
Albania           1.683644e+06           17565.0           7.540200e+05                26.25
Algeria           9.989662e+06          256927.0          4.174623e+06                9.36
Andorra           9.343000e+04           1762.0           4.183100e+04                54.08
...
Wallis and Futuna  1.023100e+04            343.0           5.016000e+03                45.21
World             6.113763e+09          43375360.0          2.563762e+09                32.56
Yemen             3.229340e+05           10240.0           1.490900e+04                0.05
Zambia            6.703610e+05           13814.0           2.919470e+05                1.54
Zimbabwe          5.218672e+06           75000.0           2.188630e+06                14.50
```

234 rows × 4 columns

The `world` row shouldn't be there. Remove it using `.drop()` and `.index()` :

```
In [14]: vacc_df.drop(vacc_df.loc[vacc_df.location == 'World'].index, inplace = True)
```

Your turn:

>What do you think `.index` does? Why is it there? How can you find out?

Before we continue, just assign this new row to a new dataframe, will be easier

```
In [15]: grouped_df = vacc_df.fillna(0).groupby('location')[['total_vaccinations', 'daily_vaccinations', 'people_fully_vaccinated', 'people_fully_vaccinated_per_hundred']]
grouped_df.tail()
```

```
Out[15]:   total_vaccinations  daily_vaccinations  people_fully_vaccinated  people_fully_vaccinated_per_hundred
location
Wales              4589226.0          33151.0          2218515.0                69.98
Wallis and Futuna  10231.0            343.0           5016.0                45.21
Yemen              322934.0          10240.0           14909.0                0.05
Zambia             670361.0          13814.0           291947.0                1.54
Zimbabwe           5218672.0          75000.0          2188630.0                14.50
```

sort the values using 'sort_values()'

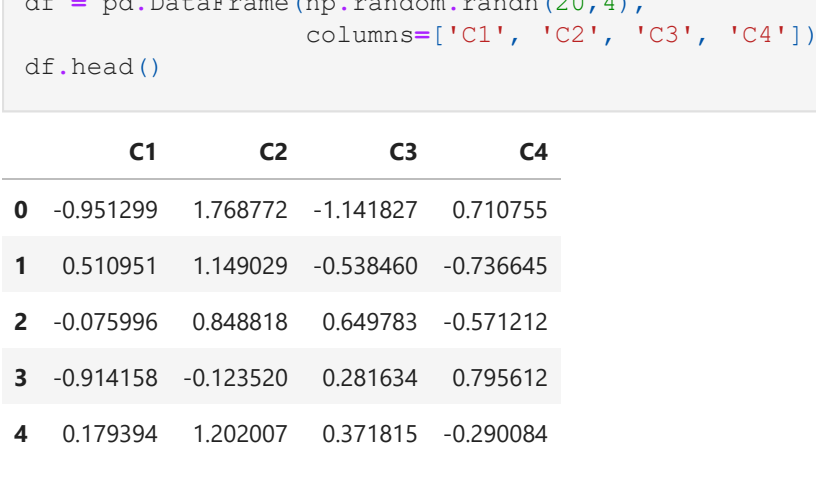
```
In [16]: grouped_df.sort_values('people_fully_vaccinated_per_hundred', ascending = False).head(10)
```

```
Out[16]:   total_vaccinations  daily_vaccinations  people_fully_vaccinated  people_fully_vaccinated_per_hundred
location
Gibraltar          79335.0            1068.0           39463.0                117.13
Pitcairn            94.0              1.0             47.0                100.00
Portugal           15904393.0          150618.0          8602244.0                84.60
Cayman Islands      109789.0           1024.0           54890.0                82.54
United Arab Emirates 19847232.0          155312.0          8169539.0                81.77
Malta               816136.0           7557.0          419919.0                81.61
Iceland            551110.0           7546.0           275493.0                80.23
Spain              69740837.0          577917.0          36335711.0                77.73
Singapore           9256975.0           47359.0          4566329.0                77.44
Qatar              4690049.0           37344.0          2218292.0                75.70
```

Histogram according to all values:

```
In [17]: grouped_df.hist(bins=50)
```

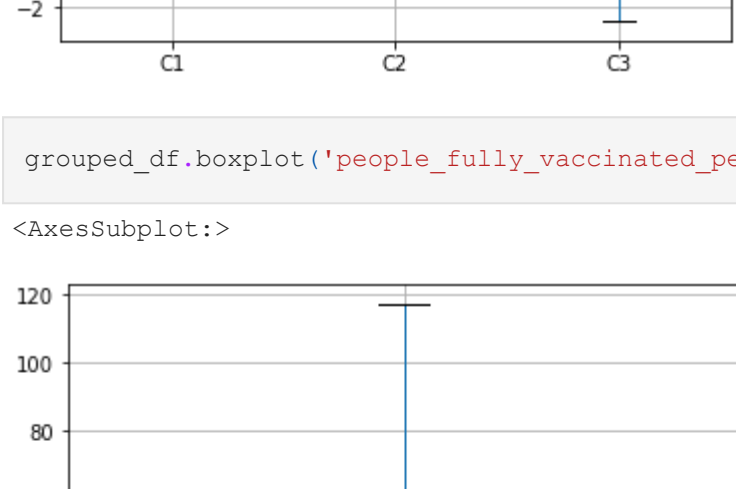
```
Out[17]: array([[<AxesSubplot:title='center':'total_vaccinations'>],
[<AxesSubplot:title='center':'daily_vaccinations'>],
[<AxesSubplot:title='center':'people_fully_vaccinated'>],
[<AxesSubplot:title='center':'people_fully_vaccinated_per_hundred'>]],
dtype=object)
```



Histogram according to `people_fully_vaccinated_per_hundred`

```
In [18]: grouped_df.hist('people_fully_vaccinated_per_hundred', bins=50)
```

```
Out[18]: array([[<AxesSubplot:title='center':'people_fully_vaccinated_per_hundred'>]],
dtype=object)
```



Remove rows with 0's

Note that this is different than changing values to 0's

```
In [19]: grouped_df.drop(grouped_df[grouped_df.people_fully_vaccinated_per_hundred == 0.0].index, inplace=True)
grouped_df.hist('people_fully_vaccinated_per_hundred', bins=50)
```

```
Out[19]: array([[<AxesSubplot:title='center':'people_fully_vaccinated_per_hundred'>]],
dtype=object)
```



Your turn:

Do the same, but for another column

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

Functions covered in this unit:

- `.std()` - standard deviation
- `.hist(data, num_bins)` - the data and the number of bins
- `.plot()` - simple line plot
- `.hist(data, num_bins, alpha)` - the data, the number of bins and the transparency (default is 10 bins, all data and not transparent)
- `.plot.hist()` - histograms on same chart
- `.drop()` - removes unwanted rows or columns
- `.index()` - index of the selected rows
- `.sort_values()` - self explanatory, it just sorts values

2. Standard deviation

A small example showing that the more the data is spread, the higher the std:

```
In [20]: df = pd.DataFrame({'height': [161, 156, 172],
'weight': [67, 65, 89],
'age': [20, 20, 20]})
df
```

```
Out[20]:   height  weight  age
0      161      67   20
1      156      65   20
2      172      89   20
```

```
In [21]: df.mean()
```

```
Out[21]: height    163.000000
weight      73.666667
age         20.000000
dtype: float64
```

```
In [22]: df.std()
```

```
Out[22]: height      8.185353
weight     13.316656
age          0.000000
dtype: float64
```

3. Boxplots

```
In [23]: np.random.seed(2345)
df = pd.DataFrame(np.random.randn(20, 4),
columns=['C1', 'C2', 'C3', 'C4'])
df.head()
```

```
Out[23]:   C1      C2      C3      C4
0 -0.951299  1.768772 -1.141827  0.710755
1  0.510951  1.149029 -0.538460 -0.736645
2 -0.075996  0.848818  0.649783 -0.571212
3 -0.914158 -0.123520  0.281634  0.795612
4  0.179394  1.202007  0.371815 -0.290084
```

```
In [24]: boxplot = df.boxplot(column=['C1', 'C2', 'C3'])
```



```
In [25]: grouped_df.boxplot('people_fully_vaccinated_per_hundred')
```

```
Out[25]: <AxesSubplot:>
```



```
In [26]: grouped_df.sort_values('people_fully_vaccinated_per_hundred')
```

```
Out[26]:   total_vaccinations  daily_vaccinations  people_fully_vaccinated  people_fully_vaccinated_per_hundred
location
Democratic Republic of Congo  134860.0          2328.0          37532.0                0.04
Yemen            322934.0          10240.0          14909.0                0.05
Chad             123052.0          4722.0           24839.0                0.15
Turkmenistan     41993.0              0.0            9753.0                0.16
Haiti            61545.0          1565.0           18987.0                0.16
...
United Arab Emirates  19847232.0          155312.0          8169539.0                81.77
Cayman Islands      109789.0           1024.0           54890.0                82.54
Portugal           15904393.0          150618.0          8602244.0                84.60
Pitcairn            94.0              1.0             47.0                100.00
Gibraltar          79335.0            1068.0           39463.0                117.13
```

232 rows × 4 columns

This doesn't seem to fit <https://ourworldindata.org/grapher/covid-vaccination-doses-per-capita>

Is this an error???

```
In [ ]:
```