

Course 2-3, part I - Data types and tidy data

Tidy data (see [here](#) for more details)

This is one of the non-tidy dataset assembled by Hadley Wickham (check out [here](#) for more datasets, explanation, and R code).

Let's take a look at this small dataset:

<https://raw.githubusercontent.com/tidyverse/tidyr/4c0a8d0fdb9372302fcc57ad995d57a43d9e4337/vignettes/pew.csv>

```
In [101... import pandas as pd
```

```
In [102... pew_df = pd.read_csv('https://raw.githubusercontent.com/tidyverse/tidyr/4c0a8d0fdb9372302fcc57ad995d57a43d9e4337/vignettes/pew.csv')  
pew_df
```

Out[102...

	religion	<\$10k	\$10-20k	\$20-30k	\$30-40k	\$40-50k	\$50-75k	\$75-100k	\$100-150k	>150k	Don't know/refused
0	Agnostic	27	34	60	81	76	137	122	109	84	96
1	Atheist	12	27	37	52	35	70	73	59	74	76
2	Buddhist	27	21	30	34	33	58	62	39	53	54
3	Catholic	418	617	732	670	638	1116	949	792	633	1489
4	Don't know/refused	15	14	15	11	10	35	21	17	18	116
5	Evangelical Prot	575	869	1064	982	881	1486	949	723	414	1529
6	Hindu	1	9	7	9	11	34	47	48	54	37
7	Historically Black Prot	228	244	236	238	197	223	131	81	78	339
8	Jehovah's Witness	20	27	24	24	21	30	15	11	6	37
9	Jewish	19	19	25	25	30	95	69	87	151	162
10	Mainline Prot	289	495	619	655	651	1107	939	753	634	1328
11	Mormon	29	40	48	51	56	112	85	49	42	69
12	Muslim	6	7	9	10	9	23	16	8	6	22
13	Orthodox	13	17	23	32	32	47	38	42	46	73
14	Other Christian	9	7	11	13	13	14	18	14	12	18
15	Other Faiths	20	33	40	46	49	63	46	40	41	71
16	Other World Religions	5	2	3	4	2	7	3	4	4	8
17	Unaffiliated	217	299	374	365	341	528	407	321	258	597

This dataset is about the relationships between income and religion, assembled from a research by the Pew Research Center. You can read more details [here](#). Is this dataset tidy or not? Why?

Yes, many of the columns are values, not variable names. How should we fix it?

Pandas provides a convenient function called `melt`. You specify the `id_vars` that are variable columns, and `value_vars` that are value columns, and provide the name for the variable as well as the name for the values.

Q: so please go ahead and tidy it up! I'd suggest to use the variable name "income" and value name "frequency"

```
In [103... pew_tidy_df = pd.melt(pew_df, id_vars=["religion"], var_name="income", value_name="frequency")
pew_tidy_df.head()
```

```
Out[103... 
```

	religion	income	frequency
0	Agnostic	<\$10k	27
1	Atheist	<\$10k	12
2	Buddhist	<\$10k	27
3	Catholic	<\$10k	418
4	Don't know/refused	<\$10k	15

If you were successful, you'll have something like this:

```
In [104... pew_tidy_df.sample(10)
```

Out[104...

	religion	income	frequency
36	Agnostic	\$20-30k	60
165	Catholic	Don't know/refused	1489
97	Historically Black Prot	\$50-75k	223
21	Catholic	\$10-20k	617
164	Buddhist	Don't know/refused	54
113	Evangelical Prot	\$75-100k	949
117	Jewish	\$75-100k	69
110	Buddhist	\$75-100k	62
111	Catholic	\$75-100k	949
0	Agnostic	<\$10k	27

Data types

Let's talk about data types briefly. Understanding data types is not only important for choosing the right visualizations, but also important for efficient computing and storage of data. A Pandas `Dataframe` is essentially a bunch of `Series`, and those `Series` are essentially `numpy` arrays. An array may contain a fixed-length items such as integers or variable length items such as strings. Putting some efforts to think about the correct data type can potentially save a lot of memory as well as time.

A nice example would be the categorical data type. If you have a variable that only has several possible values, it's essentially a categorical data. Take a look at the `income` variable.

In [105...

```
pew_tidy_df.income.value_counts()
```

```
Out[105... income
<$10k      18
$10-20k    18
$20-30k    18
$30-40k    18
$40-50k    18
$50-75k    18
$75-100k   18
$100-150k  18
>150k      18
Don't know/refused 18
Name: count, dtype: int64
```

These were the column names in the original non-tidy data. The value can take only one of these income ranges and thus it is a categorical data. What is the data type that pandas use to store this column?

```
In [106... pew_tidy_df.income.dtype
```

```
Out[106... dtype('O')
```

The `O` means that it is an object data type, which does not have a fixed size like integer or float. The series contains a sort of pointer to the actual text objects. You can actually inspect the amount of memory used by the dataset.

```
In [107... pew_tidy_df.memory_usage()
```

```
Out[107... Index      132
religion    1440
income      1440
frequency   1440
dtype: int64
```

```
In [108... pew_tidy_df.memory_usage(deep=True)
```

```
Out[108... Index      132
religion    11260
income      10260
frequency   1440
dtype: int64
```

When you don't specify `deep=True`, the memory usage method just tells you the amount of memory used by the numpy arrays in the pandas dataframe. When you pass `deep=True`, it tells you the total amount of memory by including the memory used by all the text objects. So, the `religion` and `income` columns occupies almost ten times of memory than the `frequency` column, which is simply an array of integers.

```
In [109... pew_tidy_df.frequency.dtype
```

```
Out[109... dtype('int64')
```

Is there any way to save up the memory? Note that there are only 10 categories in the income variable. That means we just need 10 numbers to represent the categories! Of course we need to store the names of each category, but that's just one-time cost. The simplest way to convert a column is using `astype` method.

```
In [110... income_categorical_series = pew_tidy_df.income.astype('category')
# you can do pew_tidy_df.income = pew_tidy_df.income.astype('category')
```

Now, this series has the `CategoricalDtype` dtype.

```
In [111... income_categorical_series.dtype
```

```
Out[111... CategoricalDtype(categories=['$10-20k', '$100-150k', '$20-30k', '$30-40k', '$40-50k',
                              '$50-75k', '$75-100k', '<$10k', '>150k',
                              'Don't know/refused'],
                  , ordered=False, categories_dtype=object)
```

How much memory do we use?

```
In [112... income_categorical_series.memory_usage(deep=True)
```

```
Out[112... 1182
```

```
In [113... pew_tidy_df.income.memory_usage(deep=True)
```

```
Out[113... 10392
```

We have reduced the memory usage by almost 10 fold! Not only that, because now the values are just numbers, it will be much faster to match, filter, manipulate. If your dataset is huge, this can save up a lot of space and time.

If the categories have ordering, you can specify the ordering too.

```
In [114... from pandas.api.types import CategoricalDtype
income_type = CategoricalDtype(categories=["Don't know/refused", '<$10k', '$10-20k', '$20-30k', '$30-40k',
                                         '$40-50k', '$50-75k', '$75-100k', '$100-150k', '>150k'], ordered=True)
income_type
```

```
Out[114... CategoricalDtype(categories=['Don't know/refused', '<$10k', '$10-20k', '$20-30k',
                              '$30-40k', '$40-50k', '$50-75k', '$75-100k', '$100-150k',
                              '>150k'],
                              , ordered=True, categories_dtype=object)
```

```
In [115... pew_tidy_df.income.astype(income_type).dtype
```

```
Out[115... CategoricalDtype(categories=['Don't know/refused', '<$10k', '$10-20k', '$20-30k',
                              '$30-40k', '$40-50k', '$50-75k', '$75-100k', '$100-150k',
                              '>150k'],
                              , ordered=True, categories_dtype=object)
```

This data type now allows you to compare and sort based on the ordering.

Q: ok, now convert both religion and income columns of `pew_tidy_df` as categorical dtype (in place) and show that `pew_tidy_df` now uses much less memory

```
In [116... # Show memory usage before conversion
print("Memory usage before conversion:")
print(pew_tidy_df.memory_usage(deep=True))

# Convert the columns to categorical dtype in place
pew_tidy_df['religion'] = pew_tidy_df['religion'].astype('category')
pew_tidy_df['income'] = pew_tidy_df['income'].astype('category')

# Show memory usage after conversion
print("\nMemory usage after conversion:")
print(pew_tidy_df.memory_usage(deep=True))
```

Memory usage before conversion:

```
Index      132
religion   11260
income     10260
frequency  1440
dtype: int64
```

Memory usage after conversion:

```
Index      132
religion   1883
income     1050
frequency  1440
dtype: int64
```

Other resources on tidy data in python

- [Jean-Nicholas Hould: Tidy Data in Python](#)

Part II - 1D data

First import basic packages and then load a dataset from `vega_datasets` package. If you don't have `vega_datasets` or `altair` installed yet, use `pip` or `conda` to install them.

```
In [117... import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from vega_datasets import data
```

```
In [118... cars = data.cars()
cars.head()
```


Out[118...

	Name	Miles_per_Gallon	Cylinders	Displacement	Horsepower	Weight_in_lbs	Acceleration	Year	Origin
0	chevrolet chevelle malibu	18.0	8	307.0	130.0	3504	12.0	1970-01-01	USA
1	buick skylark 320	15.0	8	350.0	165.0	3693	11.5	1970-01-01	USA
2	plymouth satellite	18.0	8	318.0	150.0	3436	11.0	1970-01-01	USA
3	amc rebel sst	16.0	8	304.0	150.0	3433	12.0	1970-01-01	USA
4	ford torino	17.0	8	302.0	140.0	3449	10.5	1970-01-01	USA

1D scatter plot

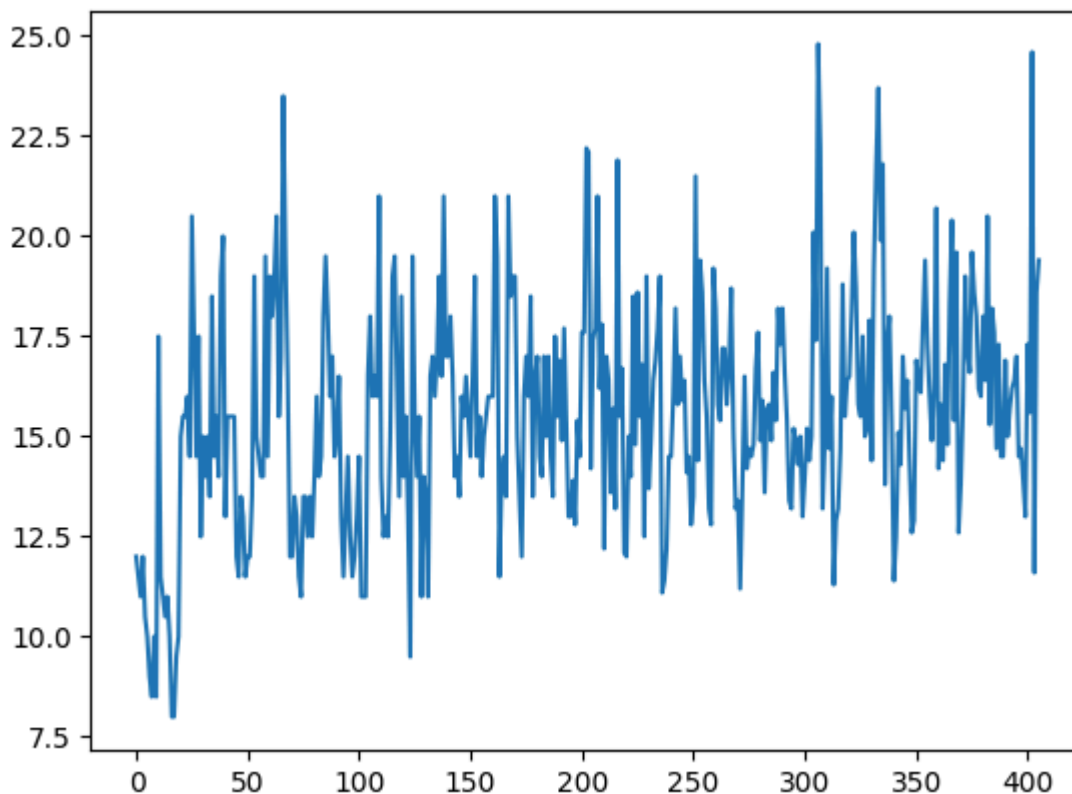
Let's consider the `Acceleration` column as our 1D data. If we ask pandas to plot this series, it'll produce a line graph where the index becomes the horizontal axis.

In [119...

```
cars.Acceleration.plot()
```

Out[119...

```
<Axes: >
```



Because the index is not really meaningful, drawing line between subsequent values is **misleading!**

It's not trivial to use pandas to create a 1-D scatter plot. Instead, we can use `matplotlib`'s `scatter` function. We can first create an array with zeros that we can use as the vertical coordinates of the points that we will plot. `np.zeros_like` returns an array with zeros that matches the shape of the input array.

```
In [120...] np.zeros_like([1,2,3])
```

```
Out[120...] array([0, 0, 0])
```

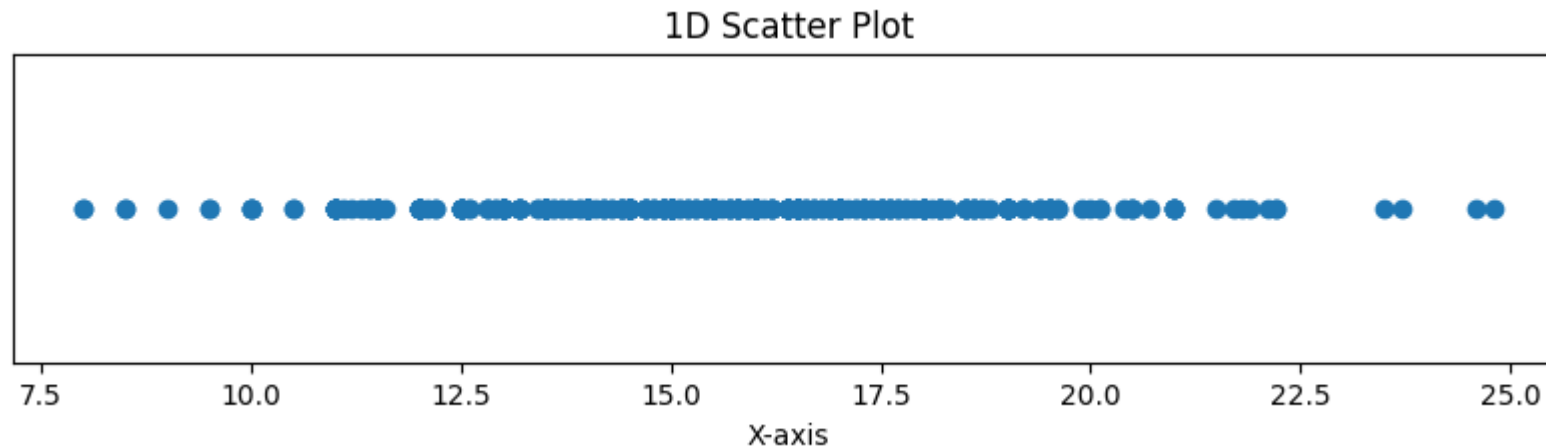
Q: now can you create an 1D scatter plot with `matplotlib`'s `scatter` function? Make the figure wide (e.g. set `figsize=(10,2)`) and then remove the y ticks.

```
In [121... y = np.zeros_like(cars.Acceleration)
```

```
plt.figure(figsize=(10,2))  
plt.scatter(cars.Acceleration,y)  
plt.yticks([])
```

```
plt.xlabel('X-axis')  
plt.title('1D Scatter Plot')
```

```
plt.show()
```



As you can see, there are lots of occlusions. So this plot cannot show the distribution properly and we would like to fix it. How about adding some jitters? You can use `numpy`'s `random.rand()` function to generate random numbers, instead of using an array with zeros.

Q: create a jittered 1D scatter plot.

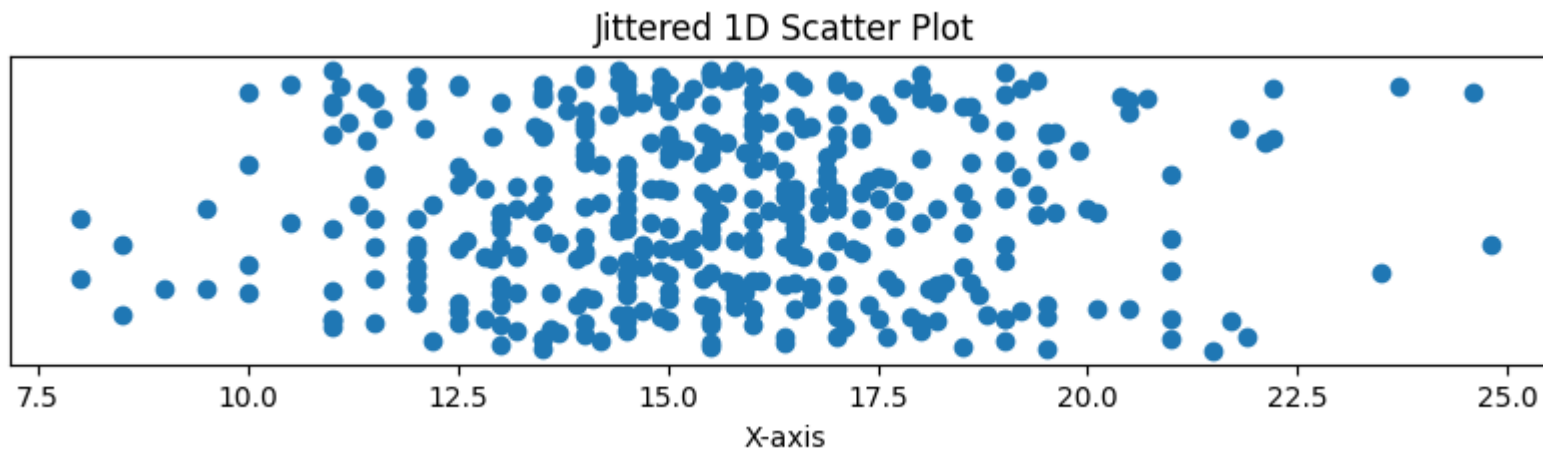
```
In [122... jittered_y = np.random.uniform(-1, 1, size= cars.Acceleration.shape)
```

```
plt.figure(figsize=(10,2))  
plt.scatter(cars.Acceleration, jittered_y)
```

```
# Remove y-axis ticks  
plt.yticks([])
```

```
plt.xlabel('X-axis')
plt.title('Jittered 1D Scatter Plot')

plt.show()
```



We can further improve this by adding transparency to the symbols. The transparency option for `scatter` function is called `alpha`. Set it to be 0.2.

Q: create a jittered 1D scatter plot with transparency (alpha=0.2)

```
In [123... jittered_y = np.random.uniform(-1, 1, size= cars.Acceleration.shape)

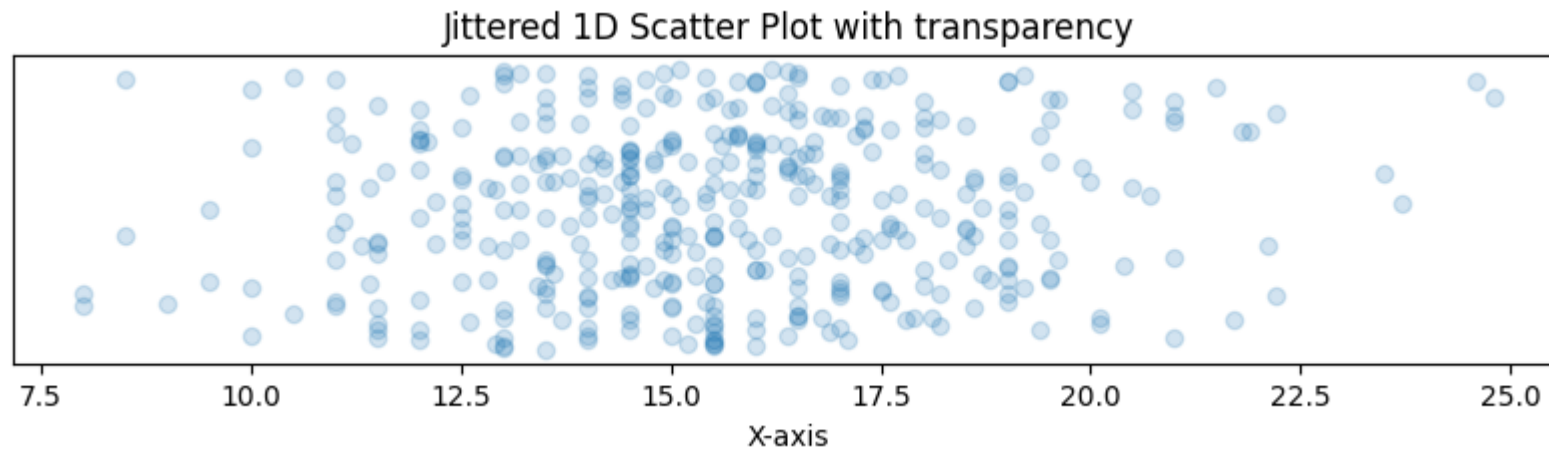
plt.figure(figsize=(10,2))

plt.scatter(cars.Acceleration, jittered_y, alpha=0.2)

plt.yticks([])

plt.xlabel('X-axis')
plt.title('Jittered 1D Scatter Plot with transparency')

plt.show()
```



Another strategy is using empty symbols. The option is `facecolors`. You can also change the stroke color (`edgecolors`).

Q: create a jittered 1D scatter plot with empty symbols.

```
In [124... jittered_y = np.random.uniform(-1, 1, size= cars.Acceleration.shape)

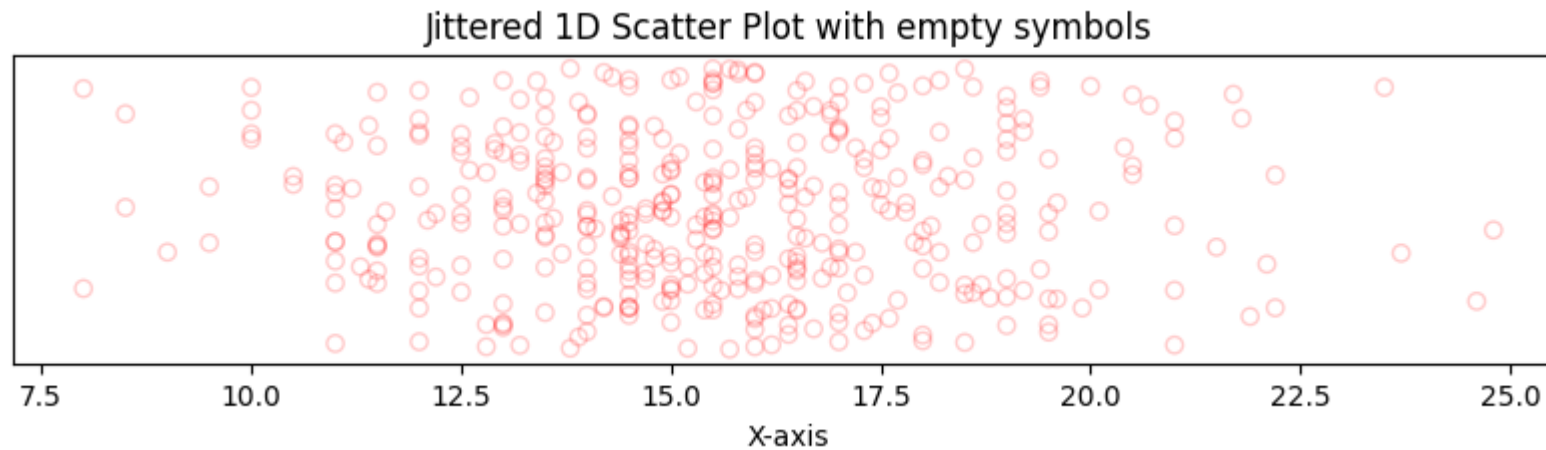
plt.figure(figsize=(10,2))

plt.scatter(cars.Acceleration, jittered_y, alpha=0.2, marker='o', facecolors='none', edgecolors='red')

# Remove y-axis ticks
plt.yticks([])

plt.xlabel('X-axis')
plt.title('Jittered 1D Scatter Plot with empty symbols')

plt.show()
```



What happens if you have lots and lots of points?

Whatever strategy that you use, it's almost useless if you have too many data points. Let's play with different number of data points and see how it looks.

It not only becomes completely useless, it also take a while to draw the plot itself.

```
In [125... Ns = [100, 1000, 10000, 100000]

fig, axs = plt.subplots(len(Ns), 1, figsize=(10, 2 * len(Ns)), squeeze=False)

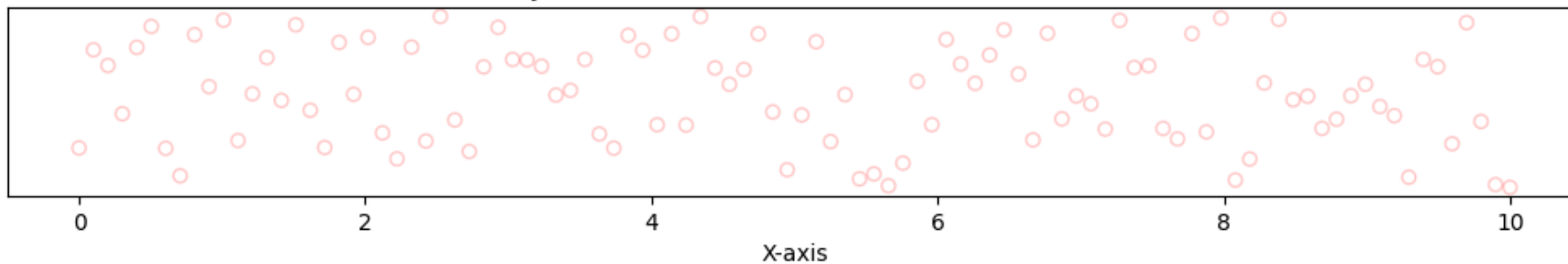
for ax, N in zip(axs.flat, Ns):
    x = np.linspace(0, 10, N)
    y = np.random.uniform(-0.1, 0.1, size=x.shape)
    ax.scatter(x, y, alpha=0.2, marker='o', facecolors='none', edgecolors='red')

    ax.set_yticks([])

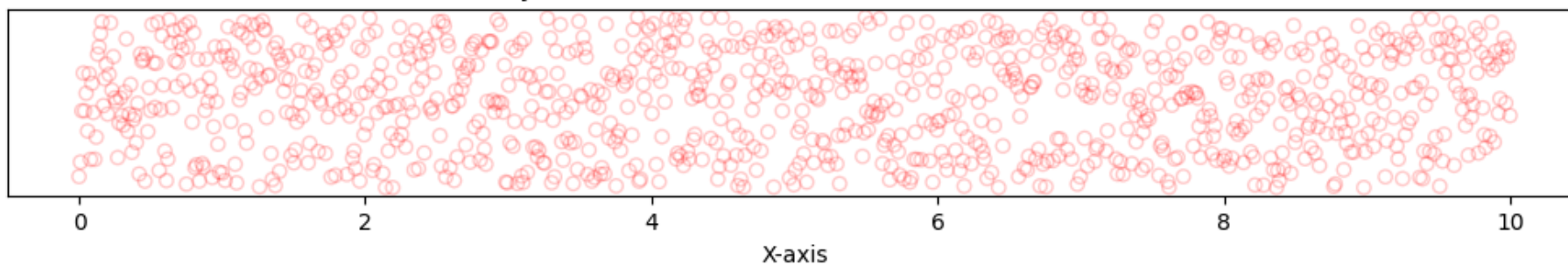
    ax.set_xlabel('X-axis')
    ax.set_title(f'Jittered 1D Scatter Plot with {N} Points')

plt.tight_layout()
plt.show()
```

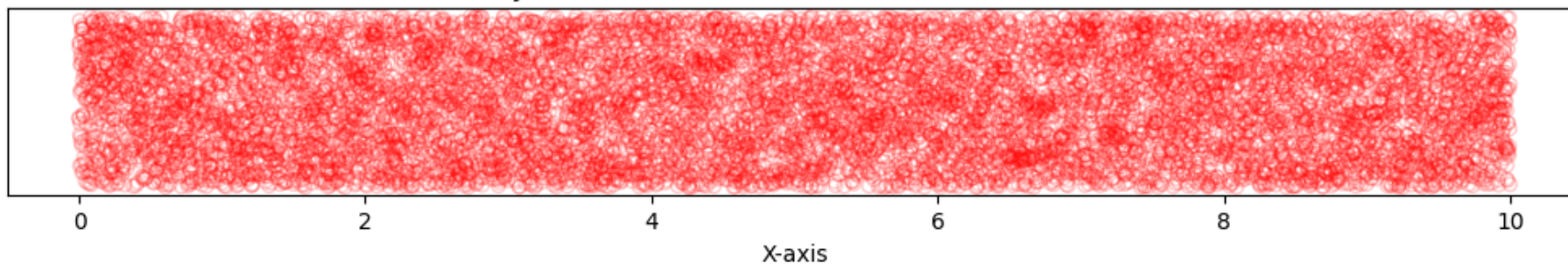
Jittered 1D Scatter Plot with 100 Points



Jittered 1D Scatter Plot with 1000 Points



Jittered 1D Scatter Plot with 10000 Points



Jittered 1D Scatter Plot with 100000 Points



X-axis

Histogram and boxplot

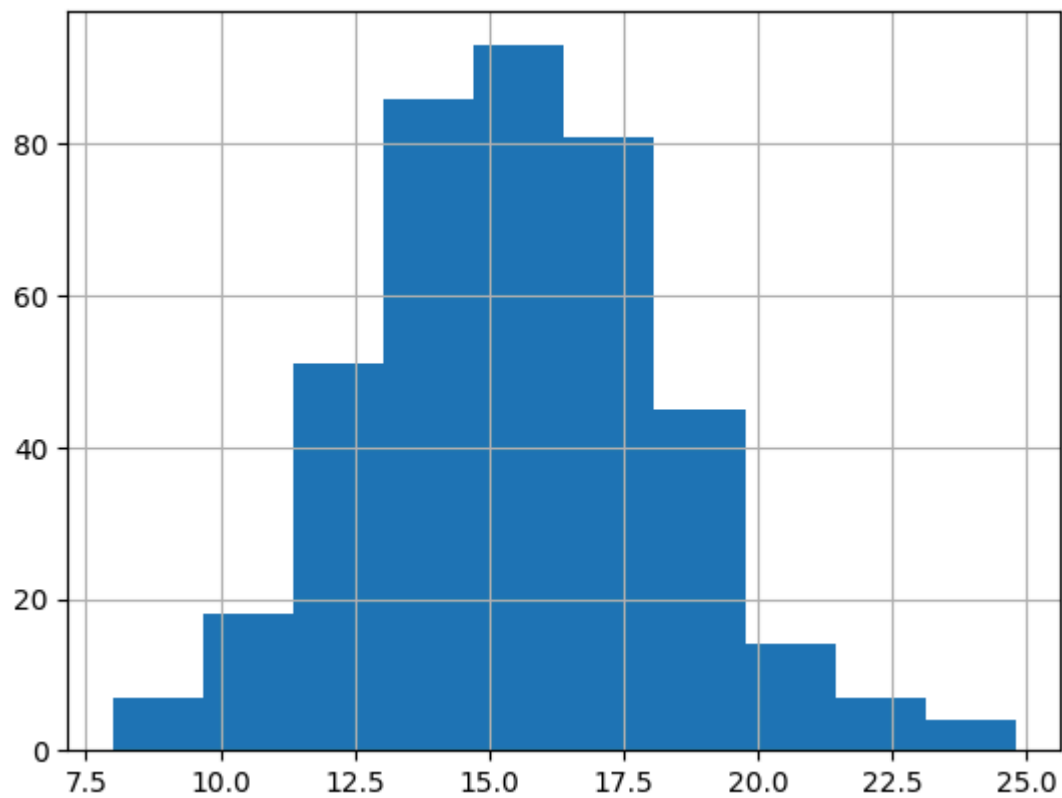
When you have lots of data points, you no longer want to use the scatter plots. Even when you don't have millions of data points, you often want to get a quick summary of the distribution rather than seeing the whole dataset. For 1-D datasets, two major approaches are histogram and boxplot. A histogram is about aggregating and counting the data while a boxplot is about summarizing the data. Let's first draw some histograms.

Histogram

It's very easy to draw a histogram with pandas.

```
In [126... cars.Acceleration.hist()
```

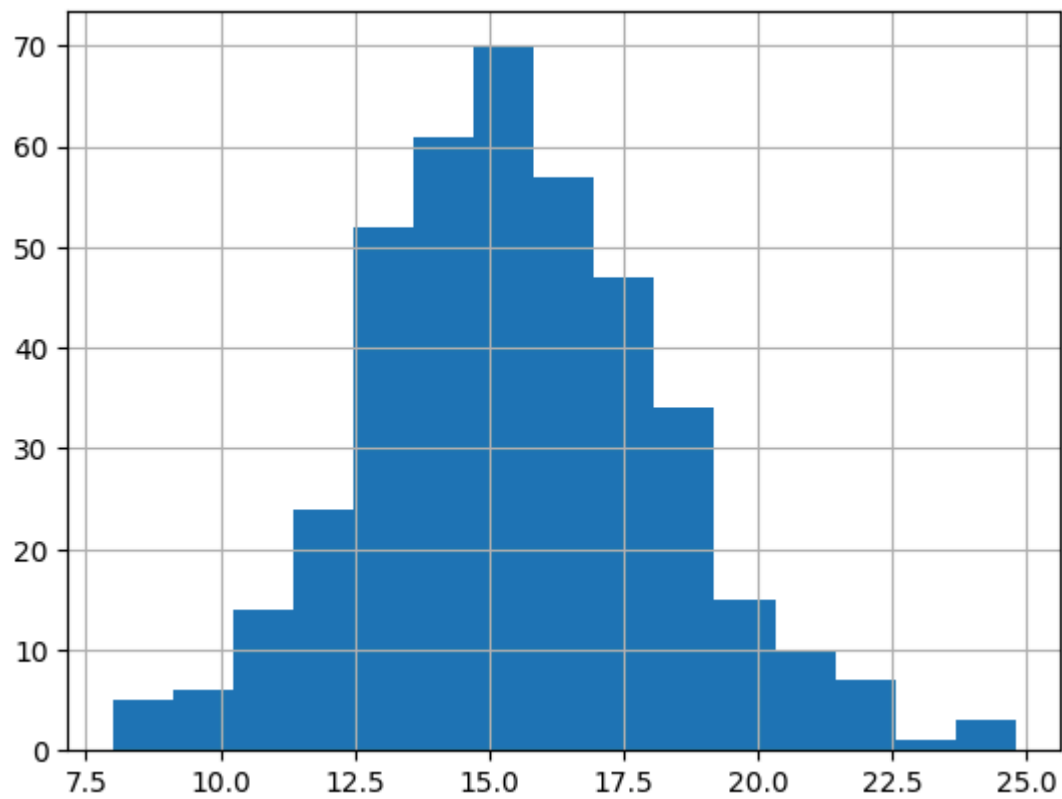
```
Out[126... <Axes: >
```

You can adjust the bin size, which is the main parameter of the histogram.

```
In [127... cars.Acceleration.hist(bins=15)
```

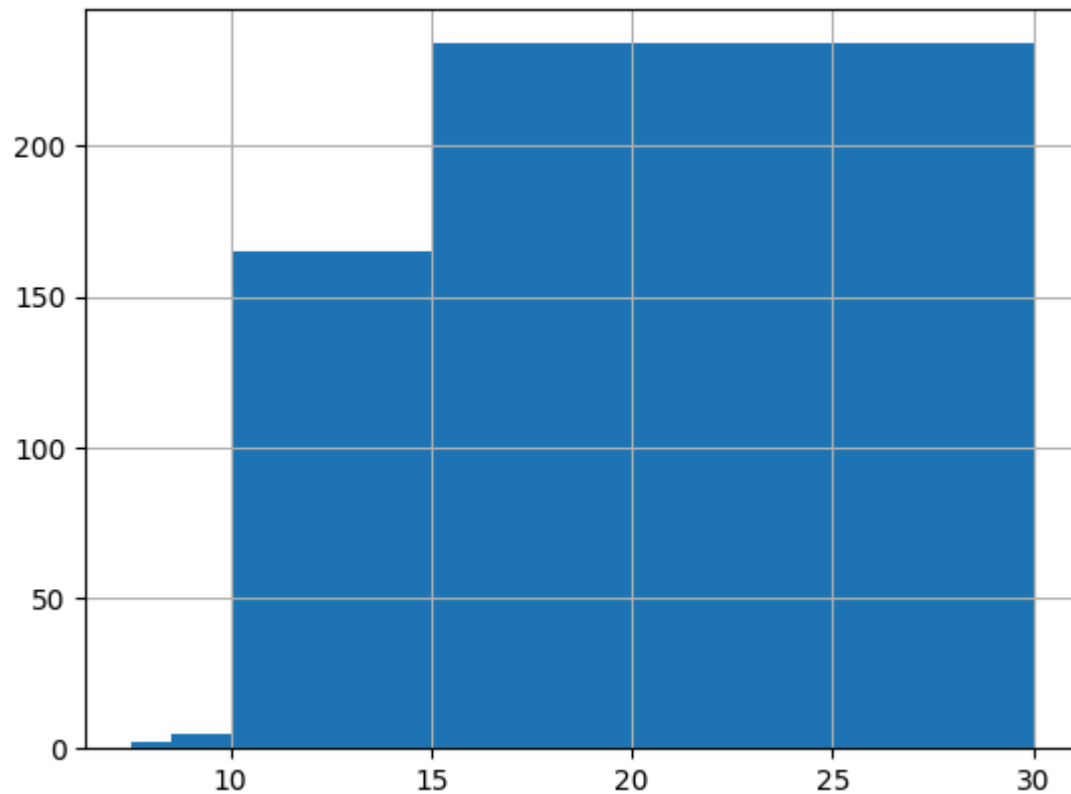
```
Out[127... <Axes: >
```



You can even specify the actual bins.

```
In [128... bins = [7.5, 8.5, 10, 15, 30]  
cars.Acceleration.hist(bins=bins)
```

```
Out[128... <Axes: >
```



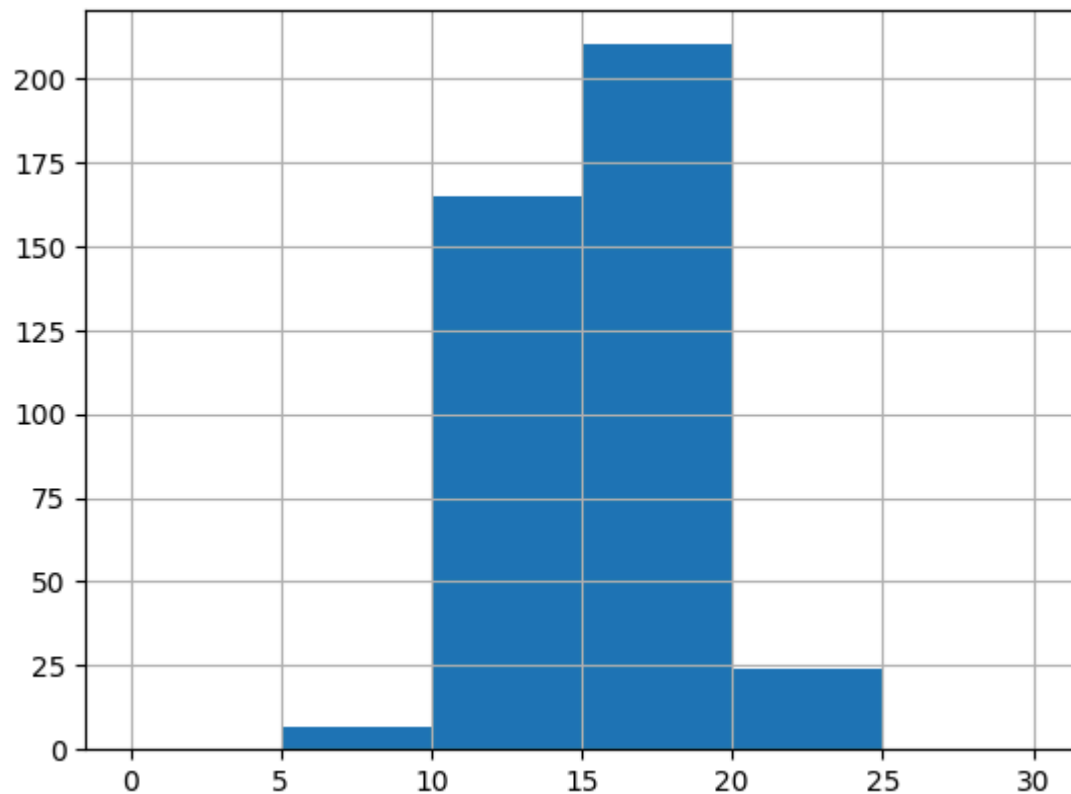
Do you see anything funky going on with this histogram? What's wrong? Can you fix it?

Q: Explain what's wrong with this histogram and fix it.

[Pandas documentation](#) does not show the option that you should use. You should take a look at the `matplotlib`'s documentation.

```
In [129... # the funky behaviour is due to the fact that the first bin is much wider than the others.  
# Also matplotlib centers bars on their bin midpoints by default, so the wide bin doesn't appear to cover the entire range 0-1  
  
# the best approach is to use uniformly spaced bins  
  
uniform_bins=[0, 5, 10, 15, 20, 25, 30]  
  
cars.Acceleration.hist(bins=uniform_bins)
```

Out[129... <Axes: >

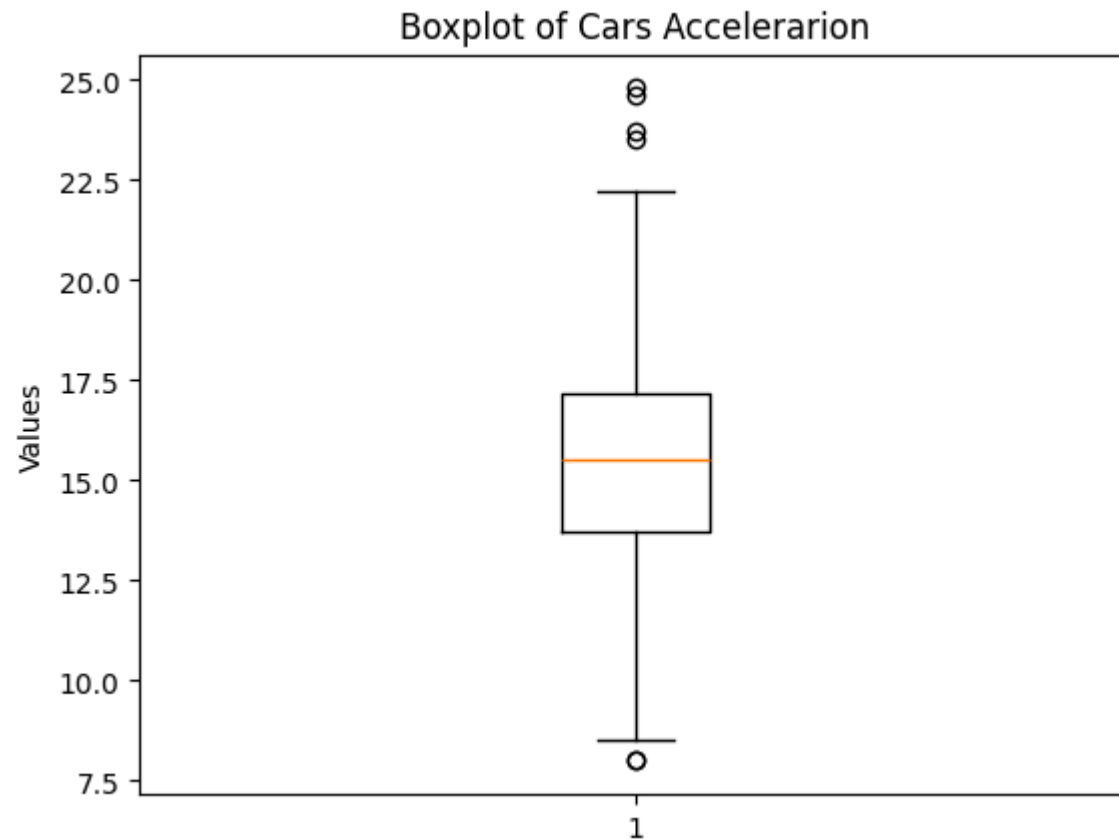


Boxplot

Boxplot can be created with pandas very easily. Check out the `plot` documentation: <https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.plot.html>

Q: create a box plot of Acceleration

```
In [130... plt.boxplot(cars.Acceleration)
plt.title("Boxplot of Cars Accelerarion")
plt.ylabel("Values")
plt.show()
```



1D scatter plot with Seaborn and Altair

As you may have noticed, it is not very easy to use `matplotlib`. The organization of plot functions and parameters are not very systematic. Whenever you draw something, you should search how to do it, what are the parameters you can tweak, etc. You need to manually tweak a lot of things when you work with `matplotlib`.

There are more systematic approaches towards data visualization, such as the "[Grammar of Graphics](#)". This idea of *grammar* led to the famous `ggplot2` (<http://ggplot2.tidyverse.org>) package in R as well as the [Vega & Vega-lite](#) for the web. The grammar-based approach lets you work with *tidy data* in a natural way, and also lets you approach the data visualization systematically.

Let's introduce two nice Python libraries. One is called `seaborn` (<https://seaborn.pydata.org>), which is focused on creating complex statistical data visualizations, and the other is called `altair` (<https://altair-viz.github.io/>) and it is a Python library that lets you *define* a visualization and translates it into vega-lite json.

Seaborn would be useful when you are doing exploratory data analysis; altair may be useful if you are thinking about creating and putting an interactive visualization on the web.

If you don't have them yet, check the [installation page of altair](#). In `conda`,

```
$ conda install -c conda-forge altair vega_datasets jupyterlab
```

Let's play with it.

```
In [131... import seaborn as sns
import altair as alt

# Uncomment the following line if you are using Jupyter notebook
#alt.renderers.enable('notebook')
```

```
In [132... cars.head()
```

```
Out[132...
```

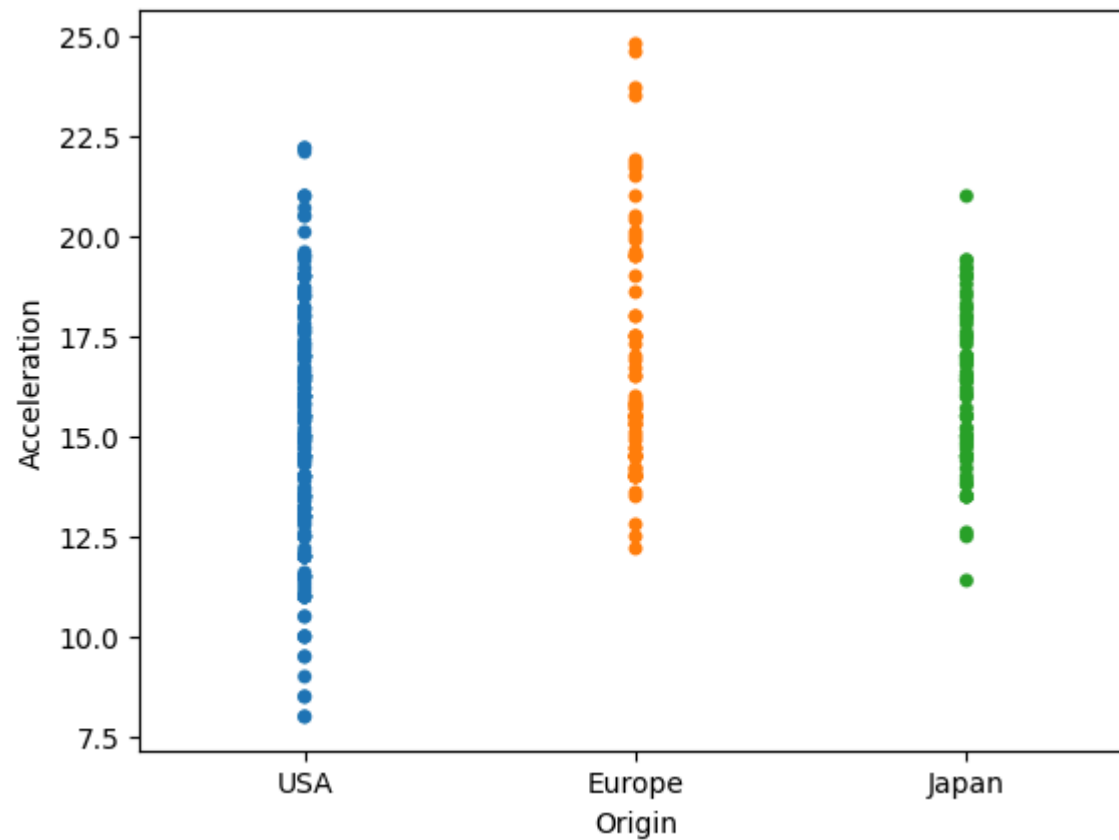
	Name	Miles_per_Gallon	Cylinders	Displacement	Horsepower	Weight_in_lbs	Acceleration	Year	Origin
0	chevrolet chevelle malibu	18.0	8	307.0	130.0	3504	12.0	1970-01-01	USA
1	buick skylark 320	15.0	8	350.0	165.0	3693	11.5	1970-01-01	USA
2	plymouth satellite	18.0	8	318.0	150.0	3436	11.0	1970-01-01	USA
3	amc rebel sst	16.0	8	304.0	150.0	3433	12.0	1970-01-01	USA
4	ford torino	17.0	8	302.0	140.0	3449	10.5	1970-01-01	USA

Beeswarm plots with seaborn

Seaborn has a built-in function to create 1D scatter plots with multiple categories, and it adds jittering by default.

```
In [133...] sns.stripplot(x='Origin', y='Acceleration', data=cars, jitter=False, hue='Origin')
```

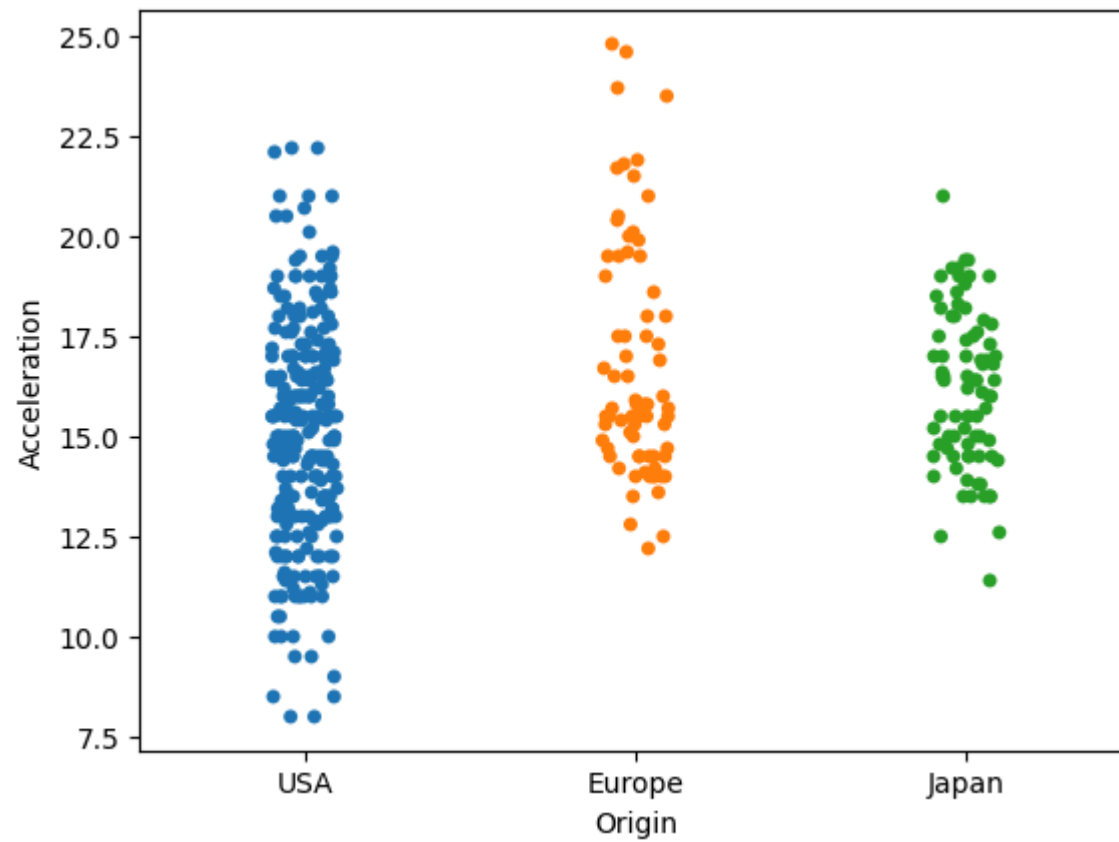
```
Out[133...] <Axes: xlabel='Origin', ylabel='Acceleration'>
```



And you can easily add jitters or even create a beeswarm plot.

```
In [134...] sns.stripplot(x='Origin', y='Acceleration', data=cars, hue='Origin')
```

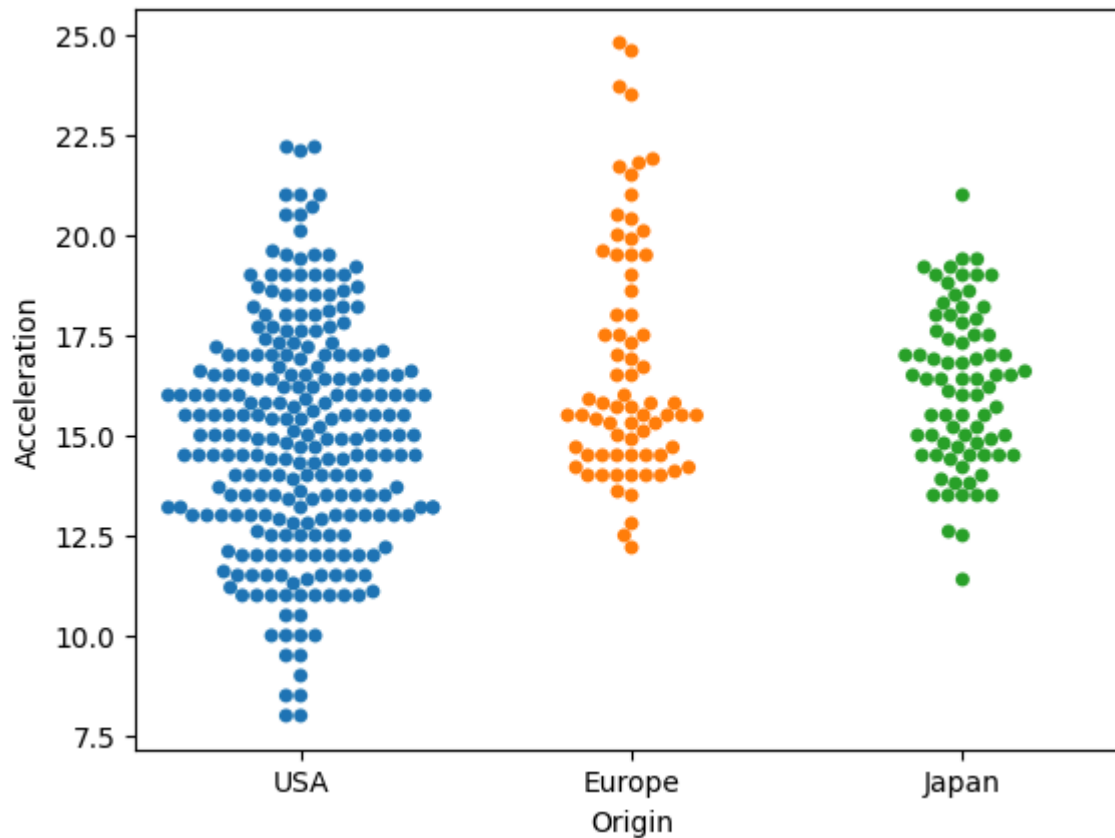
```
Out[134...] <Axes: xlabel='Origin', ylabel='Acceleration'>
```



Seems like European cars tend to have good acceleration. Let's look at the beeswarm plot, which is a pretty nice option for fairly small datasets.

```
In [135... sns.swarmplot(x='Origin', y='Acceleration', data=cars, hue='Origin')
```

```
Out[135... <Axes: xlabel='Origin', ylabel='Acceleration'>
```

Seaborn also allows you to use colors for another categorical variable. The option is `hue`.

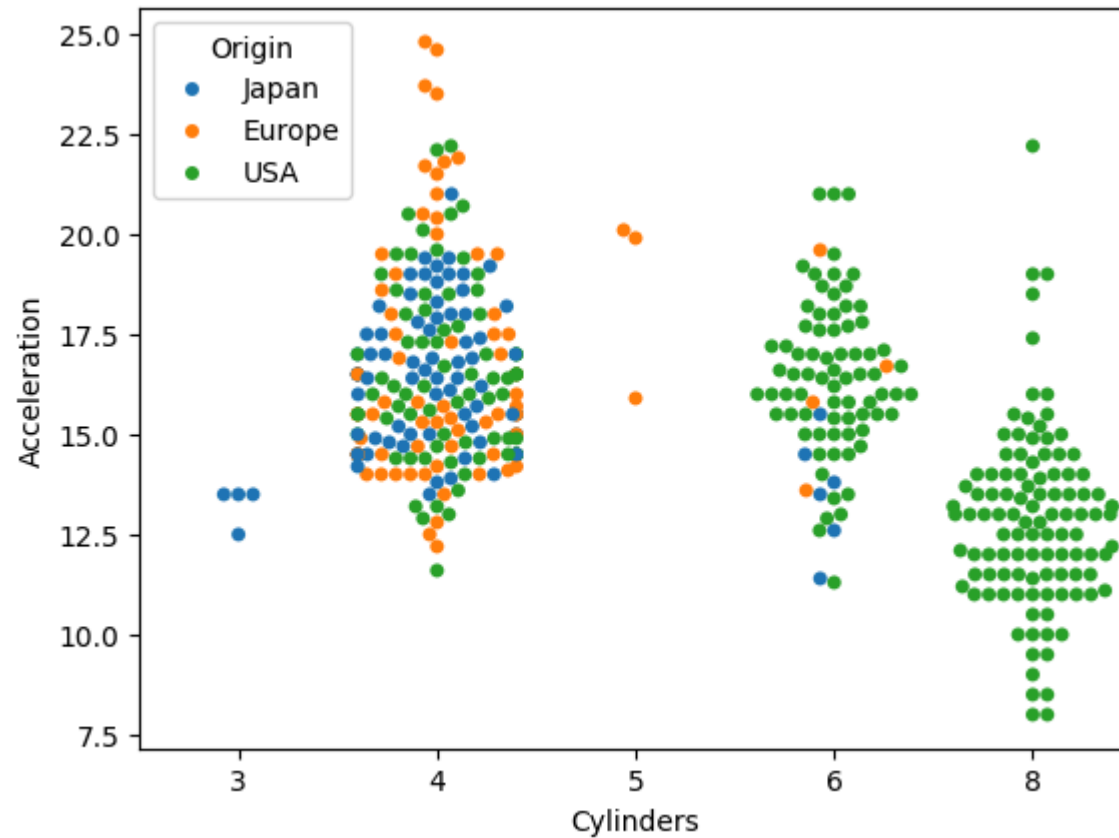
Q: can you create a beeswarm plot where the swarms are grouped by `Cylinders`, y-values are `Acceleration`, and colors represent the `Origin` ?

```
In [136...] sns.swarmplot(x='Cylinders', y='Acceleration', data=cars, hue='Origin')
```

```
d:\Masters-Projects\Semester2\ScientificDataVizualization\.venv\Lib\site-packages\seaborn\categorical.py:3399: UserWarning: 14.0% of the points cannot be placed; you may want to decrease the size of the markers or use stripplot.
  warnings.warn(msg, UserWarning)
```

```
Out[136...] <Axes: xlabel='Cylinders', ylabel='Acceleration'>
```

```
d:\Masters-Projects\Semester2\ScientificDataVizualization\.venv\Lib\site-packages\seaborn\categorical.py:3399: UserWarning: 17.9% of the points cannot be placed; you may want to decrease the size of the markers or use stripplot.  
warnings.warn(msg, UserWarning)
```

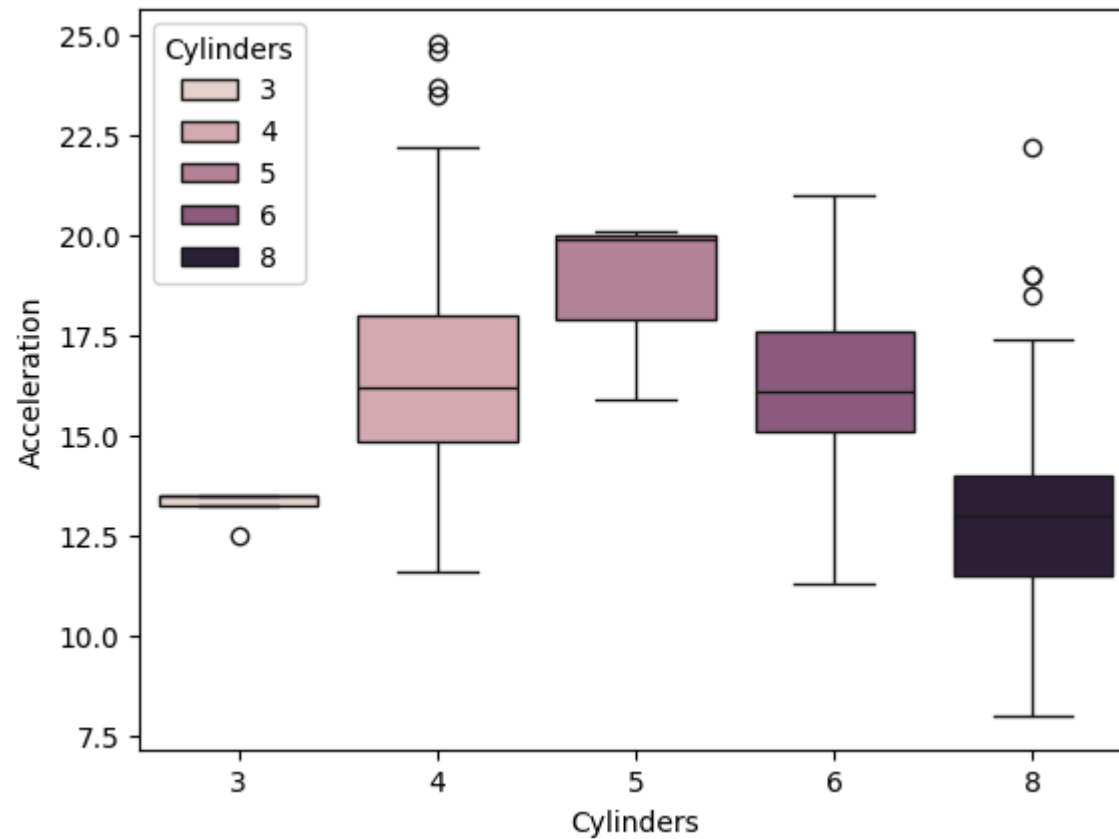


And of course you can create box plots too.

****Q:** Create boxplots to show the relationships between `Cylinders` and `Acceleration`. ******

```
In [137...] sns.boxplot(x='Cylinders', y='Acceleration', data=cars, hue='Cylinders')
```

```
Out[137...] <Axes: xlabel='Cylinders', ylabel='Acceleration'>
```



Altair basics

With `altair`, you're thinking in terms of a whole dataframe, rather than vectors for x or vectors for y. Passing the dataset to `Chart` creates an empty plot. If you try to run `alt.Chart(cars)`, it will complain. You need to say what's the visual encoding of the data.

In [138... `alt.Chart(cars)`

```

-----
SchemaValidationError                                Traceback (most recent call last)
File d:\Masters-Projects\Semester2\ScientificDataVizualization\.venv\Lib\site-packages\altair\vegalite\v5\api.py:4033, in Chart.to_dict(self, validate, format, ignore, context)
    4031     copy.data = core.InlineData(values=[{}])
    4032     return super(Chart, copy).to_dict(**kwds)
-> 4033 return super().to_dict(**kwds)

File d:\Masters-Projects\Semester2\ScientificDataVizualization\.venv\Lib\site-packages\altair\vegalite\v5\api.py:2004, in TopLevelMixin.to_dict(self, validate, format, ignore, context)
    2001 # remaining to_dict calls are not at top level
    2002 context["top_level"] = False
-> 2004 vegalite_spec: Any = _top_schema_base(super(TopLevelMixin, copy)).to_dict(
    2005     validate=validate, ignore=ignore, context=dict(context, pre_transform=False)
    2006 )
    2008 # TODO: following entries are added after validation. Should they be validated?
    2009 if is_top_level:
    2010     # since this is top-level we add $schema if it's missing

File d:\Masters-Projects\Semester2\ScientificDataVizualization\.venv\Lib\site-packages\altair\utils\schemapi.py:1169, in SchemaBase.to_dict(self, validate, ignore, context)
    1167     self.validate(result)
    1168     except jsonschema.ValidationError as err:
-> 1169     raise SchemaValidationError(self, err) from None
    1170 return result

SchemaValidationError: '{"data": {"name": "data-583e73726c1545c56c203344161a975c"}}' is an invalid value.

'mark' is a required property

```

Out[138... alt.Chart(...)

Note: If the altair plots don't show properly, use one of the following lines depending on your environment. Also check out the troubleshooting document [here](#).

```

In [139... #alt.renderers.enable('notebook')
alt.renderers.enable('jupyterlab')
#alt.renderers.enable('default')

```

Out[139... `RendererRegistry.enable('jupyterlab')`

```
In [140... alt.Chart(cars).mark_point()
```

Out[140... `<VegaLite 5 object>`

If you see this message, it means the renderer has not been properly enabled for the frontend that you are using. For more information, see https://altair-viz.github.io/user_guide/display_frontends.html#troubleshooting

So you just see one *point*. But actually this is not a single point. This is every row of the dataset represented as a point at the same location. Because there is no specification about where to put the points, it simply draws everything on top of each other. Let's specify how to spread them across the horizontal axis.

```
In [141... alt.Chart(cars).mark_point().encode(  
    x='Acceleration',  
)
```

Out[141... `<VegaLite 5 object>`

If you see this message, it means the renderer has not been properly enabled for the frontend that you are using. For more information, see https://altair-viz.github.io/user_guide/display_frontends.html#troubleshooting

This is called the "short form", and it is a simplified version of the "long form", while the long form allows more fine tuning. For this plot, they are equivalent:

```
In [142... alt.Chart(cars).mark_point().encode(  
    x=alt.X('Acceleration')  
)
```

Out[142... `<VegaLite 5 object>`

If you see this message, it means the renderer has not been properly enabled for the frontend that you are using. For more information, see https://altair-viz.github.io/user_guide/display_frontends.html#troubleshooting

There is another nice mark called `tick` :

In []:

```
In [143... alt.Chart(cars).mark_tick().encode(  
    x='Acceleration',  
    )
```

Out[143... <VegaLite 5 object>

If you see this message, it means the renderer has not been properly enabled for the frontend that you are using. For more information, see https://altair-viz.github.io/user_guide/display_frontends.html#troubleshooting

In `altair`, histogram is not a special type of visualization, but simply a plot with bars where a variable is binned and a counting aggregation function is used.

```
In [144... alt.Chart(cars).mark_bar().encode(  
    x=alt.X('Acceleration', bin=True),  
    y='count()'  
    )
```

Out[144... <VegaLite 5 object>

If you see this message, it means the renderer has not been properly enabled for the frontend that you are using. For more information, see https://altair-viz.github.io/user_guide/display_frontends.html#troubleshooting

Q: can you create a 2D scatterplot with `Acceleration` and `Horsepower` ? Use `Origin` for the colors.

```
In [145... chart = alt.Chart(cars).mark_point().encode(  
    x=alt.X('Horsepower:Q', title='Horsepower'),  
    y=alt.Y('Acceleration:Q', title='Acceleration'),  
    color=alt.Color('Origin:N', legend=alt.Legend(title='Origin'))  
    )  
chart
```

Out[145... <VegaLite 5 object>

If you see this message, it means the renderer has not been properly enabled for the frontend that you are using. For more information, see https://altair-viz.github.io/user_guide/display_frontends.html#troubleshooting

Because altair/vega-lite/vega are essentially drawing the chart using javascript (and D3.js), it is very easy to export it on the web. Probably the simplest way is just exporting it into an HTML file: https://altair-viz.github.io/getting_started/starting.html#publishing-your-visualization

Save the chart to m07.html and upload it too.

```
In [146... chart.save("m07.html")
```

Part III - Histogram and CDF

A deep dive into Histogram and boxplot.

```
In [147... import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
import altair as alt
import pandas as pd
```

```
In [148... import matplotlib
matplotlib.__version__
```

Out[148... '3.10.1'

The tricky histogram with pre-counted data

Let's consider the table

Hours	Frequency
0-1	4,300
1-3	6,900
3-5	4,900
5-10	2,000
10-24	2,100

You can draw a histogram by just providing bins and counts instead of a list of numbers. So, let's try that.

```
In [149... bins = [0, 1, 3, 5, 10, 24]
data = {0.5: 4300, 2: 6900, 4: 4900, 7: 2000, 15: 2100}
```

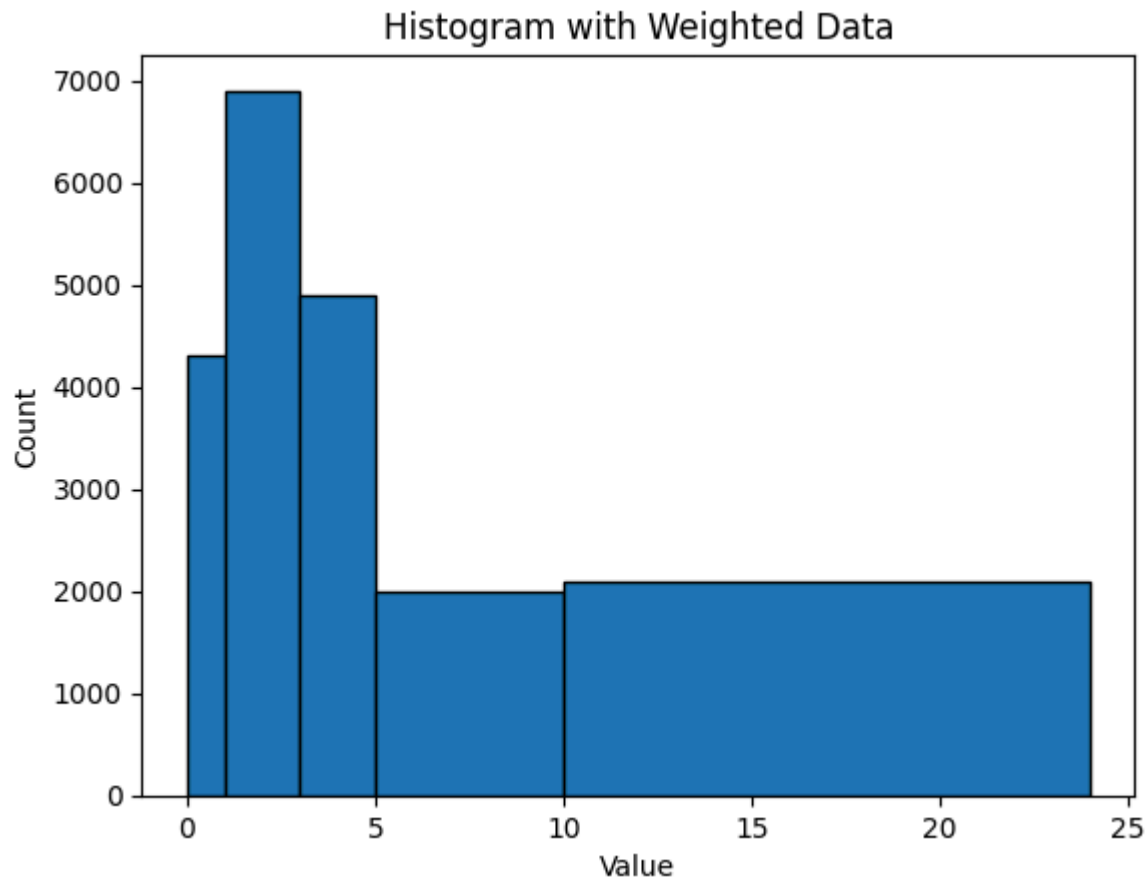
```
In [150... data.keys()
```

```
Out[150... dict_keys([0.5, 2, 4, 7, 15])
```

Q: Draw histogram using this data. Useful query: [Google search: matplotlib histogram pre-counted](#)

```
In [151... # draw a histogram with weighted data.
values= list(data.keys())
weights = list(data.values())

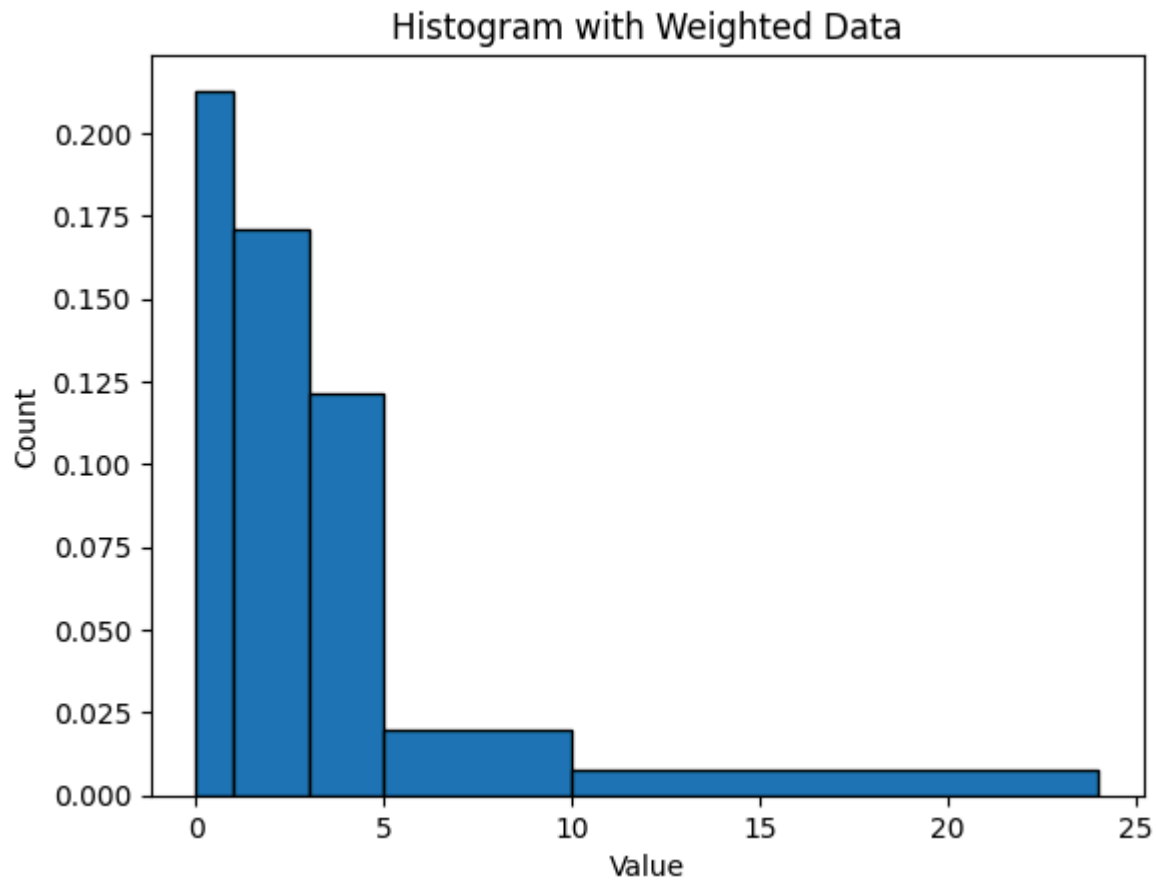
plt.hist(values, bins=bins, weights=weights, edgecolor='black')
plt.xlabel('Value')
plt.ylabel('Count')
plt.title('Histogram with Weighted Data')
plt.show()
```

As you can see, the **default histogram does not normalize with binwidth and simply shows the counts!** This can be very misleading if you are working with variable bin width (e.g. logarithmic bins). So please be mindful about histograms when you work with variable bins.

Q: You can fix this by using the `density` option.

```
In [152... plt.hist(values, bins=bins, weights=weights, edgecolor='black', density=True)
plt.xlabel('Value')
plt.ylabel('Count')
plt.title('Histogram with Weighted Data')
plt.show()
```



Let's use an actual dataset

```
In [153... import vega_datasets
```

```
In [154... vega_datasets.__version__
```


```
Out[154... '0.9.0'
```

Note: Please check your `vega_datasets` version using `vega_datasets.__version__`. If you have a version lower than `0.9.0`, you will need to check the column names in `movies.head()` and update it accordingly in the code cells below.

```
In [155... movies = vega_datasets.data.movies()
movies.head()
```

Out[155...

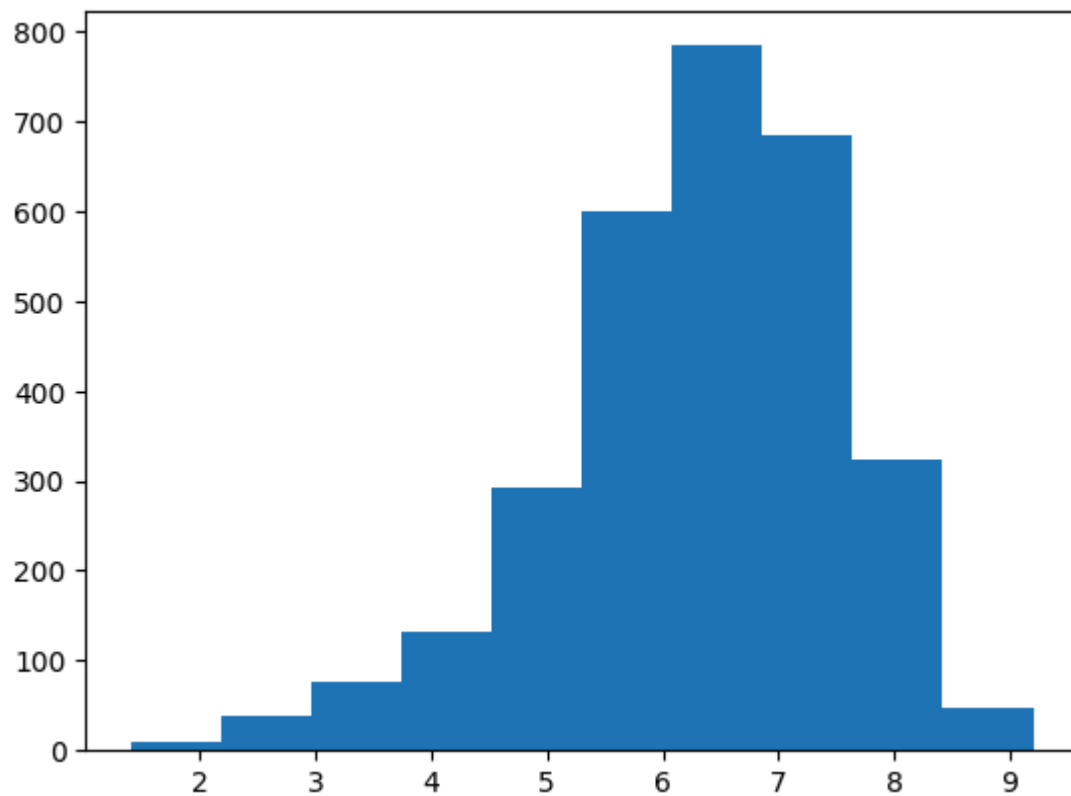
	Title	US_Gross	Worldwide_Gross	US_DVD_Sales	Production_Budget	Release_Date	MPAA_Rating	Running_Time_min	Distributor
0	The Land Girls	146083.0	146083.0	NaN	8000000.0	Jun 12 1998	R	NaN	Gramercy
1	First Love, Last Rites	10876.0	10876.0	NaN	300000.0	Aug 07 1998	R	NaN	Strand
2	I Married a Strange Person	203134.0	203134.0	NaN	250000.0	Aug 28 1998	None	NaN	Lionsgate
3	Let's Talk About Sex	373615.0	373615.0	NaN	300000.0	Sep 11 1998	None	NaN	Fine Line
4	Slam	1009819.0	1087521.0	NaN	1000000.0	Oct 09 1998	R	NaN	Trimark



Let's plot the histogram of IMDB ratings.

```
In [156... plt.hist(movies['IMDB_Rating'])
```

```
Out[156... (array([ 9., 39., 76., 133., 293., 599., 784., 684., 323., 48.]),
array([1.4 , 2.18, 2.96, 3.74, 4.52, 5.3 , 6.08, 6.86, 7.64, 8.42, 9.2 ]),
<BarContainer object of 10 artists>)
```



Did you get an error or a warning? What's going on?

The problem is that the column contains `NaN` (Not a Number) values, which represent missing data points. The following command checks whether each value is a `NaN` and returns the result.

```
In [157... movies['IMDB_Rating'].isna()
```

```
Out[157... 0      False
          1      False
          2      False
          3       True
          4      False
          ...
        3196     False
        3197       True
        3198     False
        3199     False
        3200     False
        Name: IMDB_Rating, Length: 3201, dtype: bool
```

As you can see there are a bunch of missing rows. You can count them.

```
In [158... sum(movies['IMDB_Rating'].isna())
```

```
Out[158... 213
```

or drop them.

```
In [159... IMDB_ratings_nan_dropped = movies['IMDB_Rating'].dropna()
len(IMDB_ratings_nan_dropped)
```

```
Out[159... 2988
```

```
In [160... 213 + 2988
```

```
Out[160... 3201
```

The `dropna` can be applied to the dataframe too.

Q: drop rows from `movies` dataframe where either `IMDB_Rating` or `IMDB_Votes` is `NaN`.

```
In [161... movies.dropna(subset=['IMDB_Rating', 'IMDB_Votes'], inplace=True)
```

```
In [162... # Both should be zero.
print(sum(movies['IMDB_Rating'].isna()), sum(movies['IMDB_Votes'].isna()))
```

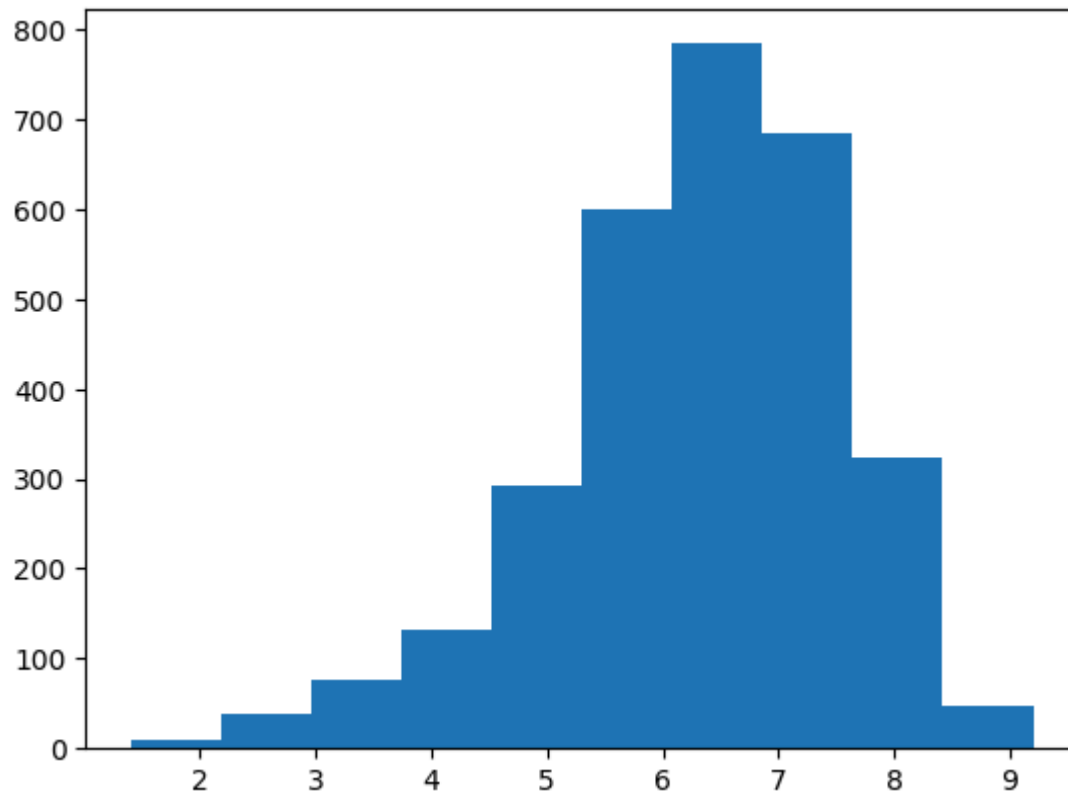
0 0

How does `matplotlib` decide the bins? Actually `matplotlib`'s `hist` function uses `numpy`'s `histogram` function under the hood.

Q: Plot the histogram of movie ratings (`IMDB_Rating`) using the `plt.hist()` function.

```
In [163... plt.hist(movies['IMDB_Rating'])
```

```
Out[163... (array([ 9., 39., 76., 133., 293., 599., 784., 684., 323., 48.]),  
array([1.4 , 2.18, 2.96, 3.74, 4.52, 5.3 , 6.08, 6.86, 7.64, 8.42, 9.2 ]),  
<BarContainer object of 10 artists>)
```

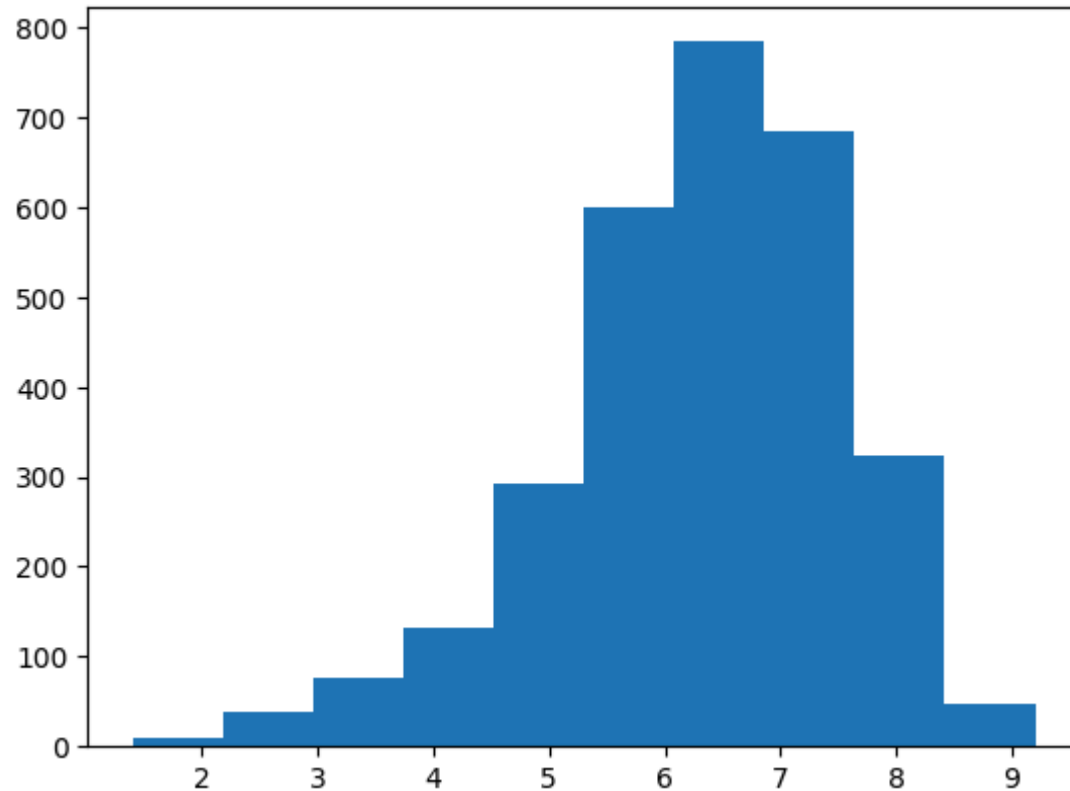


Have you noticed that this function returns three objects? Take a look at the documentation [here](#) to figure out what they are.

To get the returned three objects:

```
In [164... n_raw, bins_raw, patches = plt.hist(movies['IMDB_Rating'])
print(n_raw)
print(bins_raw)
```

```
[ 9.  39.  76. 133. 293. 599. 784. 684. 323.  48.]
[1.4  2.18 2.96 3.74 4.52 5.3  6.08 6.86 7.64 8.42 9.2 ]
```



Here, `n_raw` contains the values of histograms, i.e., the number of movies in each of the 10 bins. Thus, the sum of the elements in `n_raw` should be equal to the total number of movies.

Q: Test whether the sum of values in `n_raw` is equal to the number of movies in the `movies` dataset

```
In [165... print(n_raw.sum())
print(len(movies))
```

```
print(n_raw.sum == len(movies))
```

2988.0

2988

False

The second returned object (`bins_raw`) is a list containing the edges of the 10 bins: the first bin is [1.4, 2.18], the second [2.18, 2.96], and so on. What's the width of the bins?

```
In [166... np.diff(bins_raw)
```

```
Out[166... array([0.78, 0.78, 0.78, 0.78, 0.78, 0.78, 0.78, 0.78, 0.78, 0.78])
```

The width is same as the maximum value minus minimum value, divided by 10.

```
In [167... min_rating = min(movies['IMDB_Rating'])
max_rating = max(movies['IMDB_Rating'])
print(min_rating, max_rating)
print( (max_rating-min_rating) / 10 )
```

1.4 9.2

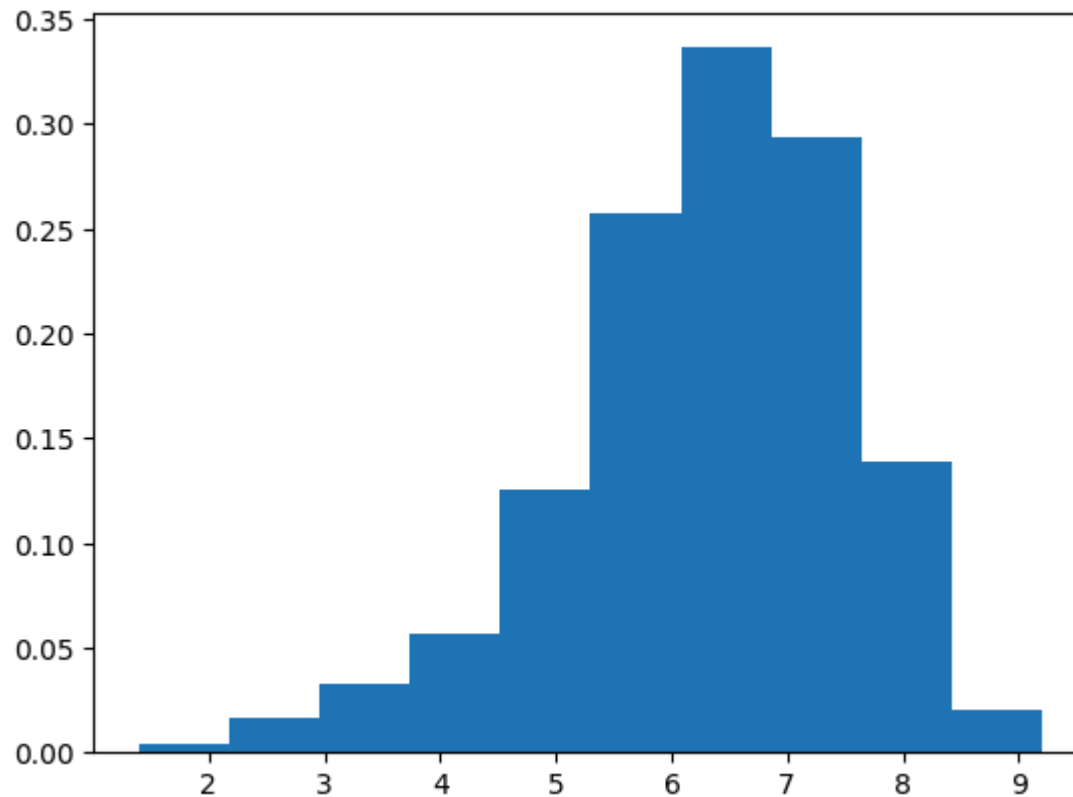
0.7799999999999999

Now, let's plot a normalized (density) histogram.

```
In [168... n, bins, patches = plt.hist(movies['IMDB_Rating'], density=True)
print(n)
print(bins)
```

```
[0.0038616  0.0167336  0.03260907 0.05706587 0.12571654 0.25701095
 0.33638829 0.29348162 0.13858854 0.0205952 ]
```

```
[1.4  2.18 2.96 3.74 4.52 5.3  6.08 6.86 7.64 8.42 9.2 ]
```

The ten bins do not change. But now `n` represents the density of the data inside each bin. In other words, the sum of the area of each bar will equal to 1.

Q: Can you verify this?

Hint: the area of each bar is calculated as height * width. You may get something like 0.9999999999999978 instead of 1.

In [169... `print(n.sum())`

1.2820512820512822

Anyway, these data generated from the `hist` function is calculated from `numpy`'s `histogram` function.

<https://docs.scipy.org/doc/numpy/reference/generated/numpy.histogram.html>

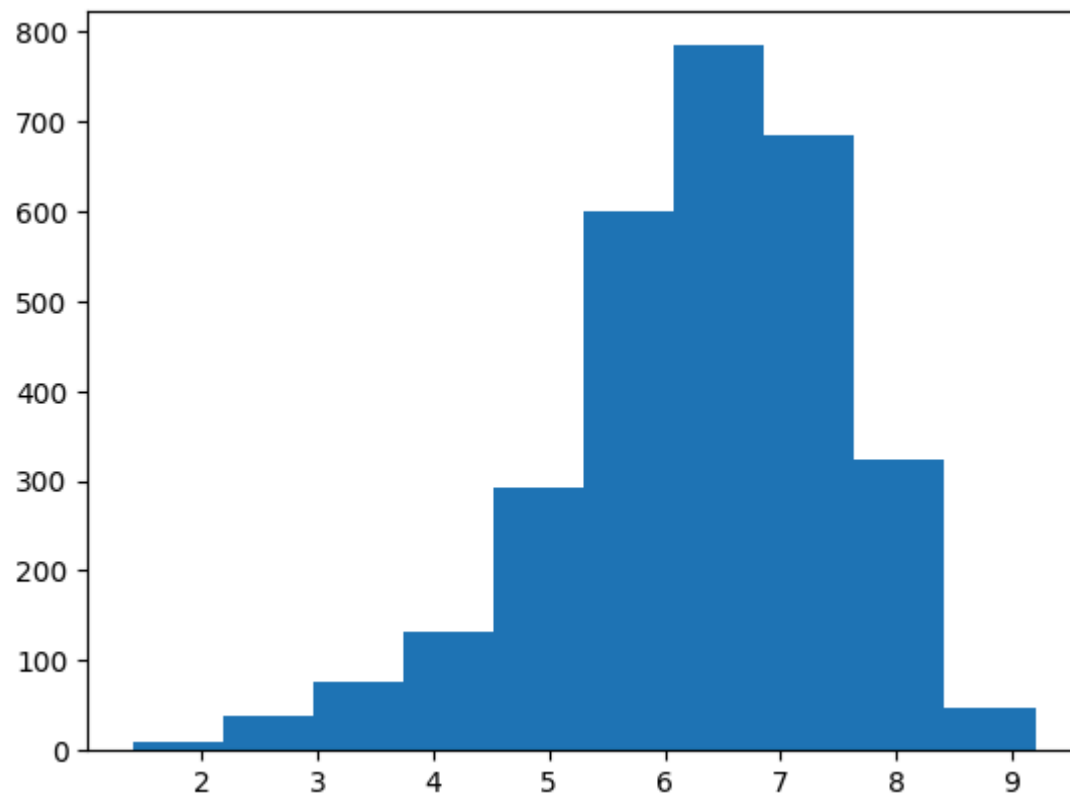
Note that the result of `np.histogram()` is same as that of `plt.hist()` .

```
In [170...] np.histogram(movies['IMDB_Rating'])
```

```
Out[170...] (array([ 9, 39, 76, 133, 293, 599, 784, 684, 323, 48]),  
            array([1.4 , 2.18, 2.96, 3.74, 4.52, 5.3 , 6.08, 6.86, 7.64, 8.42, 9.2 ]))
```

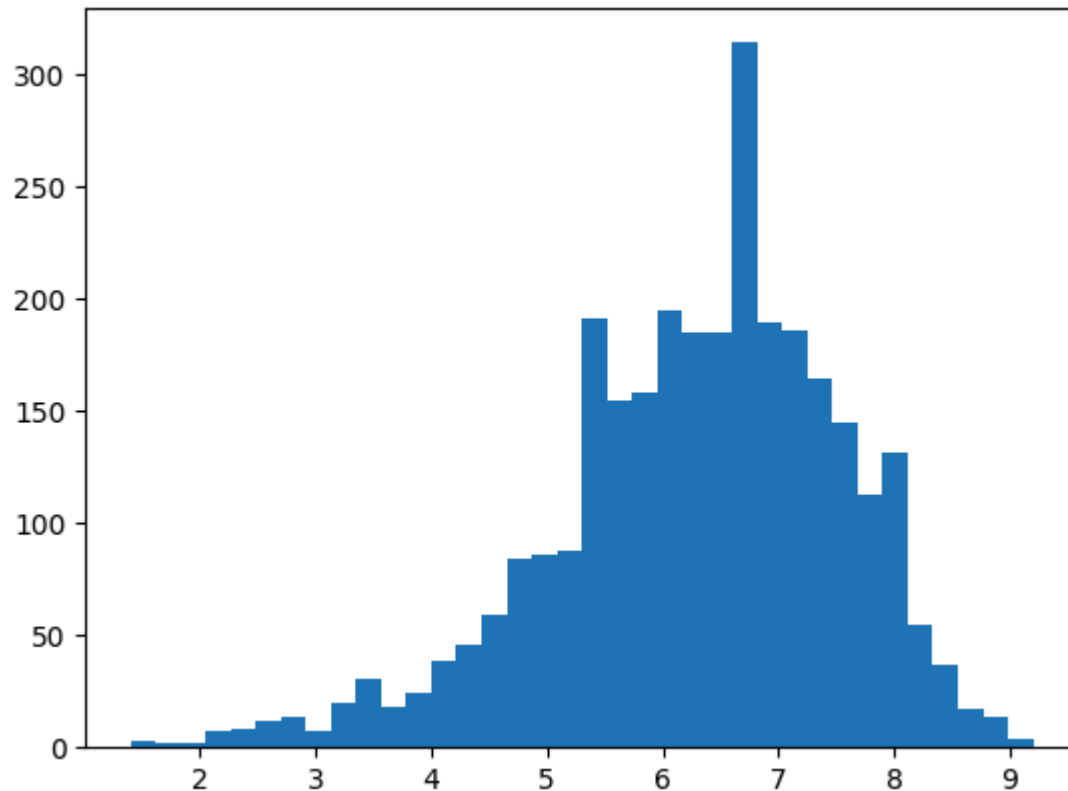
```
In [171...] plt.hist(movies['IMDB_Rating'])
```

```
Out[171...] (array([ 9., 39., 76., 133., 293., 599., 784., 684., 323., 48.]),  
            array([1.4 , 2.18, 2.96, 3.74, 4.52, 5.3 , 6.08, 6.86, 7.64, 8.42, 9.2 ]),  
            <BarContainer object of 10 artists>)
```



If you look at the documentation, you can see that `numpy` uses simply 10 as the default number of bins. But you can set it manually or set it to be `auto`, which is the "Maximum of the `sturges` and `fd` estimators.". Let's try this `auto` option.

```
In [172... _ = plt.hist(movies['IMDB_Rating'], bins='auto')
```



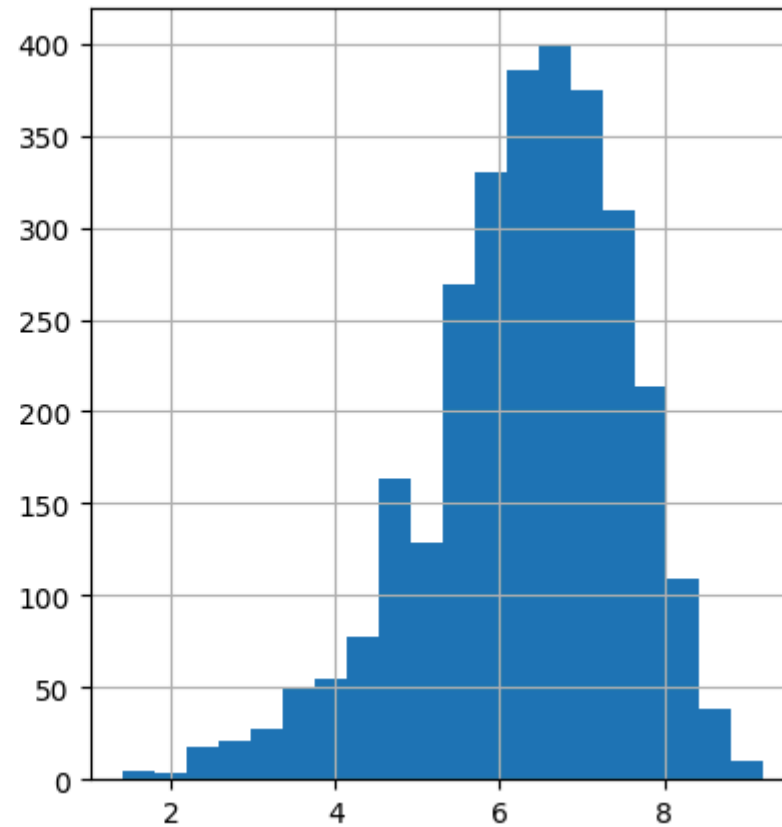
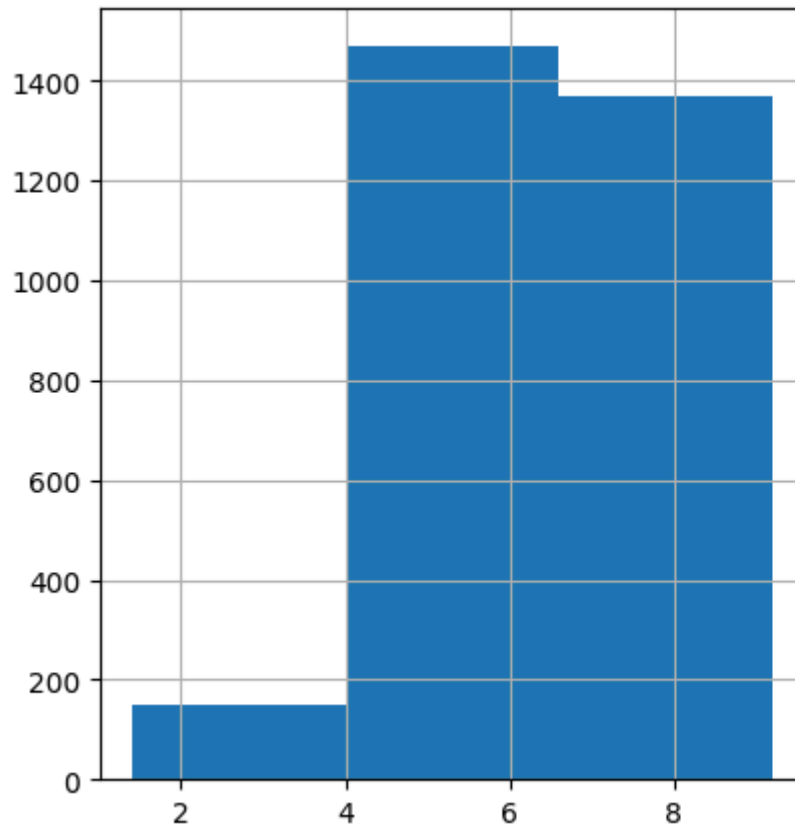
Consequences of the binning parameter

Let's explore the effect of bin size using small multiples. In `matplotlib`, you can use `subplot` to put multiple plots into a single figure.

For instance, you can do something like:

```
In [173... plt.figure(figsize=(10,5))
plt.subplot(1,2,1)
movies['IMDB_Rating'].hist(bins=3)
plt.subplot(1,2,2)
movies['IMDB_Rating'].hist(bins=20)
```

Out[173... <Axes: >

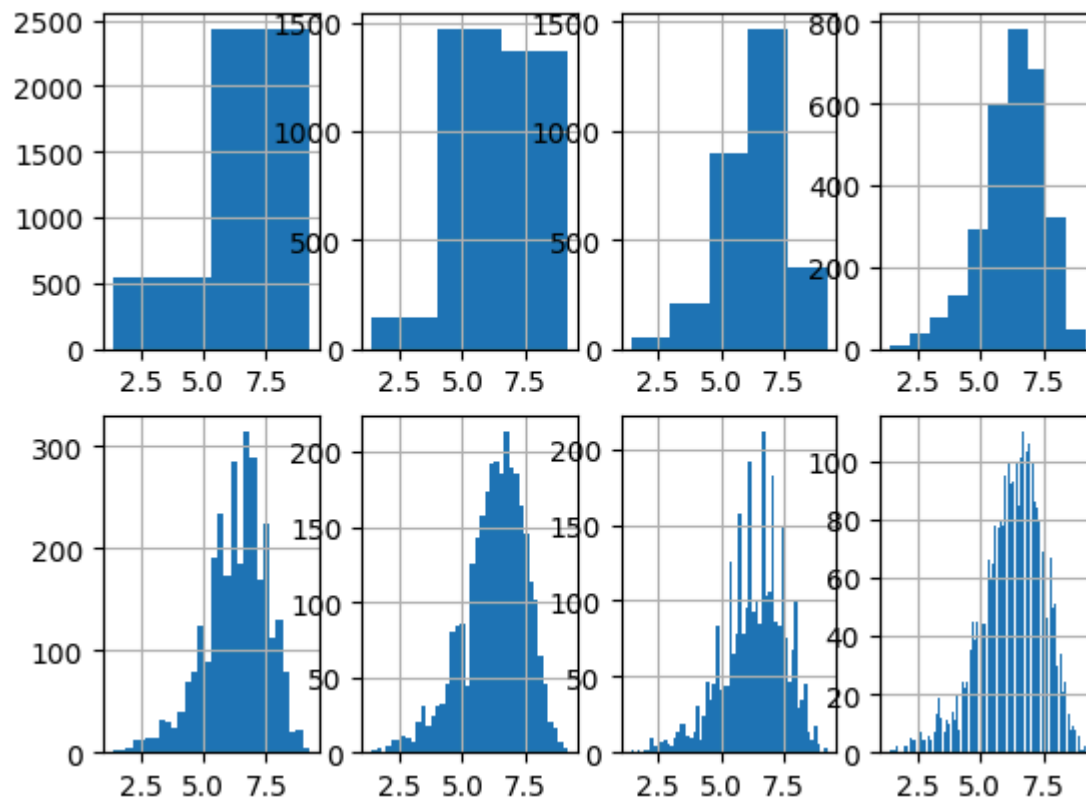


Q: create 8 subplots (2 rows and 4 columns) with the following binsizes .

```
In [174... nbins = [2, 3, 5, 10, 30, 40, 60, 100 ]
figsize = (18, 10)

for i in range(8):
```

```
plt.subplot(2,4 , i+1)
movies['IMDB_Rating'].hist(bins=nbins[i])
```



Do you see the issues with having too few bins or too many bins? In particular, do you notice weird patterns that emerge from `bins=30` ?

Q: Can you guess why do you see such patterns? What are the peaks and what are the empty bars? What do they tell you about choosing the binsize in histograms?

When setting the `bins=30`, we might notice that some bins have with very high counts(peaks) while many others are completely empty. This usually happens when the data is discrete or has only a few unique values.

- The peaks correspond to bins where the actual data points (or their weighted counts) lie. If the data has only a few distinct values, only the bins that cover those specific values will have counts.

- The empty bars are the bins that fall between these discrete values. Since no data points fall into these intervals, they show zero or near-zero counts.

What This Tells You About Choosing the Bin Size ?

- If choosing too many bins (like 30) relative to the distinct values in the data, we could potentially over-segment the data, leading to a "spiky" appearance where the true underlying pattern is obscured by the empty bins between the actual data values. It essentially exaggerates the discrete nature of the data.
- On the other hand, using too few bins can oversmooth the data, hiding important details or multimodal characteristics of the distribution

In []:

Formulae for choosing the number of bins.

We can manually choose the number of bins based on those formulae.

In [175...

```
N = len(movies)

plt.figure(figsize=(12,4))

# Sqrt
nbins = int(np.sqrt(N))

plt.subplot(1,3,1)
plt.title("SQRT, {} bins".format(nbins))
movies['IMDB_Rating'].hist(bins=nbins)

# Sturge's formula
nbins = int(np.ceil(np.log2(N) + 1))

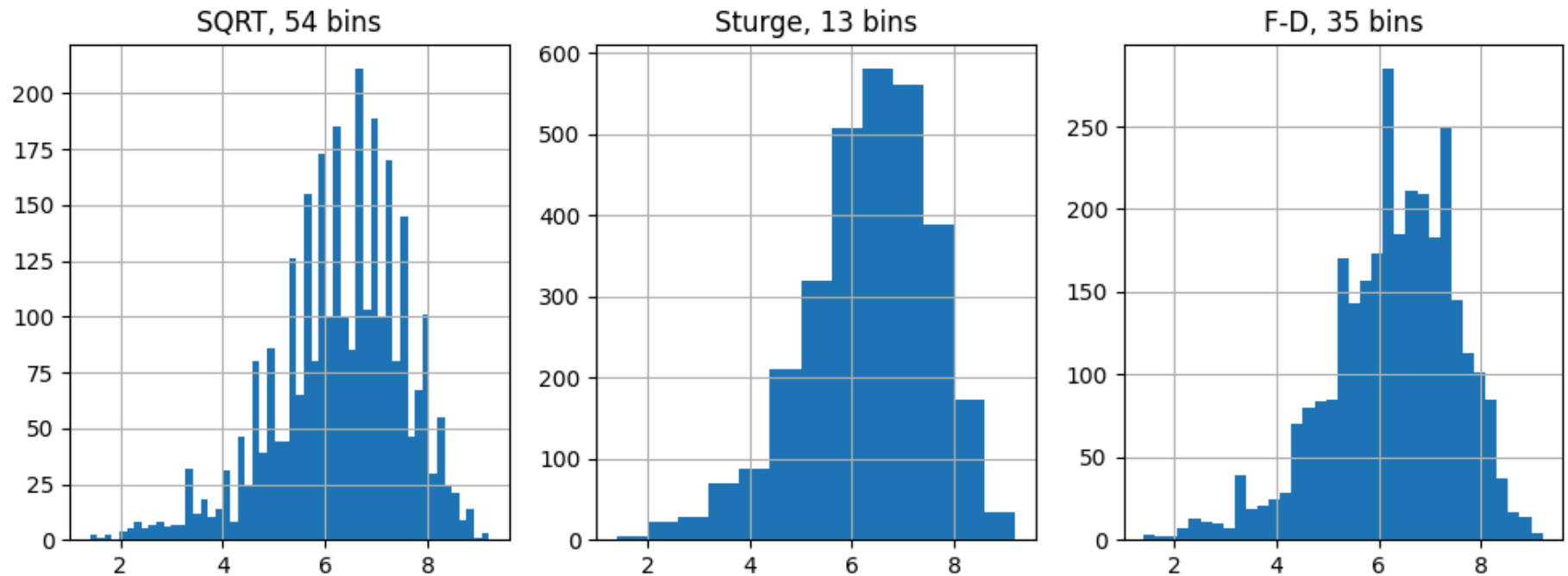
plt.subplot(1,3,2)
plt.title("Sturge, {} bins".format(nbins))
movies['IMDB_Rating'].hist(bins=nbins)

# Freedman-Diaconis
iqr = np.percentile(movies['IMDB_Rating'], 75) - np.percentile(movies['IMDB_Rating'], 25)
```

```
width = 2*iqr/np.power(N, 1/3)
nbins = int((max(movies['IMDB_Rating']) - min(movies['IMDB_Rating'])) / width)

plt.subplot(1,3,3)
plt.title("F-D, {} bins".format(nbins))
movies['IMDB_Rating'].hist(bins=nbins)
```

Out[175... <Axes: title={'center': 'F-D, 35 bins'}>



But we can also use built-in formulae too. Let's try all of them.

```
In [176... plt.figure(figsize=(20,4))

plt.subplot(161)
movies['IMDB_Rating'].hist(bins='fd')

plt.subplot(162)
movies['IMDB_Rating'].hist(bins='doane')
```

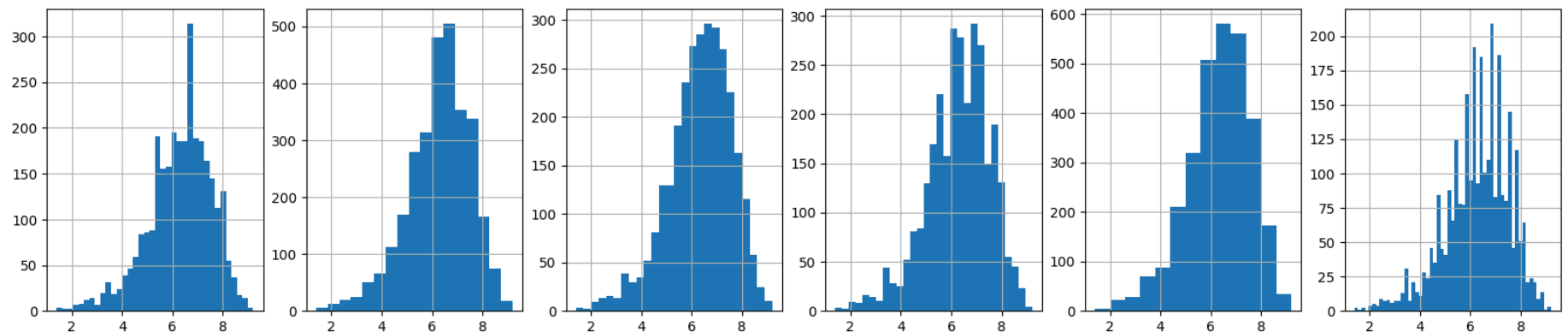
```
plt.subplot(163)
movies['IMDB_Rating'].hist(bins='scott')

plt.subplot(164)
movies['IMDB_Rating'].hist(bins='rice')

plt.subplot(165)
movies['IMDB_Rating'].hist(bins='sturges')

plt.subplot(166)
movies['IMDB_Rating'].hist(bins='sqrt')
```

Out[176... <Axes: >



Some are decent, but several of them tend to overestimate the good number of bins. As you have more data points, some of the formulae may overestimate the necessary number of bins. Particularly in our case, because of the precision issue, we shouldn't increase the number of bins too much.

Then, how should we choose the number of bins?

So what's the conclusion? use Scott's rule or Sturges' formula?

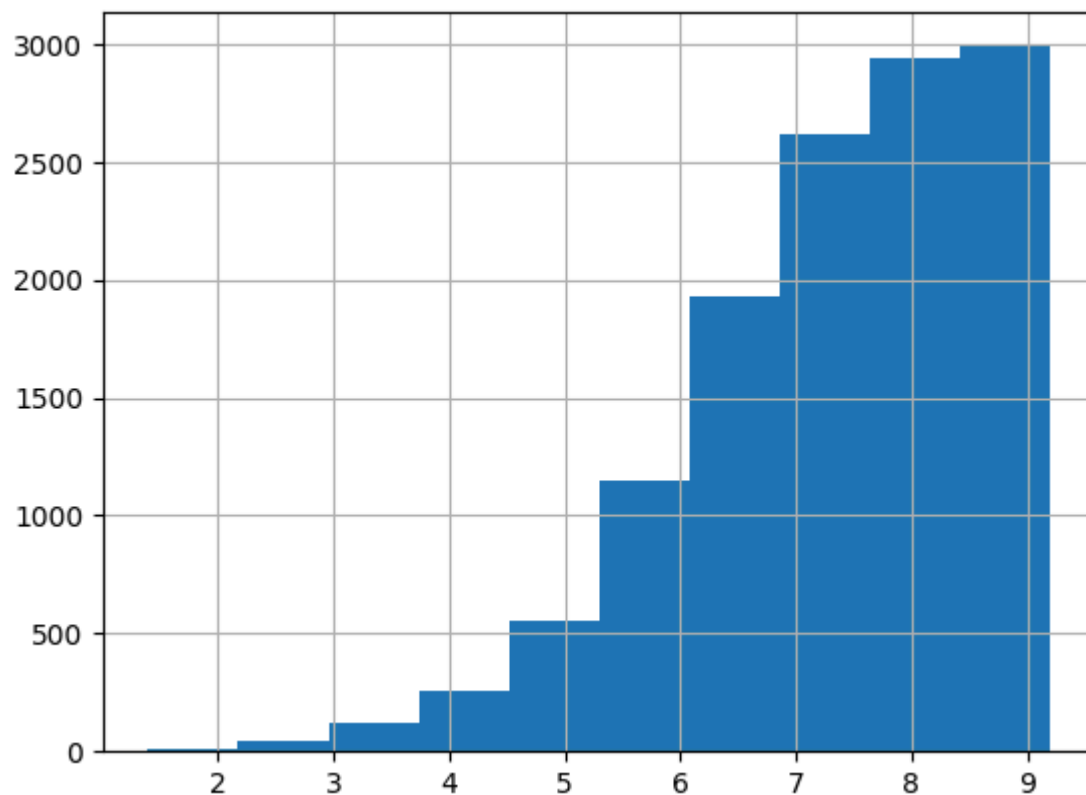
No, I think the take-away is that you **should understand how the inappropriate number of bins can mislead you** and you should **try multiple number of bins** to obtain the most accurate picture of the data. Although the 'default' may work in most cases, don't blindly trust it! Don't judge the distribution of a dataset based on a single histogram. Try multiple parameters to get the full picture!

CDF (Cumulative distribution function)

Drawing a CDF is easy. Because it's very common data visualization, histogram has an option called `cumulative`.

```
In [177... movies['IMDB_Rating'].hist(cumulative=True)
```

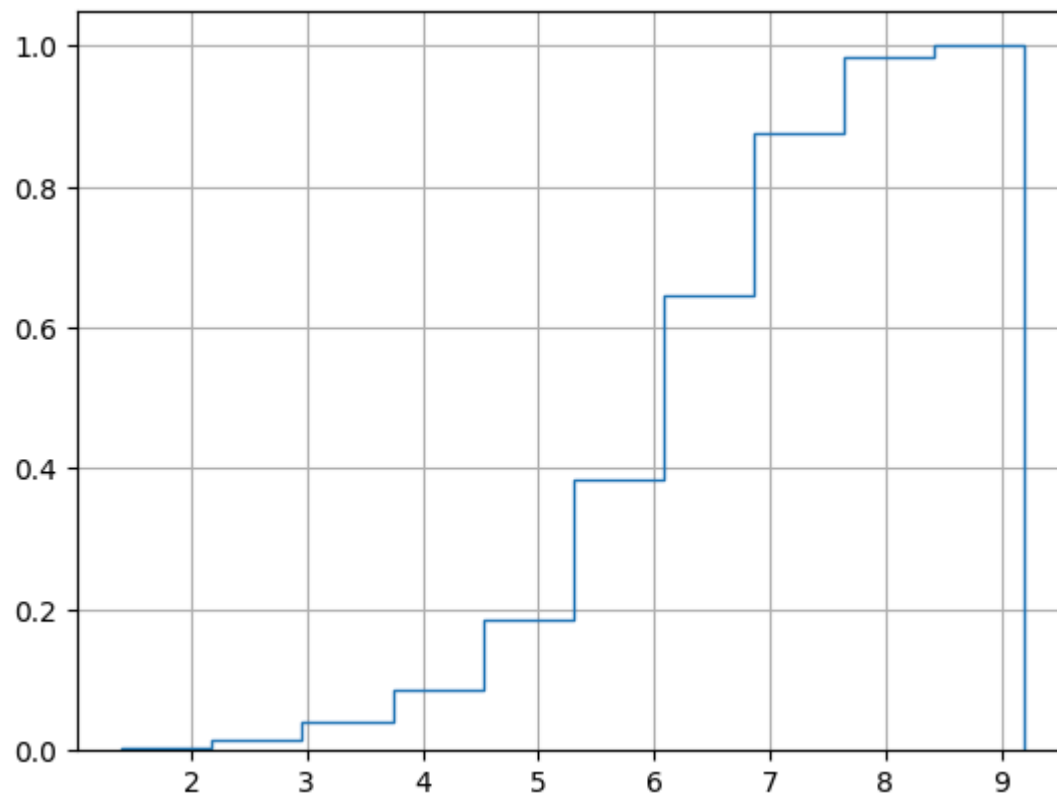
```
Out[177... <Axes: >
```



You can also combine with options such as `histtype` and `density`.

```
In [178... movies['IMDB_Rating'].hist(histtype='step', cumulative=True, density=True)
```

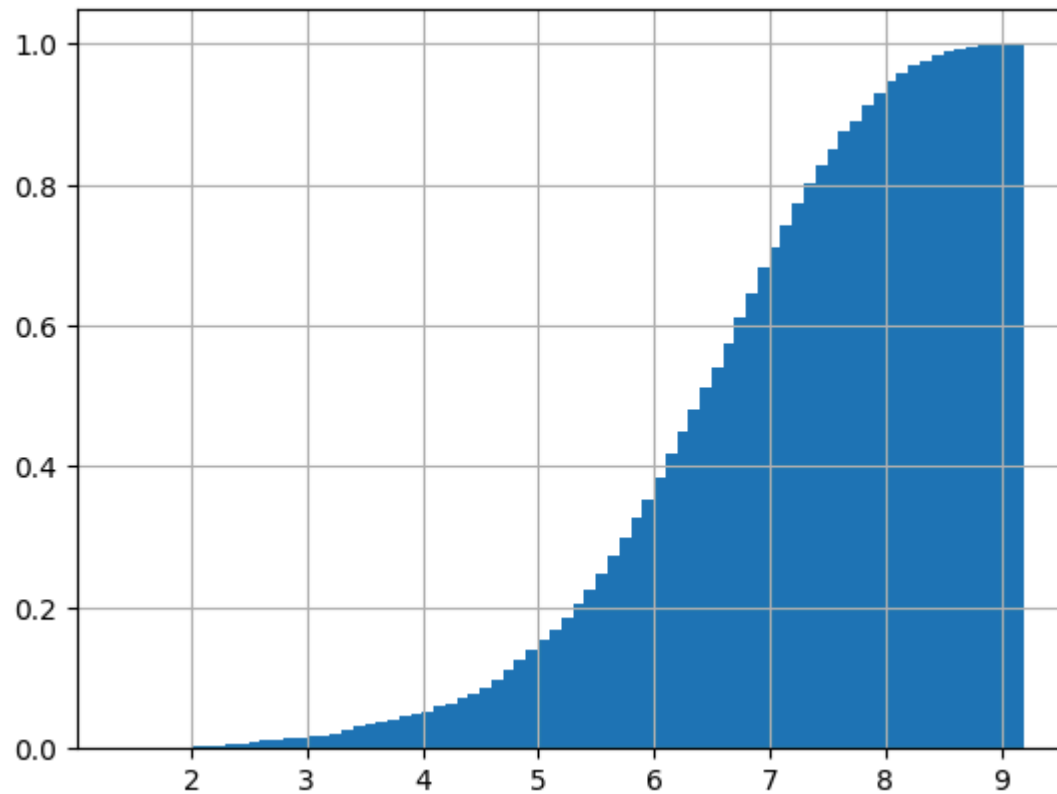
```
Out[178... <Axes: >
```



And increase the number of bins.

```
In [179... movies['IMDB_Rating'].hist(cumulative=True, density=True, bins=1000)
```

```
Out[179... <Axes: >
```



This method works fine. By increasing the number of bins, you can get a CDF in the resolution that you want. But let's also try it manually to better understand what's going on. First, we should sort all the values.

```
In [180... rating_sorted = movies['IMDB_Rating'].sort_values()  
rating_sorted.head()
```

```
Out[180... 1247    1.4  
         406    1.5  
        1754    1.6  
        1515    1.7  
        1590    1.7  
Name: IMDB_Rating, dtype: float64
```

We need to know the number of data points,

```
In [181... N = len(rating_sorted)
N
```

```
Out[181... 2988
```

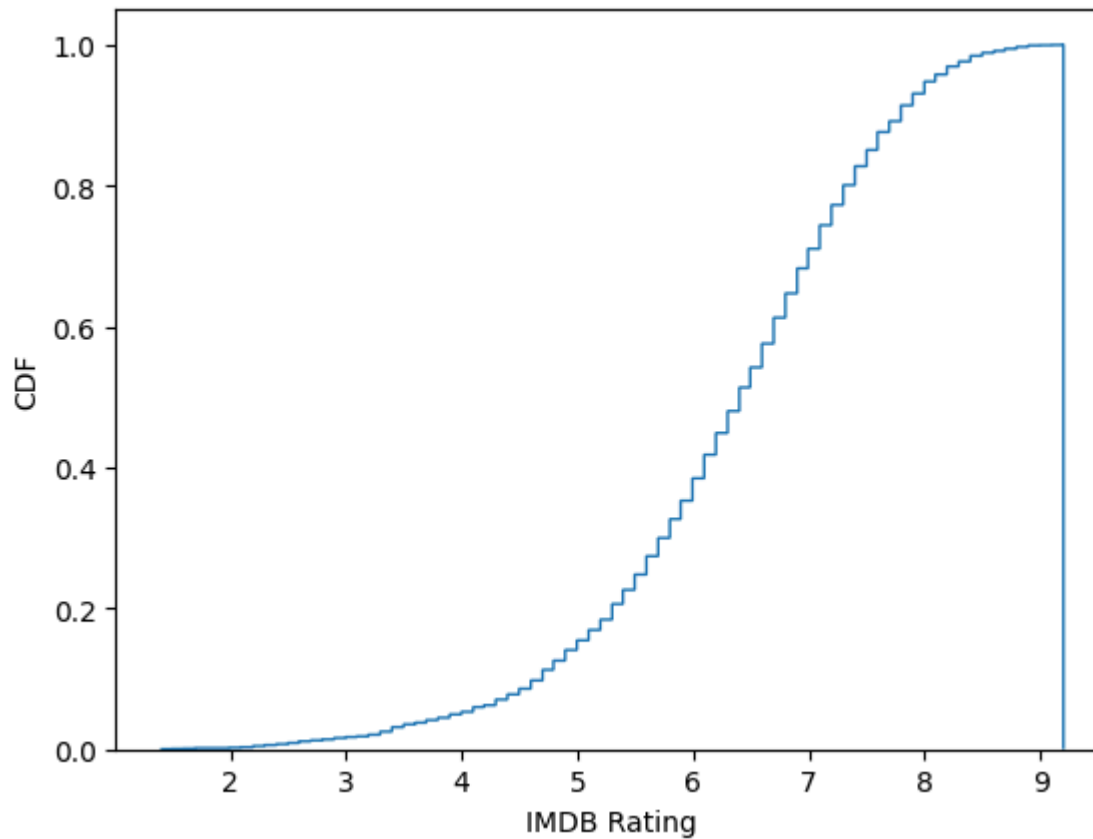
And I think this may be useful for you.

```
In [182... n = 50
np.linspace(1/n, 1.0, num=n)
```

```
Out[182... array([0.02, 0.04, 0.06, 0.08, 0.1 , 0.12, 0.14, 0.16, 0.18, 0.2 , 0.22,
        0.24, 0.26, 0.28, 0.3 , 0.32, 0.34, 0.36, 0.38, 0.4 , 0.42, 0.44,
        0.46, 0.48, 0.5 , 0.52, 0.54, 0.56, 0.58, 0.6 , 0.62, 0.64, 0.66,
        0.68, 0.7 , 0.72, 0.74, 0.76, 0.78, 0.8 , 0.82, 0.84, 0.86, 0.88,
        0.9 , 0.92, 0.94, 0.96, 0.98, 1.  ])
```

Q: now you're ready to draw a proper CDF. Draw the CDF plot of this data.

```
In [183... plt.hist(rating_sorted, bins=1000, density=True, cumulative=True, histtype='step')
plt.xlabel("IMDB Rating")
plt.ylabel("CDF")
plt.show()
```



A bit more histogram with altair

As you may remember, you can get a pandas dataframe from `vega_datasets` package and use it to create visualizations. But, if you use `altair`, you can simply pass the URL instead of the actual data.

```
In [184...] vega_datasets.data.movies.url
```

```
Out[184...] 'https://cdn.jsdelivr.net/npm/vega-datasets@v1.29.0/data/movies.json'
```

```
In [185...] # Choose based on your environment  
#alt.renderers.enable('notebook')
```

```
alt.renderers.enable('jupyterlab')
#alt.renderers.enable('default')
```

Out[185... `RendererRegistry.enable('jupyterlab')`

As mentioned before, in `altair` histogram is not special. It is just a plot that use bars (`mark_bar()`) where X axis is defined by `IMDB_Rating` with bins (`bin=True`), and Y axis is defined by `count()` aggregation function.

```
In [186... alt.Chart(vega_datasets.data.movies.url).mark_bar().encode(
    alt.X("IMDB_Rating:Q", bin=True),
    alt.Y('count()')
)
```

Out[186... `<VegaLite 5 object>`

If you see this message, it means the renderer has not been properly enabled for the frontend that you are using. For more information, see https://altair-viz.github.io/user_guide/display_frontends.html#troubleshooting

Have you noted that it is `IMDB_Rating:Q` not `IMDB_Rating` ? This is a shorthand for

```
In [187... alt.Chart(vega_datasets.data.movies.url).mark_bar().encode(
    alt.X('IMDB_Rating', type='quantitative', bin=True),
    alt.Y(aggregate='count', type='quantitative')
)
```

Out[187... `<VegaLite 5 object>`

If you see this message, it means the renderer has not been properly enabled for the frontend that you are using. For more information, see https://altair-viz.github.io/user_guide/display_frontends.html#troubleshooting

In altair, you want to specify the data types using one of the four categories: quantitative, ordinal, nominal, and temporal. https://altair-viz.github.io/user_guide/encoding.html#data-types

Although you can adjust the bins in `altair`, it does not encourage you to set the bins directly. For instance, although there is `step` parameter that directly sets the bin size, there are parameters such as `maxbins` (maximum number of bins) or `minstep` (minimum allowable step size), or `nice` (attempts to make the bin boundaries more human-friendly), that encourage you not to specify the bins directly.

```
In [188... from altair import Bin

alt.Chart(vega_datasets.data.movies.url).mark_bar().encode(
    alt.X("IMDB_Rating:Q", bin=Bin(step=0.09)),
    alt.Y('count()')
)
```

Out[188... <Vegalite 5 object>

If you see this message, it means the renderer has not been properly enabled for the frontend that you are using. For more information, see https://altair-viz.github.io/user_guide/display_frontends.html#troubleshooting

```
In [189... alt.Chart(vega_datasets.data.movies.url).mark_bar().encode(
    alt.X("IMDB_Rating:Q", bin=Bin(nice=True, maxbins=20)),
    alt.Y('count()')
)
```

Out[189... <Vegalite 5 object>

If you see this message, it means the renderer has not been properly enabled for the frontend that you are using. For more information, see https://altair-viz.github.io/user_guide/display_frontends.html#troubleshooting

Composing charts in altair

`altair` has a very nice way to compose multiple plots. Two histograms side by side? just do the following.

```
In [190... chart1 = alt.Chart(vega_datasets.data.movies.url).mark_bar().encode(
    alt.X("IMDB_Rating:Q", bin=Bin(step=0.1)),
    alt.Y('count()')
).properties(
    width=300,
    height=150
)
```

```
)  
chart2 = alt.Chart(vega_datasets.data.movies.url).mark_bar().encode(  
    alt.X("IMDB_Rating:Q", bin=Bin(nice=True, maxbins=20)),  
    alt.Y('count()')  
)  
.properties(  
    width=300,  
    height=150  
)  
)
```

In [191... chart1 | chart2

Out[191... <Vegalite 5 object>

If you see this message, it means the renderer has not been properly enabled for the frontend that you are using. For more information, see https://altair-viz.github.io/user_guide/display_frontends.html#troubleshooting

In [192... alt.hconcat(chart1, chart2)

Out[192... <Vegalite 5 object>

If you see this message, it means the renderer has not been properly enabled for the frontend that you are using. For more information, see https://altair-viz.github.io/user_guide/display_frontends.html#troubleshooting

Vertical composition?

In [193... alt.vconcat(chart1, chart2)

Out[193... <Vegalite 5 object>

If you see this message, it means the renderer has not been properly enabled for the frontend that you are using. For more information, see https://altair-viz.github.io/user_guide/display_frontends.html#troubleshooting

In [194... chart1 & chart2

Out[194... <VegaLite 5 object>

If you see this message, it means the renderer has not been properly enabled for the frontend that you are using. For more information, see https://altair-viz.github.io/user_guide/display_frontends.html#troubleshooting

Shall we avoid some repetitions? You can define a *base* empty chart first and then assign encodings later when you put together multiple charts together. Here is an example: https://altair-viz.github.io/user_guide/compound_charts.html#repeated-charts

Use base chart to produce the chart above:

```
In [195... base = alt.Chart().mark_bar().encode(
    alt.X("IMDB_Rating:Q", bin=Bin(nice=True, maxbins=20)),
    alt.Y('count()')
).properties(
    width=300,
    height=150
)

chart = alt.vconcat(data=vega_datasets.data.movies.url)
for bin_param in [Bin(step=0.1), Bin(nice=True, maxbins=20)]:
    row = alt.hconcat()
    row |= base.encode(x=alt.X("IMDB_Rating:Q", bin=bin_param), y='count()')
    chart &= row

chart
```

Out[195... <VegaLite 5 object>

If you see this message, it means the renderer has not been properly enabled for the frontend that you are using. For more information, see https://altair-viz.github.io/user_guide/display_frontends.html#troubleshooting

Q: Using the base chart approach to create a 2x2 chart where the top row shows the two histograms of `IMDB_Rating` with `maxbins = 10` and `50` respectively, and the bottom row shows another two histograms of `IMDB_Votes` with `maxbins = 10` and `50`.

```
In [196... base = alt.Chart(data=vega_datasets.data.movies.url).properties(width=250, height=250)

# Define configurations for each chart in the grid.
```

```

# Each configuration specifies the field and the desired number of bins.
configs = [
    {"field": "IMDB_Rating", "maxbins": 10},
    {"field": "IMDB_Rating", "maxbins": 50},
    {"field": "IMDB_Votes", "maxbins": 10},
    {"field": "IMDB_Votes", "maxbins": 50},
]

# Create an empty list to hold each generated chart.
charts = []

# Loop over each configuration and generate a histogram.
for config in configs:
    chart = base.encode(
        alt.X(f'{config["field"]}:Q', bin=alt.Bin(maxbins=config["maxbins"]), title=config["field"]),
        alt.Y('count()', title='Count')
    ).mark_bar()
    charts.append(chart)

# Arrange the four charts into a 2x2 grid.
top_row = alt.hconcat(charts[0], charts[1])
bottom_row = alt.hconcat(charts[2], charts[3])
final_chart = alt.vconcat(top_row, bottom_row)

final_chart.display()

```

<VegaLite 5 object>

If you see this message, it means the renderer has not been properly enabled for the frontend that you are using. For more information, see https://altair-viz.github.io/user_guide/display_frontends.html#troubleshooting