

≡ 02. Tidy Data

What is Tidy Data?

In this course, it is expected that your data is organized in some kind of tidy format. In short, a [tidy dataset](#) is a tabular dataset where:

- each variable is a column
- each observation is a row
- each type of observational unit is a table

The first three images below depict a tidy dataset. This tidy dataset is in the field of healthcare and has two tables: one for patients (with their patient ID, name, and age) and one for treatments (with patient ID, what drug that patient is taking, and the dose of that drug).

Each variable forms a column

PATIENTS			TREATMENTS		
Patient ID	Name	Age	Patient ID	Drug	Dose
101	Juan Pérez	26	101	A	60
102	Jane Citizen	43	102	B	40
103	Kwasi Mensa	75	103	C	20

Each variable in a tidy dataset must have its own column

Each variable forms a column
Each observation forms a row

PATIENTS			TREATMENTS		
Patient ID	Name	Age	Patient ID	Drug	Dose
101	Juan Pérez	26	101	A	60
102	Jane Citizen	43	102	B	40
103	Kwasi Mensa	75	103	C	20

Each observation in a tidy dataset must have its own row

Each variable forms a column
Each observation forms a row
Each observational unit forms a table

PATIENTS			TREATMENTS		
Patient ID	Name	Age	Patient ID	Drug	Dose
101	Juan Pérez	26	101	A	60
102	Jane Citizen	43	102	B	40
103	Kwasi Mensa	75	103	C	20

Each observational unit in a tidy dataset must have its own table

The next image depicts the same data but in one representation of a non-tidy format (there are other possible non-tidy representations). The *Drug A*, *Drug B*, and *Drug C* columns should form one 'Drug' column, since this is one variable. The entire table should be separated into two tables: a patients table and a treatments table.

Each variable forms a column
Each observation forms a row
Each observational unit forms a table

Patient ID	Name	Age	Drug A	Drug B	Drug C
101	Juan Pérez	26	60	—	—
102	Jane Citizen	43	—	40	—
103	Kwasi Mensa	75	—	—	20

Only the second rule of tidy data is satisfied in this non-tidy representation of the above data: each observation forms a row

While the data provided to you in the course will all be tidy, in practice, you may need to perform tidying work before exploration. You should be comfortable with reshaping your data or perform transformations to split or combine features in your data, resulting in new data columns. These operations collectively are called *data-wrangling*.

This is also not to say that tidy data is the *only* useful form that data can take. In fact, as you work with a dataset, you might need to summarize it in a non-tidy form in order to generate appropriate visualizations. You'll see one example of this (bivariate plotting) in the next lesson, where categorical counts need to put into a matrix form in order to create a heat map.

Recommended Read

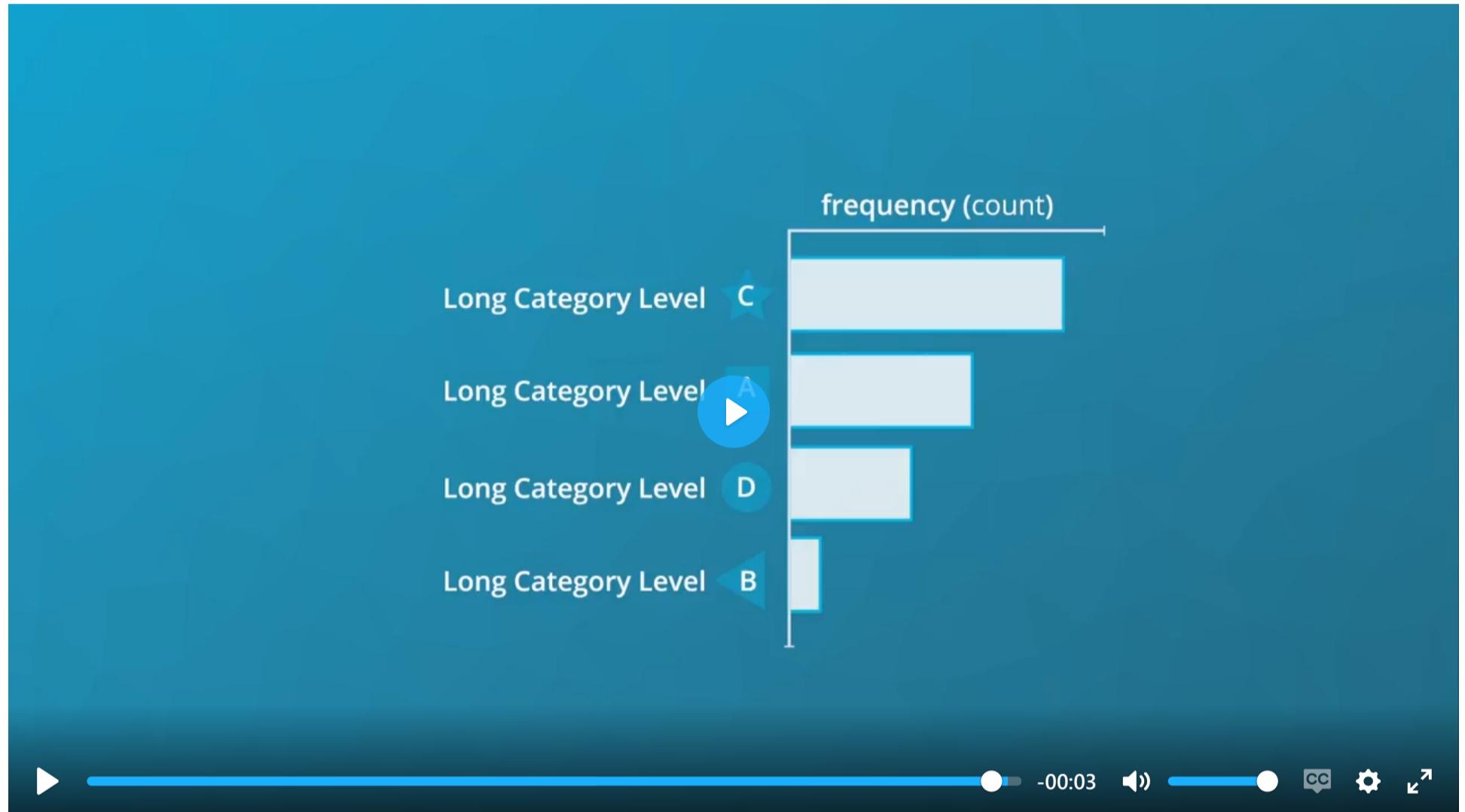
Refer to the [Data Wrangling with pandas Cheat Sheet](#) for a summary of functions helpful for data-wrangling.

Next Concept

≡ 03. Bar Charts

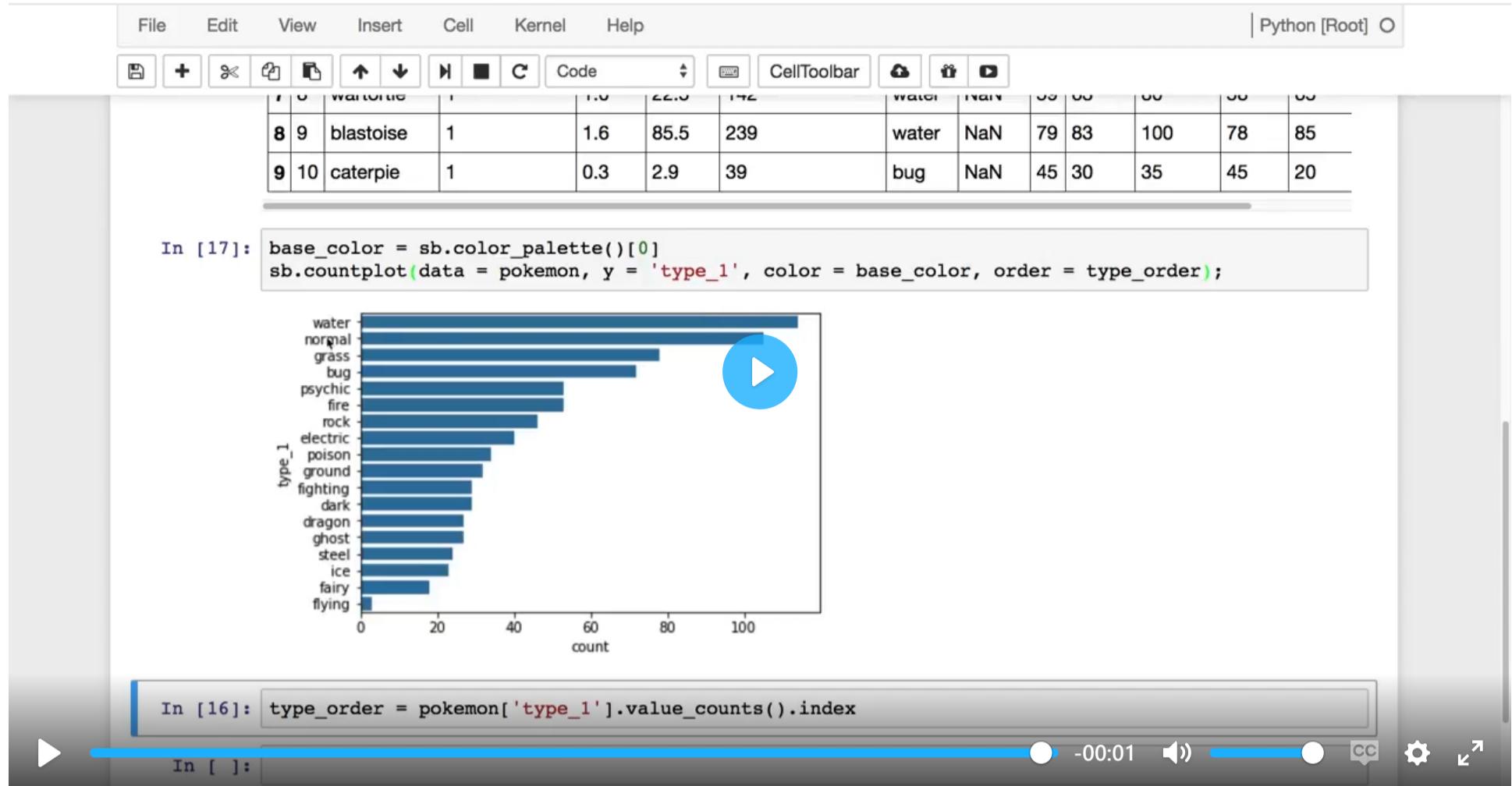
A **bar chart** depicts the distribution of a categorical variable. In a bar chart, each level of the categorical variable is depicted with a bar, whose height indicates the frequency of data points that take on that level.

L3 031 Bar Charts V3



Bar Chart - Demo

DataVis L3 03 V2



INSTRUCTOR NOTE:

The `pokemon.csv` file is available to download at the bottom of this page, though it is also present in the upcoming Jupyter notebook workspace.

Bar Chart using Seaborn

A basic bar chart of frequencies can be created through the use of seaborn's `countplot` function.

```
seaborn.countplot(*, x=None, y=None, data=None, order=None, orient=None, color=None)
```

We will see the usage of a few of the arguments of the `countplot()` function.

Example 1. Create a vertical bar chart using Seaborn, with default colors

```

# Necessary imports
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sb
%matplotlib inline

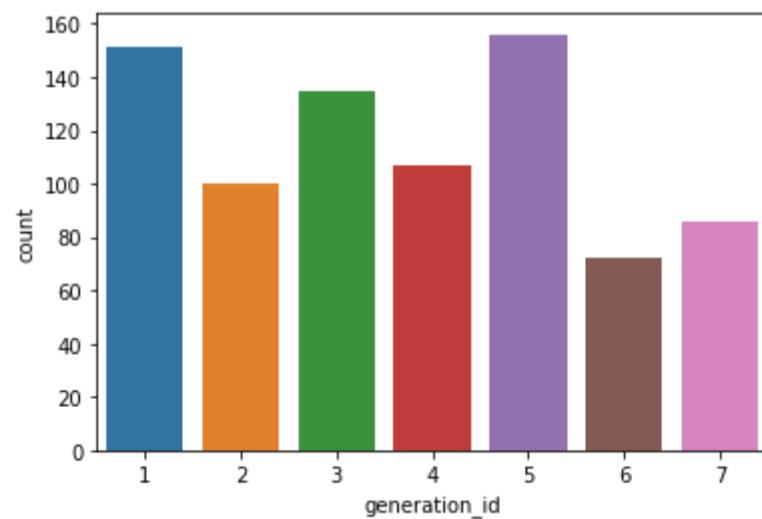
# Read the csv file, and check its top 10 rows
pokemon = pd.read_csv('pokemon.csv')
print(pokemon.shape)
pokemon.head(10)

# A semicolon (;) at the end of the statement will suppress printing the plotting information
sb.countplot(data=pokemon, x='generation_id');

```

(807, 14)

	id	species	generation_id	height	weight	base_experience	type_1	type_2	hp	attack	defense	speed	special-attack	special-defense
0	1	bulbasaur	1	0.7	6.9	64	grass	poison	45	49	49	45	65	65
1	2	ivysaur	1	1.0	13.0	142	grass	poison	60	62	63	60	80	80
2	3	venusaur	1	2.0	100.0	236	grass	poison	80	82	83	80	100	100
3	4	charmander	1	0.6	8.5	62	fire	NaN	39	52	43	65	60	50
4	5	charmeleon	1	1.1	19.0	142	fire	NaN	58	64	58	80	80	65
5	6	charizard	1	1.7	90.5	240	fire	flying	78	84	78	100	109	85
6	7	squirtle	1	0.5	9.0	63	water	NaN	44	48	65	43	50	64
7	8	wartortle	1	1.0	22.5	142	water	NaN	59	63	80	58	65	80
8	9	blastoise	1	1.6	85.5	239	water	NaN	79	83	100	78	85	105
9	10	caterpie	1	0.3	2.9	39	bug	NaN	45	30	35	45	20	20



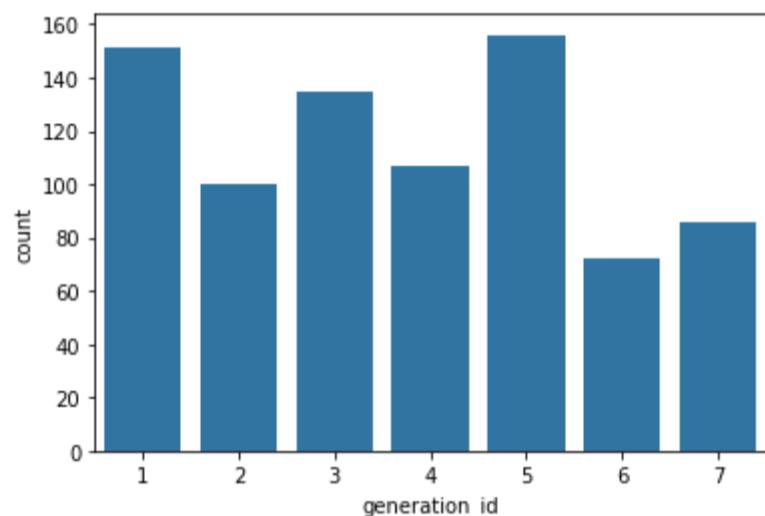
In the example above, all the bars have a different color. This might come in handy for building associations between these category labels and encodings in plots with more variables.

Otherwise, it's a good idea to simplify the plot and reduce unnecessary distractions by plotting all bars in the same color. You can

choose to have a uniform color across all bars, by using the `color` argument, as shown in the example below:

Example 2. Create a vertical bar chart using Seaborn, with a uniform single color

```
# The `color_palette()` returns the the current / default palette as a list of RGB tuples.  
# Each tuple consists of three digits specifying the red, green, and blue channel values to specify a  
color.  
# Choose the first tuple of RGB colors  
base_color = sb.color_palette()[0]  
  
# Use the `color` argument  
sb.countplot(data=pokemon, x='generation_id', color=base_color);
```



Bar Chart using the Matplotlib

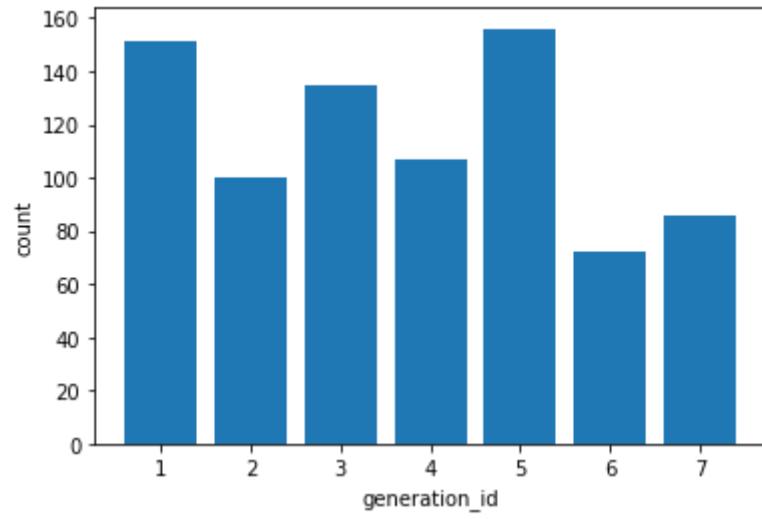
You can even create a similar bar chart using the Matplotlib, instead of Seaborn. We will use the `matplotlib.pyplot.bar()` function to plot the chart. The syntax is:

```
matplotlib.pyplot.bar(x, y, width=0.8, bottom=None, *, align='center', data=None)
```

Refer to the documentation for the details of optional arguments. In the example below, we will use `Series.value_counts()` to extract a Series from the given DataFrame object.

Example 3. Create a vertical bar chart using Matplotlib, with a uniform single color

```
# Return the Series having unique values  
x = pokemon['generation_id'].unique()  
  
# Return the Series having frequency count of each unique value  
y = pokemon['generation_id'].value_counts(sort=False)  
  
plt.bar(x, y)  
  
# Labeling the axes  
plt.xlabel('generation_id')  
plt.ylabel('count')  
  
# Dsiplay the plot  
plt.show()
```



There is a lot more you can do with both Seaborn and Matplotlib bar charts. *The remaining examples will experiment with seaborn's `countplot()` function.*

For nominal-type data, one common operation is to sort the data in terms of frequency. In the examples shown above, you can even order the bars as desirable. With our data in a pandas DataFrame, we can use various DataFrame methods to compute and extract an ordering, then set that ordering on the "order" parameter:

This can be done by using the `order` argument of the `countplot()` function.

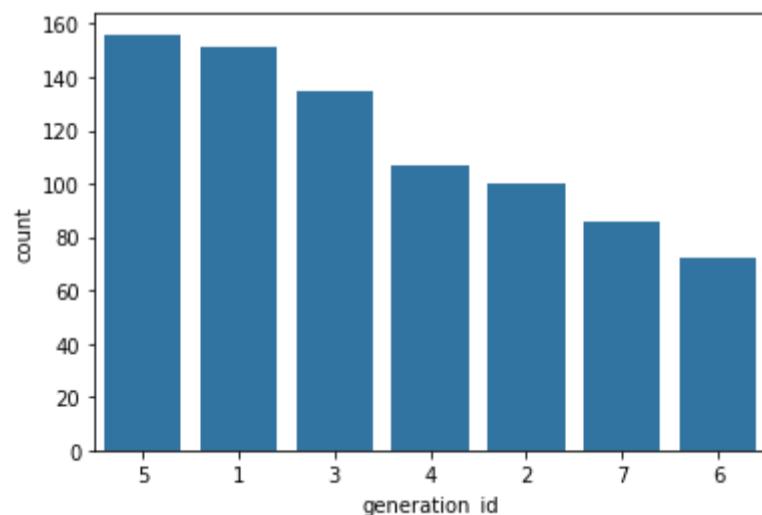
Example 4. Static and dynamic ordering of the bars in a bar chart using `seaborn.countplot()`

```
# Static-ordering the bars
sb.countplot(data=pokemon, x='generation_id', color=base_color, order=[5,1,3,4,2,7,6]);

# Dynamic-ordering the bars
# The order of the display of the bars can be computed with the following logic.
# Count the frequency of each unique value in the 'generation_id' column, and sort it in descending order
# Returns a Series
freq = pokemon['generation_id'].value_counts()

# Get the indexes of the Series
gen_order = freq.index

# Plot the bar chart in the decreasing order of the frequency of the `generation_id`
sb.countplot(data=pokemon, x='generation_id', color=base_color, order=gen_order);
```



While we could sort the levels by frequency like above, we usually care about whether the most frequent values are at high levels, low levels, etc. For ordinal-type data, we probably want to sort the bars in order of the variables. The best thing for us to do in this case is to convert the column into an ordered categorical data type.

Additional Variation - Refer to the [CategoricalDtype](#) to convert the column into an ordered categorical data type. By default, pandas reads in string data as object types, and will plot the bars in the order in which the unique values were seen. By converting the data into an ordered type, the order of categories becomes innate to the feature, and we won't need to specify an "order" parameter each time it's required in a plot.

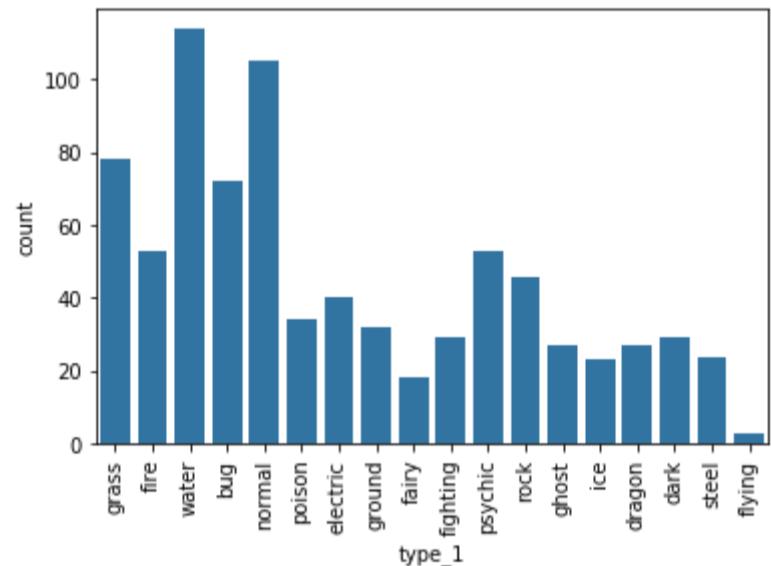
Should you find that you need to sort an ordered categorical type in a different order, you can always temporarily override the data type by setting the "order" parameter as above.

The category labels in the examples above are very small. In case, the category labels have large names, you can make use of the `matplotlib.pyplot.xticks(rotation=90)` function, which will rotate the category labels (not axes) counter-clockwise 90 degrees.

Example 5. Rotate the category labels (not axes)

```
# Plot the Pokemon type on a Vertical bar chart
sb.countplot(data=pokemon, x='type_1', color=base_color);

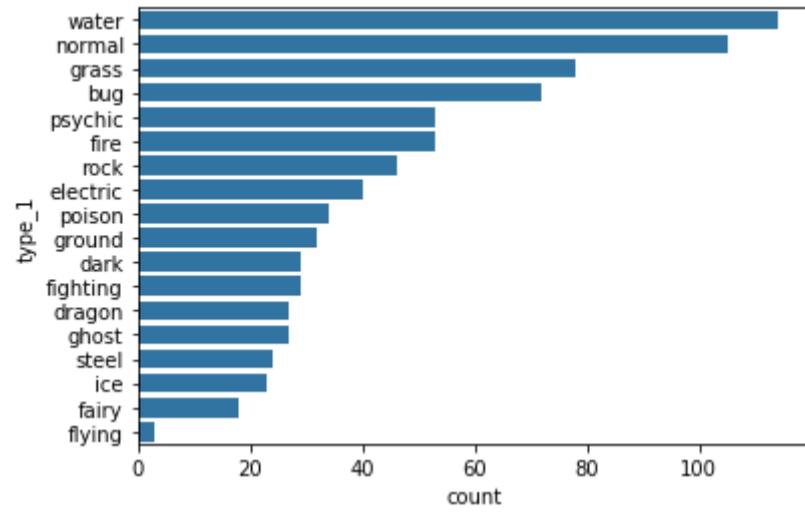
# Use xticks to rotate the category labels (not axes) counter-clockwise
plt.xticks(rotation=90)
```



Even after using the `matplotlib.pyplot.xticks(rotation=90)` function, if the category labels do not fit well, you can rotate the axes.

Example 6. Rotate the axes clockwise

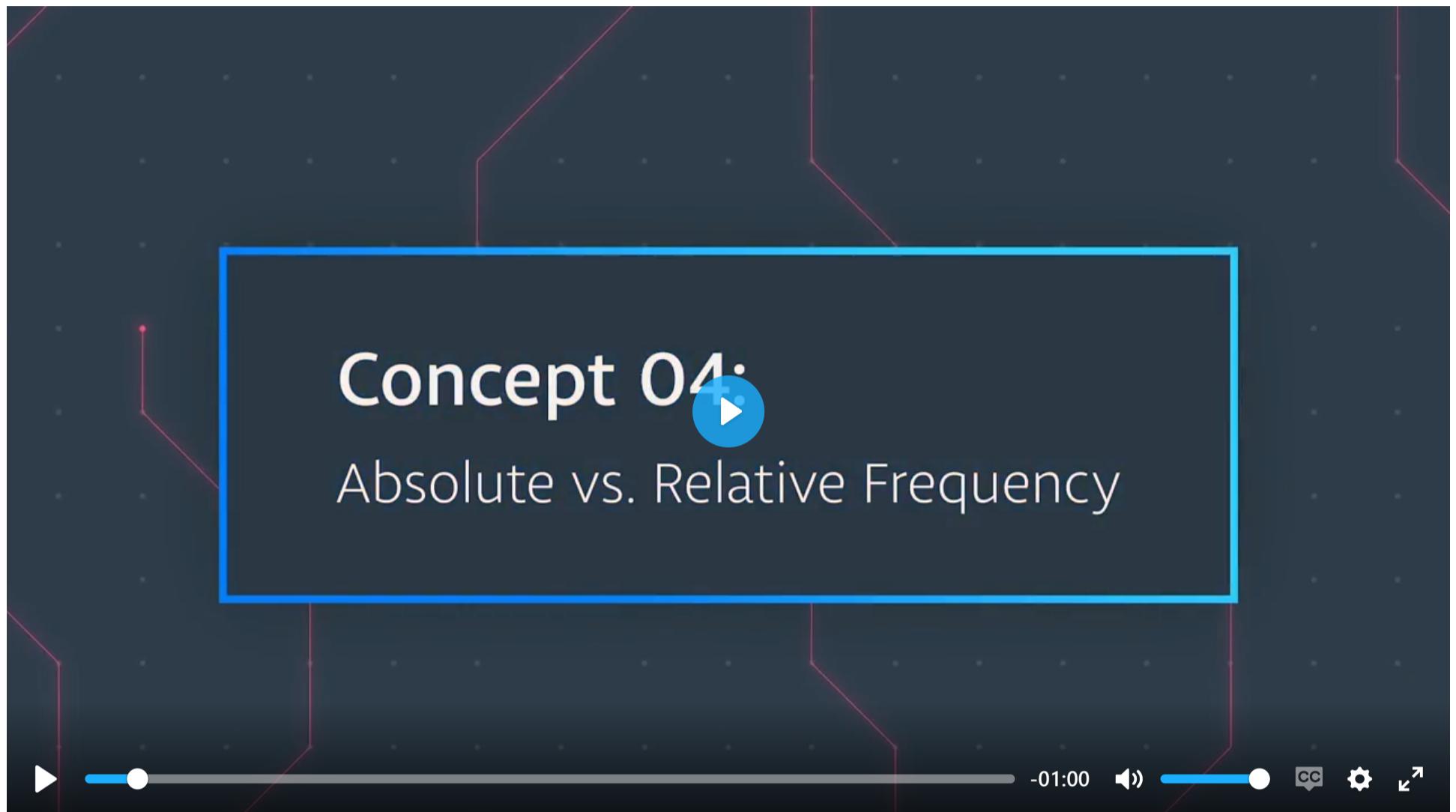
```
# Plot the Pokemon type on a Horizontal bar chart
type_order = pokemon['type_1'].value_counts().index
sb.countplot(data=pokemon, y='type_1', color=base_color, order=type_order);
```



[Next Concept](#)

≡ 04. Absolute vs. Relative Frequency

L3 041 Absolute V Relative Frequency V5



DataVis L3 04 V2



File Edit View Insert Cell Kernel Help

In [3]: `pkmn_types = pokemon.melt(id_vars = ['id', 'identifier'], value_vars = ['type_1', 'type_2'], var_name = 'type_level', value_name = 'type').dropna()`

Out[3]:

	id	identifier	type_level	type
802	803	poipole	type_1	poison
803	804	naganadel	type_1	poison
804	805	stakataka	type_1	rock
805	806	blacephalon	type_1	fire
806	807	zeraora	type_1	electric
807	1	bulbasaur	type_2	poison
808	2	ivysaur	type_2	poison
809	3	venusaur	type_2	poison
812	6	charizard	type_2	flying
818	12	butterfree	type_2	flying

In []: `type_counts = pkmn_types['type'].value_counts()`

Absolute vs. Relative Frequency

By default, seaborn's `countplot` function will summarize and plot the data in terms of **absolute frequency**, or pure counts. In certain cases, you might want to understand the distribution of data or want to compare levels in terms of the proportions of the whole. In this case, you will want to plot the data in terms of **relative frequency**, where the height indicates the proportion of data taking each level, rather than the absolute count.

One method of plotting the data in terms of relative frequency on a bar chart is to just relabel the count's axis in terms of proportions. The underlying data will be the same, it will simply be the scale of the axis ticks that will be changed.

Example 1. Demonstrate data wrangling, and plot a horizontal bar chart.

Example 1 - Step 1. Make the necessary import

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sb
%matplotlib inline

# Read the data from a CSV file
pokemon = pd.read_csv('pokemon.csv')
print(pokemon.shape)
pokemon.head(10)
```

	id	species	generation_id	height	weight	base_experience	type_1	type_2	hp	attack	defense	speed	special-attack	special-defense
0	1	bulbasaur	1	0.7	6.9	64	grass	poison	45	49	49	45	65	65
1	2	ivysaur	1	1.0	13.0	142	grass	poison	60	62	63	60	80	80
2	3	venusaur	1	2.0	100.0	236	grass	poison	80	82	83	80	100	100
3	4	charmander	1	0.6	8.5	62	fire	NaN	39	52	43	65	60	50
4	5	charmeleon	1	1.1	19.0	142	fire	NaN	58	64	58	80	80	65
5	6	charizard	1	1.7	90.5	240	fire	flying	78	84	78	100	109	85
6	7	squirtle	1	0.5	9.0	63	water	NaN	44	48	65	43	50	64
7	8	wartortle	1	1.0	22.5	142	water	NaN	59	63	80	58	65	80
8	9	blastoise	1	1.6	85.5	239	water	NaN	79	83	100	78	85	105
9	10	caterpie	1	0.3	2.9	39	bug	NaN	45	30	35	45	20	20

Last time we created the bar chart of pokemon by their `type_1`. Let's club the rows of both `type_1` and `type_2`, so that the resulting dataframe has **new** column, `type_level`.

This operation will double the number of rows in pokemon from 807 to 1614.

Data Wrangling Step

We will use the `pandas.DataFrame.melt()` method to unpivot a DataFrame from wide to long format, optionally leaving identifiers set. The syntax is:

```
DataFrame.melt(id_vars, value_vars, var_name, value_name, col_level, ignore_index)
```

It is essential to understand the parameters involved:

1. `id_vars` - It is a tuple representing the column(s) to use as identifier variables.
2. `value_vars` - It is tuple representing the column(s) to unpivot (remove, out of place).
3. `var_name` - It is a name of the **new** column.
4. `value_name` - It is a name to use for the 'value' of the columns that are unpivoted.

Refer [here](#) for more details on the parameters.

The function below will do the following in the pokemon dataframe *out of place*:

1. Select the 'id', and 'species' columns from pokemon.
2. Remove the 'type_1', 'type_2' columns from pokemon
3. Add a new column 'type_level' that can have a value either 'type_1' or 'type_2'
4. Add another column 'type' that will contain the actual value contained in the 'type_1', 'type_2' columns. For example, the first row in the pokemon dataframe having `id=1` and `species=bulbasaur` will now occur twice in the resulting dataframe after the `melt()` operation. The first occurrence will have `type=grass`, whereas, the second occurrence will have `type=poison`.

Example 1 - Step 2. Data wrangling to reshape the pokemon dataframe

```
pkmn_types = pokemon.melt(id_vars=['id', 'species'],
                           value_vars=['type_1', 'type_2'],
                           var_name='type_level',
                           value_name='type')

pkmn_types.head(10)
#pkmn_types.shape
```

	id	species	type_level	type
0	1	bulbasaur	type_1	grass
1	2	ivysaur	type_1	grass
2	3	venusaur	type_1	grass
3	4	charmander	type_1	fire
4	5	charmeleon	type_1	fire
5	6	charizard	type_1	fire
6	7	squirtle	type_1	water
7	8	wartortle	type_1	water
8	9	blastoise	type_1	water
9	10	caterpie	type_1	bug

Example 1 - Step 3. Find the frequency of unique values in the `type` column

```
# Count the frequency of unique values in the `type` column of pkmn_types dataframe.  
# By default, returns the decreasing order of the frequency.  
type_counts = pkmn_types['type'].value_counts()  
type_counts
```

```
water          131
normal         109
flying         98
grass          97
psychic        82
bug            77
poison          66
ground          64
fire            64
rock            60
fighting        54
electric         48
fairy           47
steel           47
dark            46
dragon          45
ghost           43
ice             34
```

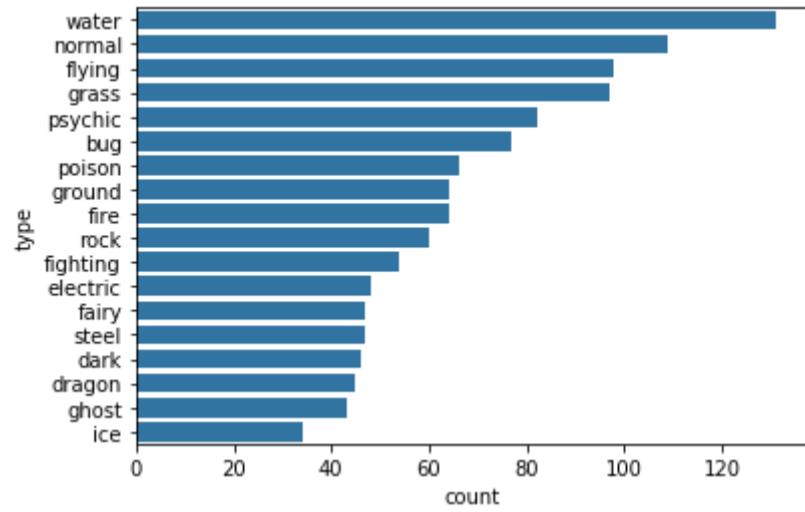
Name: type, dtype: int64

```
# Get the unique values of the `type` column, in the decreasing order of the frequency.
type_order = type_counts.index
type_order
```

```
Index(['water', 'normal', 'flying', 'grass', 'psychic', 'bug', 'poison',
'ground', 'fire', 'rock', 'fighting', 'electric', 'fairy', 'steel', 'dark', 'dragon', 'ghost', 'ice'],
dtype='object')
```

Example 1 - Step 4. Plot the horizontal bar charts

```
base_color = sb.color_palette()[0]
sb.countplot(data=pkmn_types, y='type', color=base_color, order=type_order);
```



Example 2. Plot a bar chart having the proportions, instead of the actual count, on one of the axes.

Example 2 - Step 1. Find the maximum proportion of bar

```
# Returns the sum of all not-null values in `type` column
n_pokemon = pkmn_types['type'].value_counts().sum()

# Return the highest frequency in the `type` column
max_type_count = type_counts[0]

# Return the maximum proportion, or in other words,
# compute the length of the longest bar in terms of the proportion
max_prop = max_type_count / n_pokemon
print(max_prop)
```

0.1623296158612144

Example 2 - Step 2. Create an array of evenly spaced proportioned values

```
# Use numpy.arange() function to produce a set of evenly spaced proportioned values
# between 0 and max_prop, with a step size 2\%
tick_props = np.arange(0, max_prop, 0.02)
tick_props
```

array([0., 0.02, 0.04, 0.06, 0.08, 0.1, 0.12, 0.14, 0.16])

We need x-tick labels that must be evenly spaced on the x-axis. For this purpose, we must have a list of labels ready with us, before using it with `plt.xticks()` function.

Example 2 - Step 3. Create a list of String values that can be used as tick labels.

```

# Use a list comprehension to create tick_names that we will apply to the tick labels.
# Pick each element `v` from the `tick_props`, and convert it into a formatted string.
# `{:0.2f}` denotes that before formatting, we 2 digits of precision and `f` is used to represent floating
point number.
# Refer [here](https://docs.python.org/2/library/string.html#format-string-syntax) for more details
tick_names = ['{:0.2f}'.format(v) for v in tick_props]
tick_names

```

```
['0.00', '0.02', '0.04', '0.06', '0.08', '0.10', '0.12', '0.14', '0.16']
```

The `xticks` and `yticks` functions aren't only about rotating the tick labels. You can also get and set their locations and labels as well. The first argument takes the tick locations: in this case, the tick proportions multiplied back to be on the scale of counts. The second argument takes the tick names: in this case, the tick proportions formatted as strings to two decimal places.

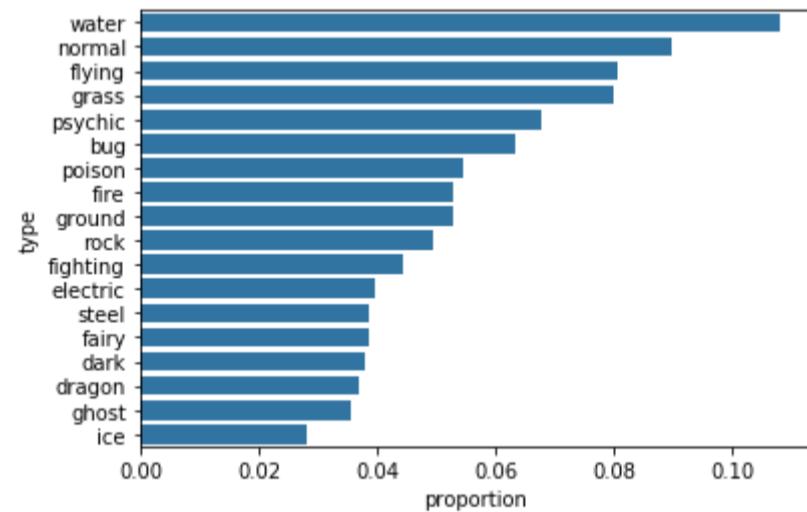
I've also added a `ylabel` call to make it clear that we're no longer working with straight counts.

Example 2 - Step 4. Plot the bar chart, with new x-tick labels

```

sb.countplot(data=pkmn_types, y='type', color=base_color, order=type_order);
# Change the tick locations and labels
plt.xticks(tick_props * n_pokemon, tick_names)
plt.xlabel('proportion');

```



Additional Variation

Rather than plotting the data on a relative frequency scale, you might use text annotations to label the frequencies on bars instead. This requires writing a loop over the tick locations and labels and adding one text element for each bar.

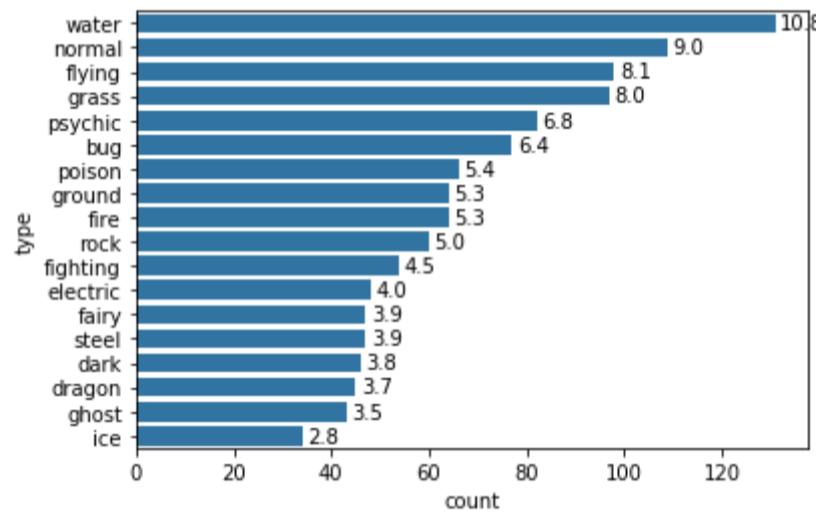
Example 3. Print the text (proportion) on the bars of a horizontal plot.

```

# Considering the same chart from the Example 1 above, print the text (proportion) on the bars
base_color = sb.color_palette()[0]
sb.countplot(data=pknn_types, y='type', color=base_color, order=type_order);

# Logic to print the proportion text on the bars
for i in range (type_counts.shape[0]):
    # Remember, type_counts contains the frequency of unique values in the `type` column in decreasing
    # order.
    count = type_counts[i]
    # Convert count into a percentage, and then into string
    pct_string = '{:0.1f}'.format(100*count/n_pokemon)
    # Print the string value on the bar.
    # Read more about the arguments of text() function [here]
    # (https://matplotlib.org/3.1.1/api/_as_gen/matplotlib.pyplot.text.html)
    plt.text(count+1, i, pct_string, va='center')

```



Example 4. Print the text (proportion) below the bars of a Vertical plot.

```

# Considering the same chart from the Example 1 above, print the text (proportion) BELOW the bars
base_color = sb.color_palette()[0]
sb.countplot(data=pknn_types, x='type', color=base_color, order=type_order);

# Recalculating the type_counts just to have clarity.
type_counts = pknn_types['type'].value_counts()

# get the current tick locations and labels
locs, labels = plt.xticks(rotation=90)

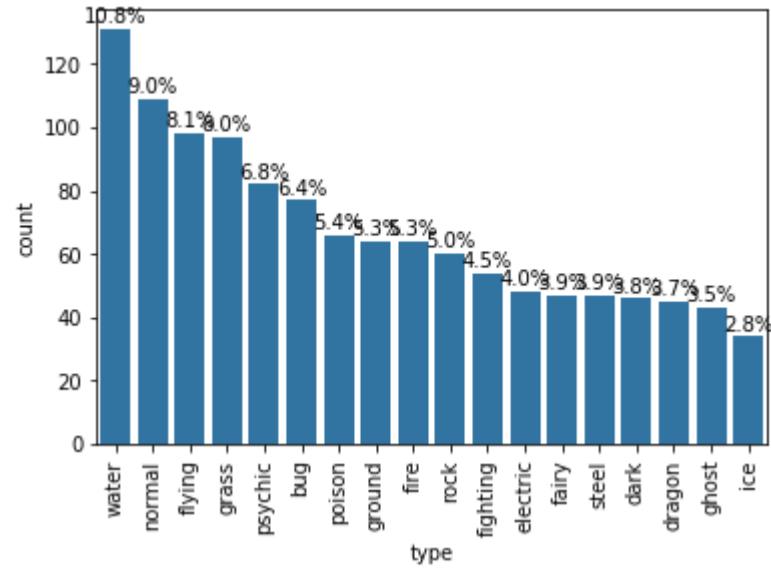
# loop through each pair of locations and labels
for loc, label in zip(locs, labels):

    # get the text property for the label to get the correct count
    count = type_counts[label.get_text()]
    pct_string = '{:0.1f}%'.format(100*count/n_pokemon)

    # print the annotation just below the top of the bar
    plt.text(loc, count+2, pct_string, ha = 'center', color = 'black')

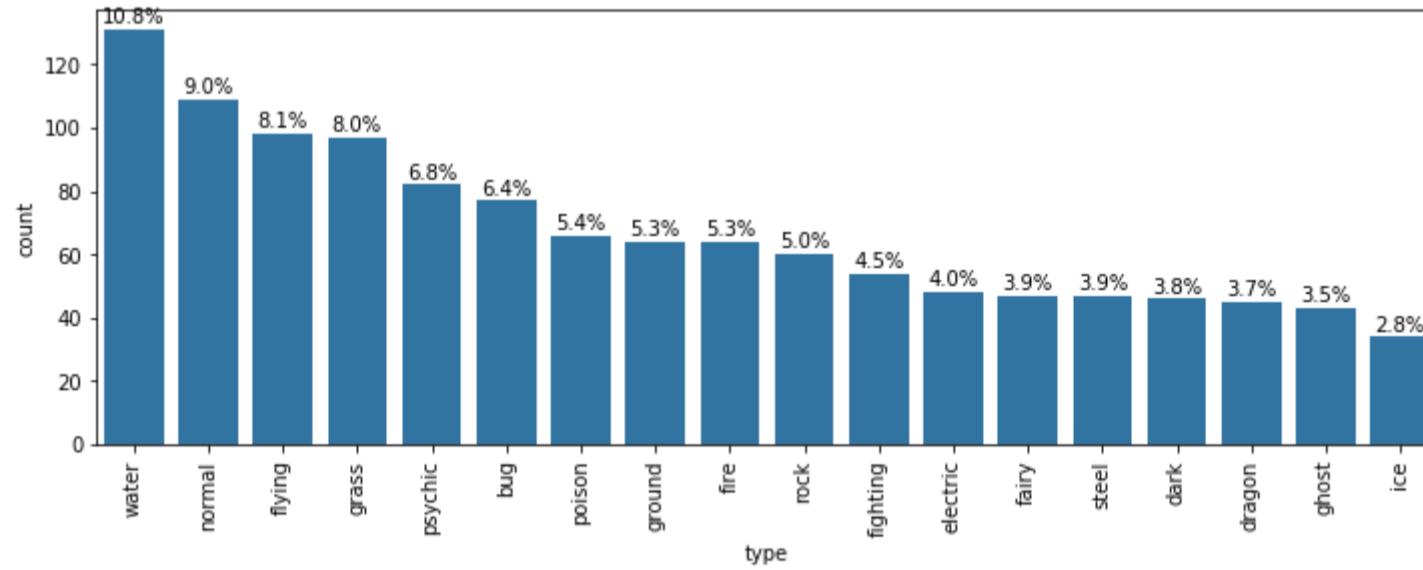
```

I use the `.get_text()` method to obtain the category name, so I can get the count of each category level. At the end, I use the `text` function to print each percentage, with the x-position, y-position, and string as the three main parameters to the function.



Tip - Is the text on the bars not readable clearly? Consider changing the size of the plot by using the following:

```
from matplotlib import rcParams  
# Specify the figure size in inches, for both X, and Y axes  
rcParams['figure.figsize'] = 12,4
```



[Next Concept](#)

≡ 05. Counting Missing Data

Counting Missing Data

If you have a large dataframe, and it contains a few missing values (`None` or a `numpy.Nan`), then you can find the count of such missing value across the given label.

For this purpose, you can use either of the following two analogous functions :

1. [pandas.DataFrame.isna\(\)](#)
2. [pandas.DataFrame.isnull\(\)](#)

The functions above are alias of each other and detect missing values by returning the same sized object as that of the calling dataframe, made up of boolean True/False.

Step 1. Load the dataset

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sb
%matplotlib inline

# Read the data from a CSV file
# Original source of data: https://www.kaggle.com/manjeetsingh/retaildataset available under CC0 1.0
Universal (CC0 1.0) Public Domain Dedication License
sales_data = pd.read_csv('sales_data.csv')
sales_data.head(10)
```

	Store	Date	Temperature	Fuel_Price	MarkDown1	MarkDown2	MarkDown3	MarkDown4	MarkDown5	CPI	Unemployment	IsHoliday
0	1	05/02/2010	42.31	2.572	NaN	NaN	NaN	NaN	NaN	211.096358	8.106	False
1	1	12/02/2010	38.51	2.548	NaN	NaN	NaN	NaN	NaN	211.242170	8.106	True
2	1	19/02/2010	39.93	2.514	NaN	NaN	NaN	NaN	NaN	211.289143	8.106	False
3	1	26/02/2010	46.63	2.561	NaN	NaN	NaN	NaN	NaN	211.319643	8.106	False
4	1	05/03/2010	46.50	2.625	NaN	NaN	NaN	NaN	NaN	211.350143	8.106	False
5	1	12/03/2010	57.79	2.667	NaN	NaN	NaN	NaN	NaN	211.380643	8.106	False
6	1	19/03/2010	54.58	2.720	NaN	NaN	NaN	NaN	NaN	211.215635	8.106	False
7	1	26/03/2010	51.45	2.732	NaN	NaN	NaN	NaN	NaN	211.018042	8.106	False
8	1	02/04/2010	62.27	2.719	NaN	NaN	NaN	NaN	NaN	210.820450	7.808	False
9	1	09/04/2010	65.86	2.770	NaN	NaN	NaN	NaN	NaN	210.622857	7.808	False

```
sales_data.shape
```

```
(8190, 12)
```

```
# Use either of the functions below
# sales_data.isna()
sales_data.isnull()
```

	Store	Date	Temperature	Fuel_Price	MarkDown1	MarkDown2	MarkDown3	MarkDown4	MarkDown5	CPI	Unemployment	IsHoliday
0	False	False	False	False	True	True	True	True	True	False	False	False
1	False	False	False	False	True	True	True	True	True	False	False	False
2	False	False	False	False	True	True	True	True	True	False	False	False
3	False	False	False	False	True	True	True	True	True	False	False	False
4	False	False	False	False	True	True	True	True	True	False	False	False
...
8185	False	False	False	False	False	False	False	False	False	True	True	True
8186	False	False	False	False	False	False	False	False	False	True	True	True
8187	False	False	False	False	False	False	False	False	False	True	True	True
8188	False	False	False	False	False	False	False	False	False	True	True	True
8189	False	False	False	False	False	False	False	False	False	True	True	True

We can use pandas functions to create a table with the number of missing values in each column. Once, you have the label-wise count of missing values, you try plotting the tabular data in the form of a bar chart.

```
sales_data.isna().sum()
```

```
Store          0
Date           0
Temperature    0
Fuel_Price     0
MarkDown1     4158
MarkDown2     5269
MarkDown3     4577
MarkDown4     4726
MarkDown5     4140
CPI            585
Unemployment  585
IsHoliday      0
dtype: int64
```

What if we want to visualize these missing value counts?

One interesting way we can apply bar charts is through the visualization of missing data. We could treat the variable names as levels of a categorical variable, and create a resulting bar plot. However, since the data is not in its tidy, unsummarized form, we need to make use of a different plotting function. Seaborn's `barplot` function is built to depict a summary of one quantitative variable against levels of a second, qualitative variable, but can be used here.

Step 2 - Prepare a NaN tabular data

```
# Let's drop the column that do not have any NaN/None values
na_counts = sales_data.drop(['Date', 'Temperature', 'Fuel_Price'], axis=1).isna().sum()
print(na_counts)
```

```
Store          0
MarkDown1     4158
MarkDown2     5269
MarkDown3     4577
MarkDown4     4726
MarkDown5     4140
CPI           585
Unemployment 585
IsHoliday      0
dtype: int64
```

Use `seaborn.barplot()`

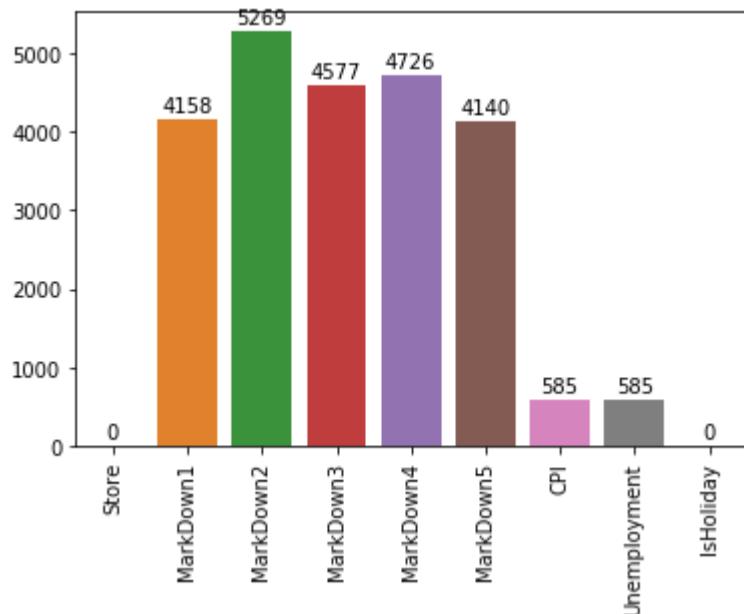
Step 3 - Plot the bar chart from the NaN tabular data, and also print values on each bar

```
# The first argument to the function below contains the x-values (column names), the second argument the
y-values (our counts).
# Refer to the syntax and more example here - https://seaborn.pydata.org/generated/seaborn.barplot.html
sb.barplot(na_counts.index.values, na_counts)

# get the current tick locations and labels
plt.xticks(rotation=90)

# Logic to print value on each bar
for i in range (na_counts.shape[0]):
    count = na_counts[i]

    # Refer here for details of the text() -
https://matplotlib.org/3.1.1/api/\_as\_gen/matplotlib.pyplot.text.html
    plt.text(i, count+300, count, ha = 'center', va='top')
```

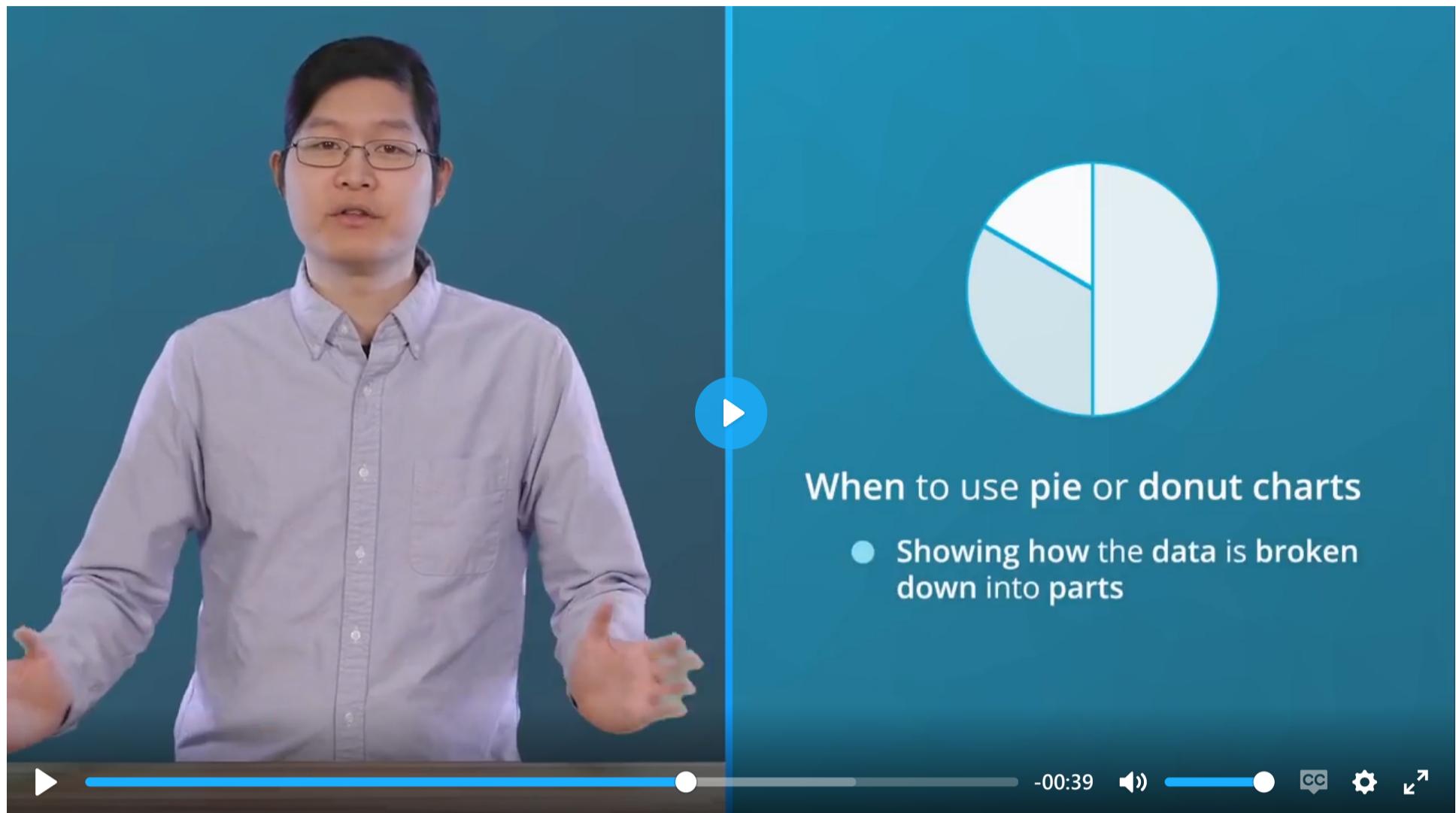


Note - The `seaborn.barplot()` is a useful function to keep in mind if your data is summarized and you still want to build a bar chart. If your data is not yet summarized, however, just use the `countplot` function so that you don't need to do extra summarization work. In addition, you'll see what `barplot`'s main purpose is in the next lesson when we discuss adaptations of univariate plots for plotting bivariate data.

[Next Concept](#)

≡ 07. Pie Charts

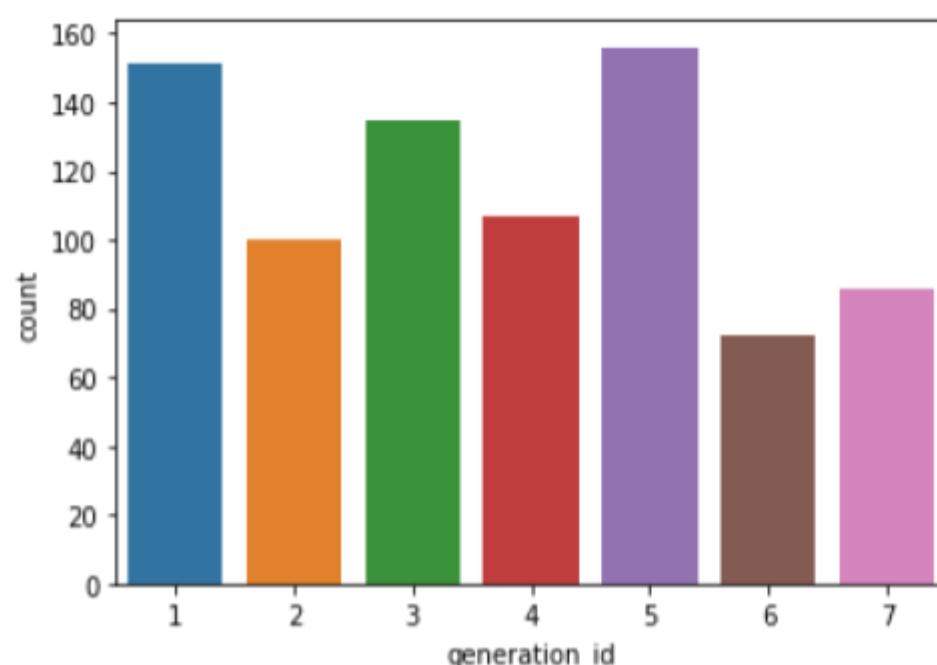
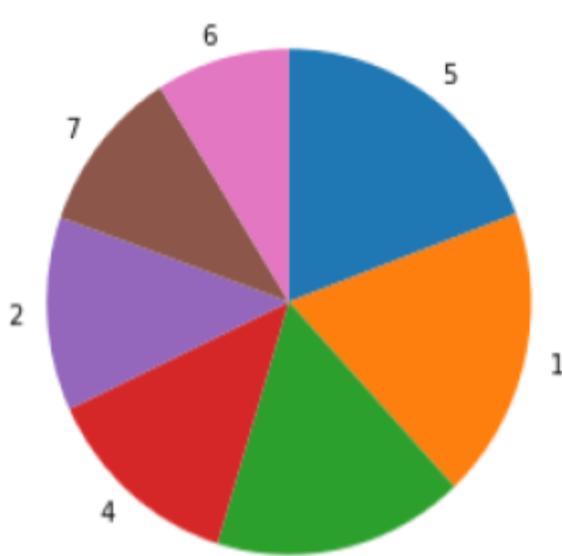
L3 071 Pie Charts V3



A video player interface showing a person in a light blue shirt speaking. To the right of the video, there is a pie chart divided into two equal halves. Below the pie chart, the text "When to use pie or donut charts" is displayed, followed by a bullet point: "Showing how the data is broken down into parts". At the bottom of the video player, there is a progress bar, a play button, and various control icons.

Pie Charts

A **pie chart** is a common univariate plot type that is used to depict relative frequencies for levels of a categorical variable. Frequencies in a pie chart are depicted as wedges drawn on a circle: the larger the angle or area, the more common the categorical value taken. Use a Pie chart only when the number of categories is less, and you'd like to see the proportion of each category on a chart.



Pie chart (left) and bar chart (right) displaying the same categorical counts, 'generation_id' in the pokemon dataset.

Unfortunately, pie charts are a fairly limited plot type in the range of scenarios where they can be used, and it is easy for chart makers to try and spice up pie charts in a way that makes them more difficult to read.

Guidelines to Use a Pie Chart

If you want to use a pie chart, try to follow certain guidelines:

- Make sure that your interest is in *relative* frequencies. Areas should represent parts of a whole, rather than measurements on a second variable (unless that second variable can logically be summed up into some whole).
- Limit the number of slices plotted. A pie chart works best with two or three slices, though it's also possible to plot with four or five slices as long as the wedge sizes can be distinguished. If you have a lot of categories, or categories that have small proportional representation, consider grouping them together so that fewer wedges are plotted, or use an 'Other' category to handle them.
- Plot the data systematically. One typical method of plotting a pie chart is to start from the top of the circle, then plot each categorical level clockwise from most frequent to least frequent. If you have three categories and are interested in the comparison of two of them, a common plotting method is to place the two categories of interest on either side of the 12 o'clock direction, with the third category filling in the remaining space at the bottom.

If these guidelines cannot be met, then you should probably make use of a bar chart instead. A bar chart is a safer choice in general. The bar heights are more precisely interpreted than areas or angles, and a bar chart can be displayed more compactly than a pie chart. There's also more flexibility with a bar chart for plotting variables with a lot of levels, like plotting the bars horizontally.

Plot a Pie Chart

matplotlib.pyplot.pie()

You can create a pie chart with matplotlib's `matplotlib.pyplot.pie()` function. A basic syntax is:

```
matplotlib.pyplot.pie(x_data, labels, colors, startangle, counterclock, wedgeprops)
```

This function requires that the data be in a summarized form: the primary argument to the function will be the wedge sizes. Refer to the [function syntax](#) for details about all other arguments.

matplotlib.pyplot.axis()

We also need to know about the `matplotlib.pyplot.axis()` function to set some axis properties. It optionally accepts the axis limits in the form of `xmin, xmax, ymin, ymax` floats, and returns the updated values.

```
matplotlib.pyplot.axis(*args, emit=True, **kwargs)
```

In the function above, the `*args` represents any number of arguments that you can pass to the function, whereas `**kwargs` stands for keyword arguments, generally passed in the form of a dictionary.

Refer to the [function syntax](#) for in-depth details on the all possible values of the arguments.

Example 1. Plot a simple Pie chart

```
# Use the same pokemon dataset
sorted_counts = pokemon['generation_id'].value_counts()

plt.pie(sorted_counts, labels = sorted_counts.index, startangle = 90, counterclock = False);

# We have the used option `Square`.
# Though, you can use either one specified here -
https://matplotlib.org/api/_as_gen/matplotlib.pyplot.axis.html?highlight=pyplot%20axis#matplotlib-pyplot-
axis
plt.axis('square')
```

To follow the guidelines in the bullet points above, I include the "startangle = 90" and "counterclock = False" arguments to start the first slice at vertically upwards, and will plot the sorted counts in a clockwise fashion. The `axis` function call and 'square' argument makes it so that the scaling of the plot is equal on both the x- and y-axes. Without this call, the pie could end up looking oval-shaped, rather than a circle.

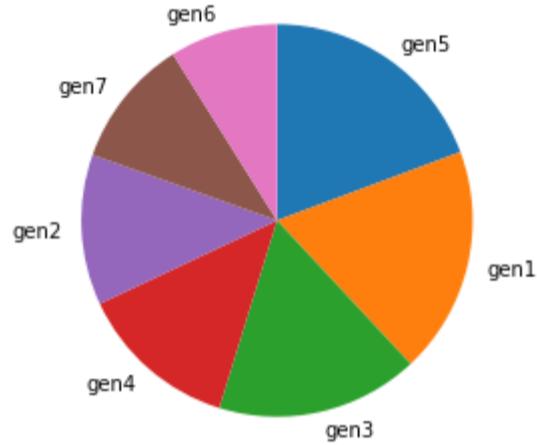
TO DO

Did you notice the various arguments in the `plt.pie()` function? Particularly, the `labels = sorted_counts.index` argument represents a list of strings serving as labels for each wedge. In the example above, the labels have used the following list:

```
sorted_counts.index
```

```
([5, 1, 3, 4, 2, 7, 6], dtype='int64')
```

Can you try using another list of strings, `['gen5', 'gen1', 'gen3', 'gen4', 'gen2', 'gen7', 'gen6']`, to display labels on each wedge? Notice that the labels are arranged in the decreasing order of the frequency. The expected output is shown below.



A bar chart with an updated label on each wedge.

Donut Plot

A sister plot to the pie chart is the **donut plot**. It's just like a pie chart, except that there's a hole in the center of the plot. Perceptually, there's not much difference between a donut plot and a pie chart, and donut plots should be used with the same guidelines as a pie chart. Aesthetics might be one of the reasons why you would choose one or the other. For instance, you might see statistics reported in the hole of a donut plot to better make use of available space.

To create a donut plot, you can add a `wedgeprops` argument to the `pie` function call. By default, the radius of the pie (circle) is 1; setting the wedges' width property to less than 1 removes coloring from the center of the circle.

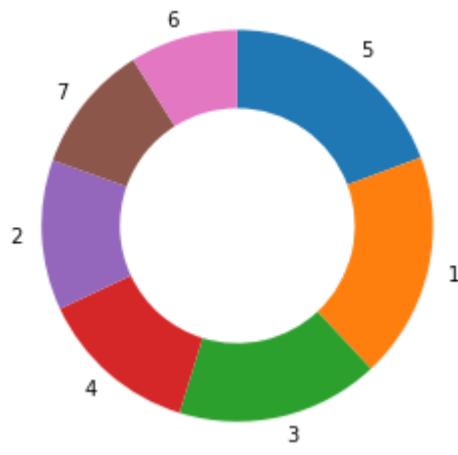
Example 2. Plot a simple Donut plot

```
sorted_counts = pokemon['generation_id'].value_counts()

plt.pie(sorted_counts, labels = sorted_counts.index, startangle = 90,
        counterclock = False, wedgeprops = {'width' : 0.4});
plt.axis('square')
```

Recommended Read

Refer to the documentation: [Wedge patches](#)



Donut plot representing the 'generation_id' in the pokemon dataset

Further Reading

- Eager Eyes: [Understanding Pie Charts](#)
- Eager Eyes: [An Illustrated Tour of the Pie Chart Study Results](#) - how accurately do people perceive different formulations of the pie chart?
- Datawrapper: [What to Consider when Creating a Pie Chart](#)

[Next Concept](#)

☰ 08. Histograms

L3 081 Histograms V2



DataVis L3 08 V2

(807, 14)

Out[2]:

	id	identifier	generation_id	height	weight	base_experience	type_1	type_2	hp	attack	defense	speed	special_attack
0	1	bulbasaur	1	0.7	6.9	64	grass	poison	45	49	49	45	65
1	2	ivysaur	1	1.0	13.0	142	grass	poison	60	62	63	60	80
2	3	venusaur	1	2.0	100.0	236	grass	poison	80	82	83	80	100
3	4	charmander	1	0.6	8.5	62	fire	NaN	39	52	43	65	60
4	5	charmeleon	1	1.1	19.0	142	fire	NaN	58	64	58	80	80
5	6	charizard	1	1.7	90.5	240	fire	flying	78	84	78	100	109
6	7	squirtle	1	0.5	9.0	63	water	NaN	44	48	65	43	50
7	8	wartortle	1	1.0	22.5	142	water	NaN	59	63	80	58	65
8	9	blastoise	1	1.6	85.5	239	water	NaN	79	83	100	78	85
9	10	caterpie	1	0.3	2.9	39	bug	NaN	45	30	35	45	20

In []:

In []:

03:17 CC

Histograms

A **histogram** is used to plot the distribution of a numeric variable. It's the quantitative version of the bar chart. However, rather than plot one bar for each unique numeric value, values are grouped into continuous bins, and one bar for each bin is plotted to depict the number. You can use either Matplotlib or Seaborn to plot the histograms. There is a mild variation in the specifics, such as plotting gaussian-estimation line along with bars in Seabron's distplot(), and the arguments that you can use in either case.

Matplotlib.pyplot.hist()

You can use the default settings for matplotlib's [hist\(\)](#) function to plot a histogram with 10 bins:

Example 1. Plot a default histogram

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sb
%matplotlib inline

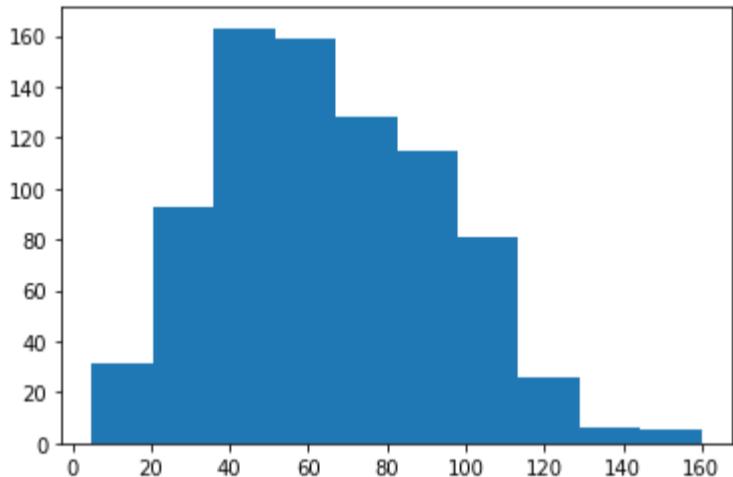
pokemon = pd.read_csv('pokemon.csv')
print(pokemon.shape)
pokemon.head(10)
```

	id	species	generation_id	height	weight	base_experience	type_1	type_2	hp	attack	defense	speed	special-attack	special-defense
0	1	bulbasaur	1	0.7	6.9	64	grass	poison	45	49	49	45	65	65
1	2	ivysaur	1	1.0	13.0	142	grass	poison	60	62	63	60	80	80
2	3	venusaur	1	2.0	100.0	236	grass	poison	80	82	83	80	100	100
3	4	charmander	1	0.6	8.5	62	fire	NaN	39	52	43	65	60	50
4	5	charmeleon	1	1.1	19.0	142	fire	NaN	58	64	58	80	80	65
5	6	charizard	1	1.7	90.5	240	fire	flying	78	84	78	100	109	85
6	7	squirtle	1	0.5	9.0	63	water	NaN	44	48	65	43	50	64
7	8	wartortle	1	1.0	22.5	142	water	NaN	59	63	80	58	65	80
8	9	blastoise	1	1.6	85.5	239	water	NaN	79	83	100	78	85	105
9	10	caterpie	1	0.3	2.9	39	bug	NaN	45	30	35	45	20	20

We will plot the `speed` column on histogram

```
# We have intentionally not put a semicolon at the end of the statement below to see the bar-width
plt.hist(data = pokemon, x = 'speed')
```

```
(array([ 31.,  93., 163., 159., 128., 115., 81., 26., 6., 5.]), array([ 5. , 20.5, 36. , 51.5, 67. , 82.5, 98. , 113.5, 129.
, 144.5, 160. ]))
```



You can see a non-uniform distribution of data points in different bins.

Overall, a generally bimodal distribution is observed (one with two peaks or humps). The direct adjacency of the bars in the histogram, in contrast to the separated bars in a bar chart, emphasizes the fact that the data takes on a continuous range of values. When a data value is on a bin edge, it is counted in the bin to its right. The exception is the rightmost bin edge, which places data values equal to the uppermost limit into the right-most bin (to the upper limit's left).

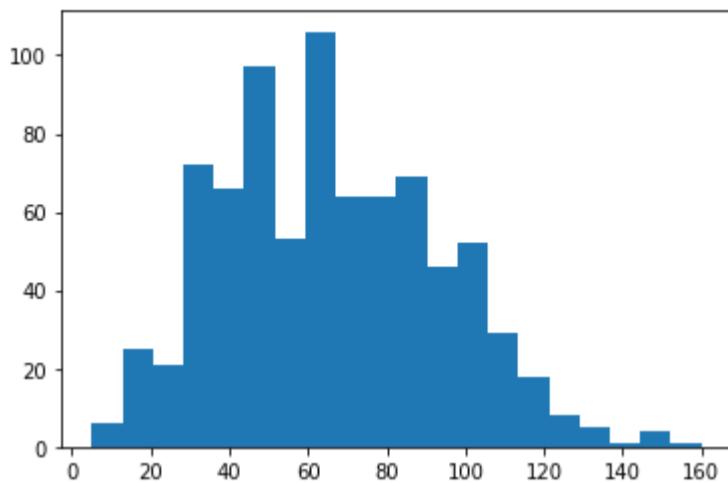
By default, the `hist` function divides the data into 10 bins, based on the range of values taken. In almost every case, we will want to change these settings. Usually, having only ten bins is too few to really understand the distribution of the data. And the default tick marks are often not on nice, 'round' values that make the ranges taken by each bin easy to interpret.

Wouldn't it be better if I said "between 0 and 2.5" instead of "between *about* 0 and 2.5", and "from 2.5 to 5" instead of "from *about* 2.5 to 5" above?

You can use descriptive statistics (e.g. via `dataframe['column'].describe()`) to gauge what minimum and maximum bin limits might be appropriate for the plot. These bin edges can be set using numpy's `arange` function:

Example 2. Histogram with fixed number of bins

```
plt.hist(data = pokemon, x = 'speed', bins = 20)
```

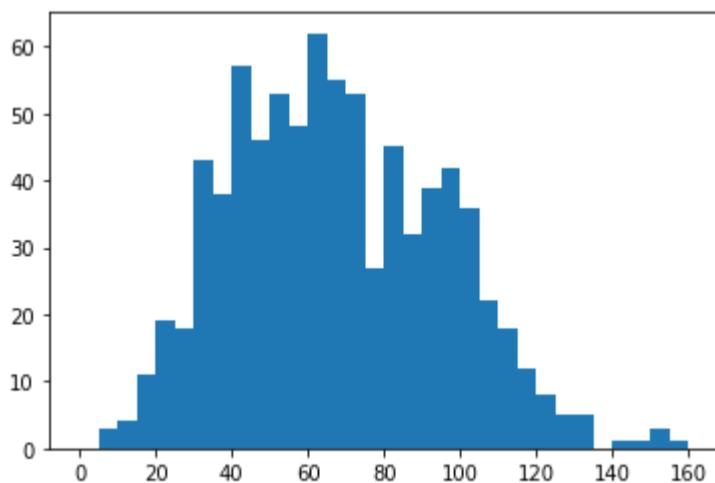


A histogram having 20 bins

Example 3. Histogram with dynamic number of bins

```
# Create bins with step-size 5
bins = np.arange(0, pokemon['speed'].max()+5, 5)
plt.hist(data = pokemon, x = 'speed', bins = bins)
```

The first argument to `arange` is the leftmost bin edge, the second argument the upper limit, and the third argument the bin width. Note that even though I've specified the "max" value in the second argument, I've added a "+5" (the bin width). That is because `arange` will only return values that are strictly less than the upper limit. Adding in "+5" is a safety measure to ensure that the rightmost bin edge is at least the maximum data value, so that all of the data points are plotted. The leftmost bin is set as a hardcoded value to get a nice, interpretable value, though you could use functions like numpy's `around` if you wanted to approach that end programmatically.



A histogram with a dynamic number of bins, each with a step-size 5.

Alternative Approach - Seaborn's `distplot()`

This function can also plot histograms, as similar to the `pyploy.hist()` function, and is integrated with other univariate plotting functions. This is in contrast to our ability to specify a data source and column as separate arguments, like we've seen with and `countplot` and `hist`.

The basic syntax is:

```
seaborn.distplot(Series, bins, kde, hist_kws)
```

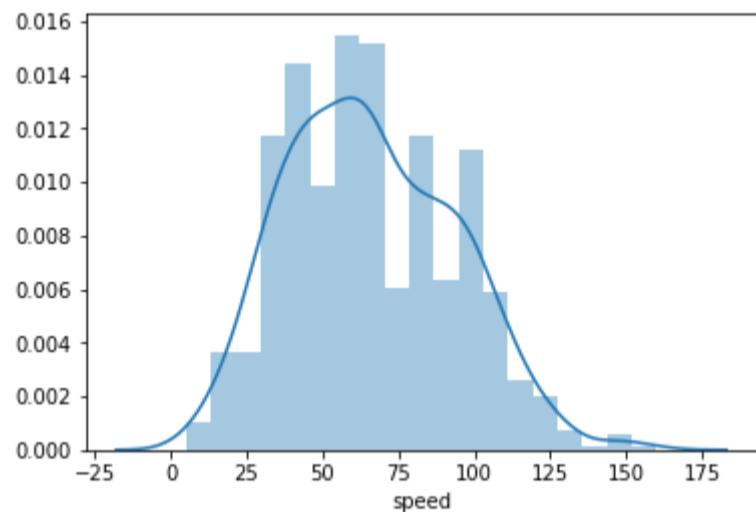
Let's see the sample usage of the arguments mentioned above. However, there are many other arguments that you can explore in the [syntax definition](#).

Note: From the Seaborn v0.11.0 onwards, this function is deprecated and will be removed in a future version. You can use either of the following two functions: [displot\(\)](#) or [histplot\(\)](#) to plot histograms using Seaborn.

Example 4. Plot the similar histogram with Seaborn's `distplot()`

```
sb.distplot(pokemon['speed'])

# Set the argument `kde=False` to remove the estimate-line representing the Gaussian kernel density
# estimate.
sb.distplot(pokemon['speed'], kde=False)
```



A default histogram plotted using Seaborn's `distplot()` function.

The `distplot` function has built-in rules for specifying histogram bins, and by default plots a curve depicting the kernel density estimate (KDE) on top of the data. The vertical axis is based on the KDE, rather than the histogram: you shouldn't expect the total heights of the bars to equal 1, but the area under the curve should equal 1. If you want to learn more about KDEs, check out the extra page at the end of the lesson.

Despite the fact that the default bin-selection formula used by `distplot` might be better than the choice of ten bins that `.hist` uses, you'll still want to do some tweaking to align the bins to 'round' values. You can use other parameter settings to plot just the histogram and specify the bins like before:

```
bin_edges = np.arange(0, df['num_var'].max()+1, 1)
sb.distplot(df['num_var'], bins = bin_edges, kde = False,
            hist_kws = {'alpha' : 1})
```

The alpha (transparency) setting must be associated as a dictionary to "hist_kws" since there are other underlying plotting functions, like the KDE, that have their own optional keyword parameters.

Plot two histograms side-by-side

When creating histograms, it's useful to play around with different bin widths to see what represents the data best. Too many bins, and you may see too much noise that interferes with the identification of the underlying signal. Too few bins, and you may not be able to see the true signal in the first place.

Let's see a new example demonstrating a few new functions, `pyplot.subplot()` and `pyplot.figure()`. We will learn more in the upcoming concepts.

Example 5. Plot two histograms side-by-side

```
'''python
```

Resize the chart, and have two plots side-by-side

Set a larger figure size for subplots

```
plt.figure(figsize = [20, 5])
```

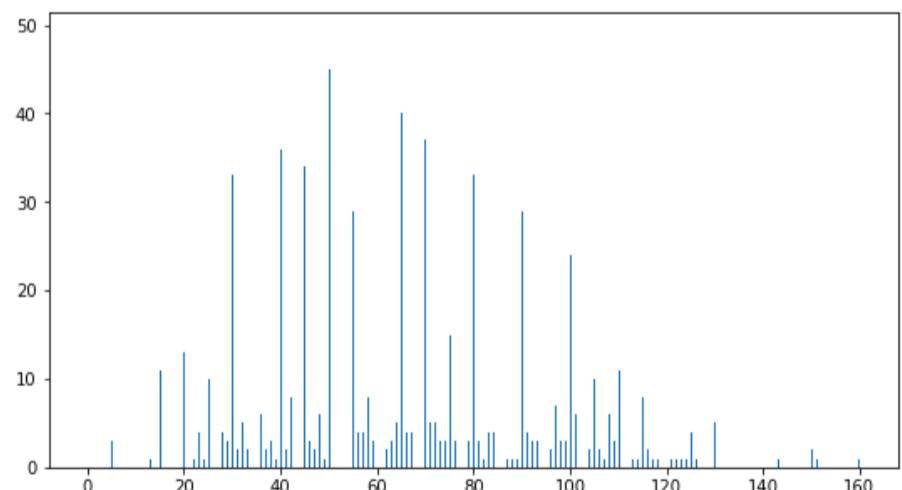
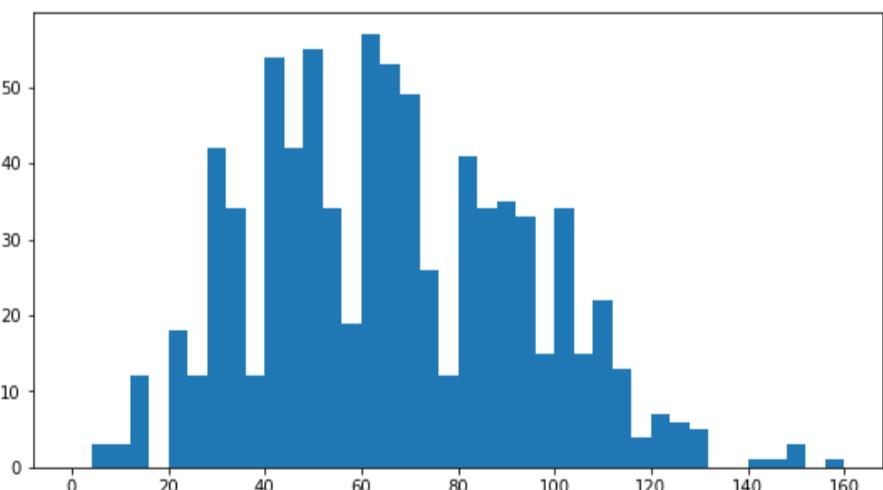
histogram on left, example of too-large bin size

1 row, 2 cols, subplot 1

```
plt.subplot(1, 2, 1)  
bins = np.arange(0, pokemon['speed'].max() + 4, 4)  
plt.hist(data = pokemon, x = 'speed', bins = bins);
```

histogram on right, example of too-small bin size

```
plt.subplot(1, 2, 2) # 1 row, 2 cols, subplot 2  
bins = np.arange(0, pokemon['speed'].max() + 1/4, 1/4)  
plt.hist(data = pokemon, x = 'speed', bins = bins);  
  
``` This example puts two plots side by side through use of the [ subplot ]  
[https://matplotlib.org/api/_as_gen/matplotlib.pyplot.subplot.html] function, whose arguments specify the number of rows,
columns, and index of the active subplot (in that order). The [figure()]
[https://matplotlib.org/api/_as_gen/matplotlib.pyplot.figure.html] function is called with the "figsize" parameter so that we can have a
larger figure to support having multiple subplots. (More details on figures and subplots are coming up next in the lesson.)
```



## Summary of Histograms

In summary, if your exploration is only interested in the histogram-depiction of the data, and not the additional functionality offered by `distplot`, then you might be better off with just using Matplotlib's `hist` function for simplicity. On the other hand, if you want a quick start on choosing a representative bin size for histogram plotting, you might take a quick look at the basic `distplot` first before getting into the customization.

[Next Concept](#)

# ≡ 10. Figures, Axes, and Subplots

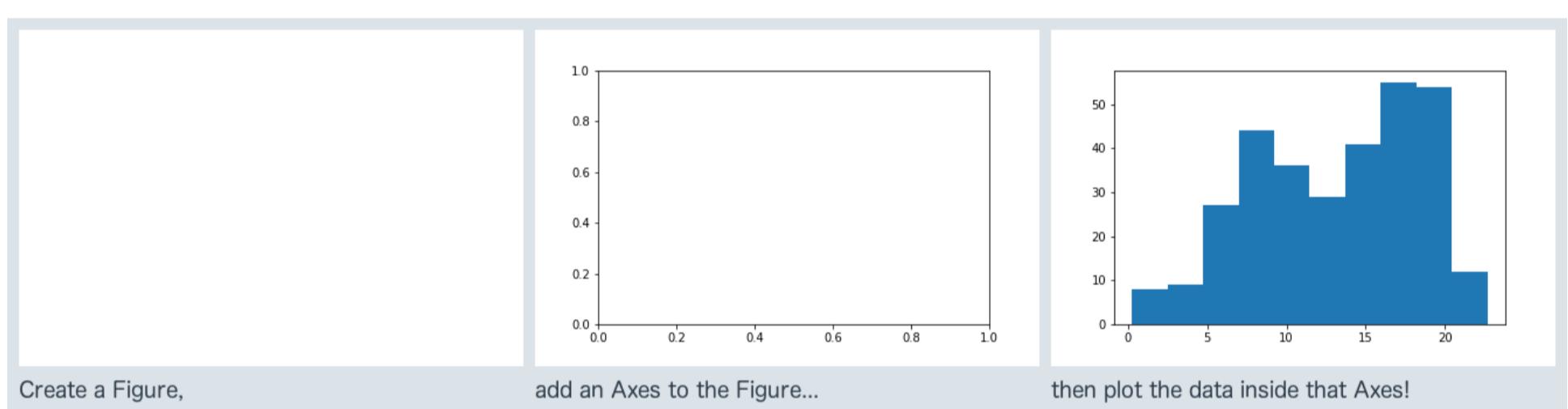
## Figures, Axes, and Subplots

At this point, you've seen and had some practice with some basic plotting functions using matplotlib and seaborn. The previous page introduced something a little bit new: creating two side-by-side plots through the use of matplotlib's `subplot()` function. If you have any questions about how that or the `figure()` function worked, then read on. This page will discuss the basic structure of visualizations using matplotlib and how subplots work in that structure.

The base of visualization in matplotlib is a [Figure](#) object. Contained within each Figure will be one or more [Axes](#) objects, each Axes object containing a number of other elements that represent each plot. In the earliest examples, these objects have been created implicitly. Let's say that the following expression is run inside a Jupyter notebook to create a histogram:

```
plt.hist(data=pokemon, x='speed');
```

Since we don't have a Figure area to plot inside, Python first creates a Figure object. And since the Figure doesn't start with any Axes to draw the histogram onto, an Axes object is created inside the Figure. Finally, the histogram is drawn within that Axes.



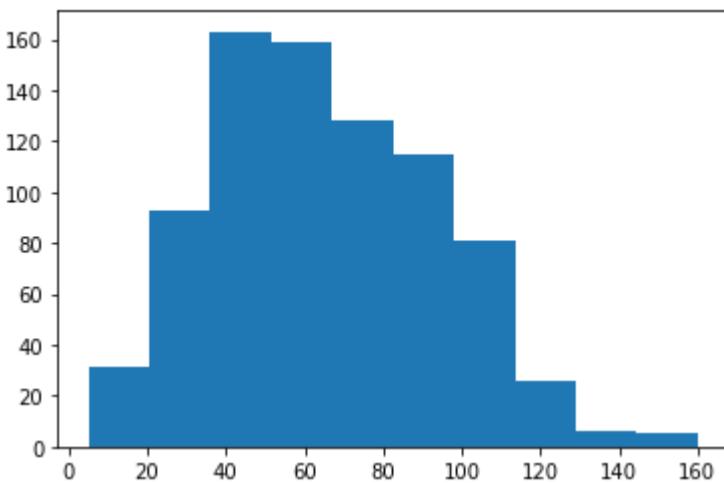
This hierarchy of objects is useful to know about so that we can take more control over the layout and aesthetics of our plots. One alternative way we could have created the histogram is to explicitly set up the Figure and Axes like this:

### Example 2. Demonstrate `figure.add_axes()` and `axes.hist()`

```
Create a new figure
fig = plt.figure()

The argument of add_axes represents the dimensions [left, bottom, width, height] of the new axes.
All quantities are in fractions of figure width and height.
ax = fig.add_axes([.125, .125, .775, .755])
ax.hist(data=pokemon, x='speed');
```

`figure()` creates a new Figure object, a reference to which has been stored in the variable `fig`. One of the Figure methods is `.add_axes()`, which creates a new Axes object in the Figure. The method requires one list as argument specifying the dimensions of the Axes: the first two elements of the list indicate the position of the lower-left hand corner of the Axes (in this case one quarter of the way from the lower-left corner of the Figure) and the last two elements specifying the Axes width and height, respectively. We refer to the Axes in the variable `ax`. Finally, we use the Axes method `.hist()` just like we did before with `plt.hist()`.



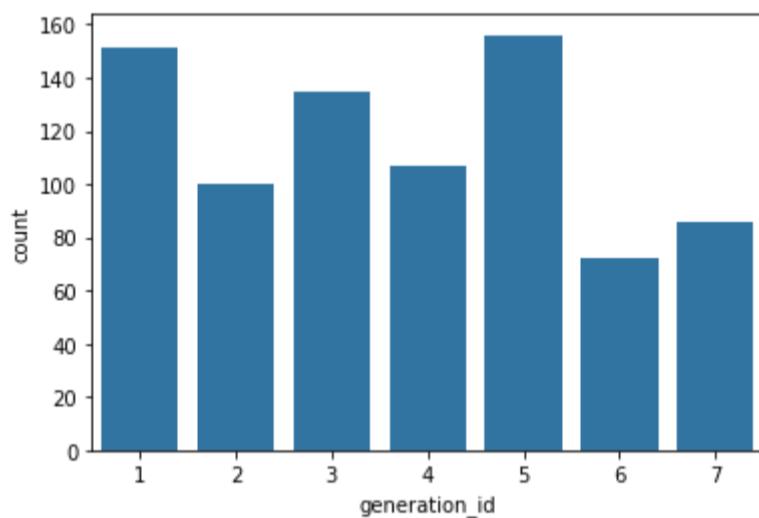
This plot is just like the first histogram on the Histograms page.

---

To use Axes objects with seaborn, seaborn functions usually have an "ax" parameter to specify upon which Axes a plot will be drawn.

## Example 2. Use axes with seaborn.countplot()

```
fig = plt.figure()
ax = fig.add_axes([.125, .125, .775, .755])
base_color = sb.color_palette()[0]
sb.countplot(data = pokemon, x = 'generation_id', color = base_color, ax = ax)
```



This is the same as the second plot on the Bar Charts page.

---

In the above two cases, there was no purpose to explicitly go through the Figure and Axes creation steps. And indeed, in most cases, you can just use the basic matplotlib and seaborn functions as is. Each function targets a Figure or Axes, and they'll automatically target the most recent Figure or Axes worked with. As an example of this, let's review in detail how `subplot()` was used on the Histograms page:

## Example 3. Sub-plots

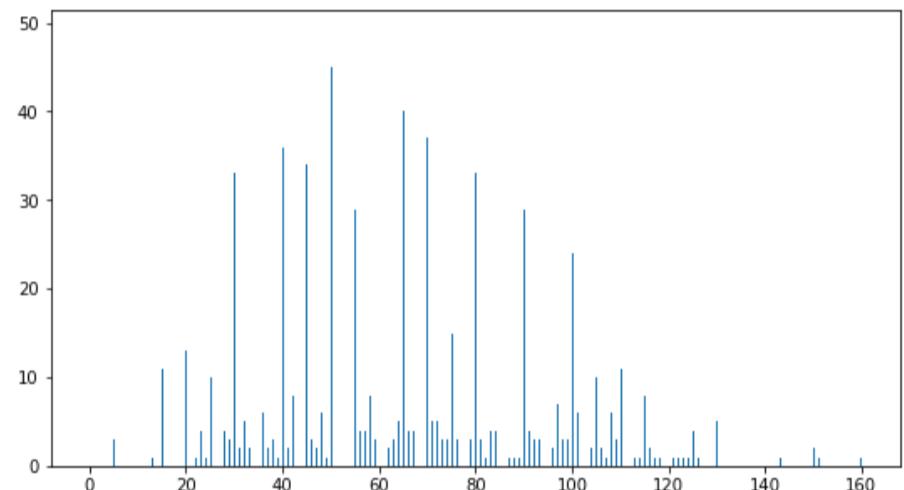
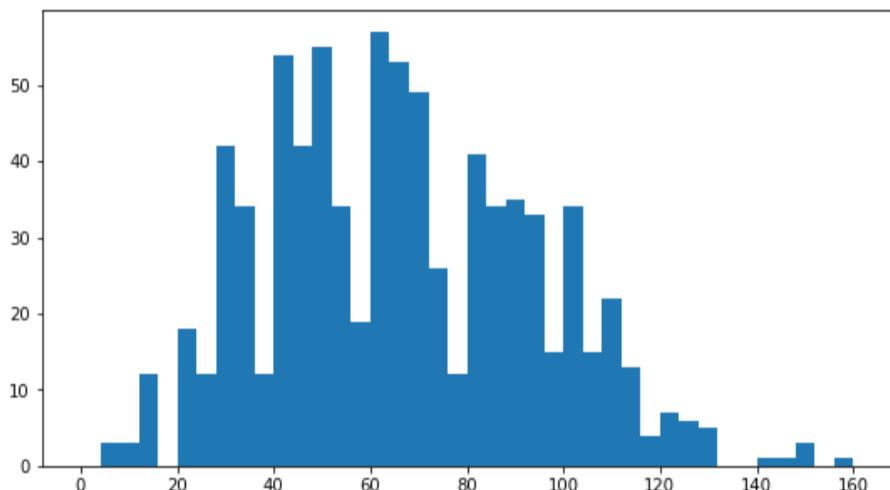
```

Resize the chart, and have two plots side-by-side
set a larger figure size for subplots
plt.figure(figsize = [20, 5])

histogram on left, example of too-large bin size
1 row, 2 cols, subplot 1
plt.subplot(1, 2, 1)
bins = np.arange(0, pokemon['speed'].max()+4, 4)
plt.hist(data = pokemon, x = 'speed', bins = bins);

histogram on right, example of too-small bin size
plt.subplot(1, 2, 2) # 1 row, 2 cols, subplot 2
bins = np.arange(0, pokemon['speed'].max()+1/4, 1/4)
plt.hist(data = pokemon, x = 'speed', bins = bins);

```



First of all, `plt.figure(figsize = [20, 5])` creates a new Figure, with the "figsize" argument setting the width and height of the overall figure to 20 inches by 5 inches, respectively. Even if we don't assign any variable to return the function's output, Python will still implicitly know that further plotting calls that need a Figure will refer to that Figure as the active one.

Then, `plt.subplot(1, 2, 1)` creates a new Axes in our Figure, its size determined by the `subplot()` function arguments. The first two arguments say to divide the figure into one row and two columns, and the third argument says to create a new Axes in the first slot. Slots are numbered from left to right in rows from top to bottom. Note in particular that the index numbers start at 1 (rather than the usual Python indexing starting from 0). (You'll see the indexing a little better in the example at the end of the page.) Again, Python will implicitly set that Axes as the current Axes, so when the `plt.hist()` call comes, the histogram is plotted in the left-side subplot.

Finally, `plt.subplot(1, 2, 2)` creates a new Axes in the second subplot slot, and sets that one as the current Axes. Thus, when the next `plt.hist()` call comes, the histogram gets drawn in the right-side subplot.

## Subplots Exploration Quiz

**[Multiple Choice Question]** - What if we remove one statement `plt.subplot(1, 2, 2)` from the above code block, and just ran the rest of the lines? What would the outcome plot look like? (**HINT:** Try playing around with some code for yourself to come up with an answer!)

- We would see only one set of bars, for the first `.hist()` call.

We would see only one set of bars, for the second `.hist()` call.

We would see two sets of bars, plotted one on top of the other.

We would see one set of axes, occupying the left side of the figure.

We would see one set of axes, occupying the right side of the figure.

We would see one set of axes, filling the full figure length.

## Additional Techniques

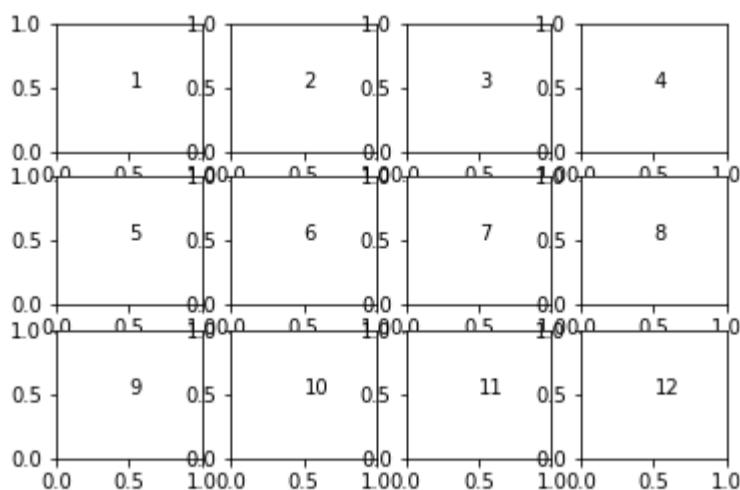
To close this page, we'll quickly run through a few other ways of dealing with Axes and subplots. The techniques above should suffice for basic plot creation, but you might want to keep the following in the back of your mind as additional tools to break out as needed.

If you don't assign Axes objects as they're created, you can retrieve the current Axes using `ax = plt.gca()`, or you can get a list of all Axes in a Figure `fig` by using `axes = fig.get_axes()`. As for creating subplots, you can use `fig.add_subplot()` in the same way as `plt.subplot()` above. If you already know that you're going to be creating a bunch of subplots, you can use the `plt.subplots()` function:

### Example 4. Demonstrate pyplot.sca() and pyplot.text() to generate a grid of subplots

```
fig, axes = plt.subplots(3, 4) # grid of 3x4 subplots
axes = axes.flatten() # reshape from 3x4 array into 12-element vector
for i in range(12):
 plt.sca(axes[i]) # set the current Axes
 plt.text(0.5, 0.5, i+1) # print conventional subplot index number to middle of Axes
```

As a special note for the text, the Axes limits are [0,1] on each Axes by default, and we increment the iterator counter `i` by 1 to get the subplot index, if we were creating the subplots through `subplot()`. (Reference: `plt.sca()`, `plt.text()`)



## Documentation

Documentation pages for Figure and Axes objects are linked below. Note that they're pretty dense, so I don't suggest reading them until you need to dig down deeper and override matplotlib or seaborn's default behavior. Even then, they *are* just reference pages, so they're better for skimming or searching in case other internet resources don't provide enough detail.

- [Figure](#)
- [Axes](#)

---

[Next Concept](#)

# ≡ 11. Choosing a Plot for Discrete Data

## Choosing a Plot for Discrete Data

If you want to plot a **discrete quantitative variable**, it is possible to select either a histogram or a bar chart to depict the data.

- Here, the **discrete** means [non-continuous](#) values. In general, a discrete variable can be assigned to any of the limited (countable) set of values from a given set/range, for example, the number of family members, number of football matches in a tournament, number of departments in a university.
- The **quantitative** term shows that it is the outcome of the measurement of a quantity.

The histogram is the most immediate choice since the data is numeric, but there's one particular consideration to make regarding the bin edges. Since data points fall on set values (bar-width), it can help to reduce ambiguity by putting bin edges between the actual values taken by the data.

## An example describing the ambiguity

For example, assume a given bar falls in a range [10-20], and there is an observation with value 20. This observation will lie on the *next* bar because the given range [10-20] does not include the upper limit 20. Therefore, your readers may not know that values on bin edges end up in the bin to their right, so this can bring potential confusion when they interpret the plot.

Compare the two visualizations of 100 random die rolls below (in `die_rolls`), with bin edges *falling on* the observation values in the left subplot, and bin edges *in between* the observation values in the right subplot.

## Preparatory Step

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sb
%matplotlib inline

die_rolls = pd.read_csv('die_rolls.csv')

A fair dice has six-faces having numbers [1-6].
There are 100 dices, and two trials were conducted.
In each trial, all 100 dices were rolled down, and the outcome [1-6] was recorded.
The `Sum` column represents the sum of the outcomes in the two trials, for each given dice.
die_rolls.head(10)
```

The `die_rolls.csv` file is available to download at the bottom of this page.

---

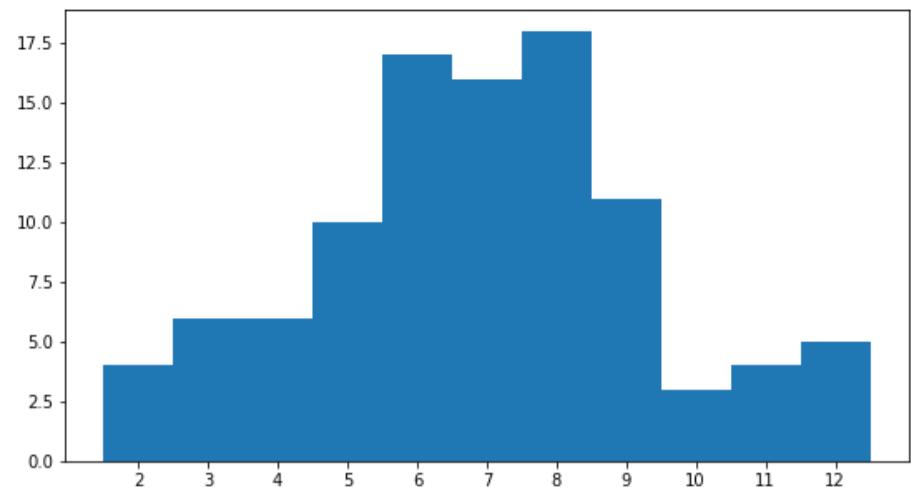
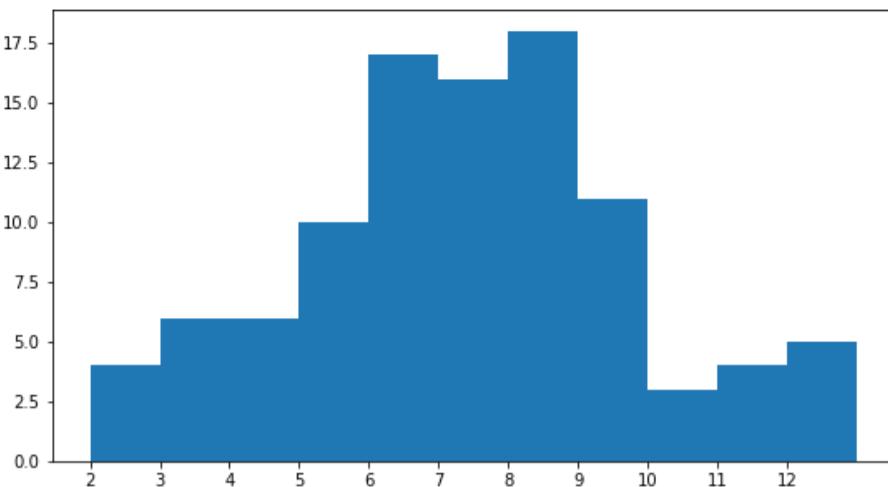
Dice	Trial 1	Trial 2	Sum	
0	1	4	1	5
1	2	4	5	9
2	3	2	6	8
3	4	6	3	9
4	5	3	6	9
5	6	6	6	12
6	7	3	3	6
7	8	3	2	5
8	9	2	6	8
9	10	6	6	12

**Example 1. Shifting the edges of the bars can remove ambiguity in the case of Discrete data**

```
plt.figure(figsize = [20, 5])

Histogram on the left, bin edges on integers
plt.subplot(1, 2, 1)
bin_edges = np.arange(2, 12+1.1, 1) # note `+1.1`, see below
plt.hist(data=die_rolls, x='Sum', bins = bin_edges);
plt.xticks(np.arange(2, 12+1, 1));

Histogram on the right, bin edges between integers
plt.subplot(1, 2, 2)
bin_edges = np.arange(1.5, 12.5+1, 1)
plt.hist(data=die_rolls, x='Sum', bins = bin_edges);
plt.xticks(np.arange(2, 12+1, 1));
```



The same data is plotted in both subplots, but the alignment of the bin edges is different.

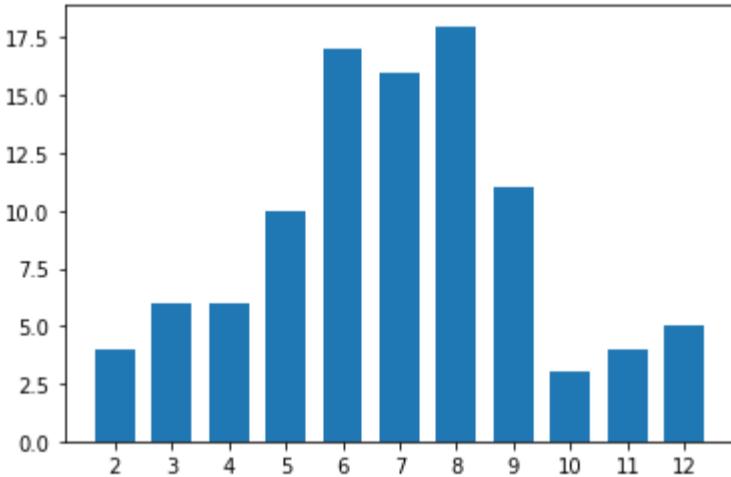
You'll notice for the left histogram, in a deviation from the examples that have come before, I've added 1.1 to the max value (12) for setting the bin edges, rather than just the desired bin width of 1. Recall that data that is equal to the rightmost bin edge gets lumped in to the last bin. This presents a potential problem in perception when a lot of data points take the maximum value, and so is especially relevant when the data takes on discrete values. The 1.1 adds an additional bin to the end to store the die rolls of value 12 alone, to avoid having the last bar catch both 11 and 12.

Alternatively to the histogram, consider if a bar chart with non-connected bins might serve your purposes better. The plot below takes the code from before, but adds the "rwidth" parameter to set the proportion of the bin widths that will be filled by each histogram bar.

## Example 2. Making gaps between individual bars

```
bin_edges = np.arange(1.5, 12.5+1, 1)
plt.hist(data=die_rolls, x='Sum', bins = bin_edges, rwidth = 0.7)
plt.xticks(np.arange(2, 12+1, 1));
```

With "rwidth" set to 0.7, the bars will take up 70% of the space allocated by each bin, with 30% of the space left empty. This changes the default display of the histogram (which you could think of as "rwidth = 1") into a bar chart.



Gaps between bars makes it clear that the data is discrete in nature.

By adding gaps between bars, you emphasize the fact that the data is discrete in value. On the other hand, plotting your quantitative data in this manner might cause it to be interpreted as ordinal-type data, which can have an effect on overall perception.

For continuous numeric data, you should not make use of the "rwidth" parameter, since the gaps imply discreteness of value. As another caution, it might be tempting to use seaborn's `countplot` function to plot the distribution of a discrete numeric variable as bars. Be careful about doing this, since each unique numeric value will get a bar, regardless of the spacing in values between bars. (For example, if the unique values were {1, 2, 4, 5}, missing 3, `countplot` would only plot four bars, with the bars for 2 and 4 right next to one another.) Also, even if your data is technically discrete numeric, you should probably not consider either of the variants depicted on this

page unless the number of unique values is small enough to allow for the half-unit shift or discrete bars to be interpretable. If you have a large number of unique values over a large enough range, it's better to stick with the standard histogram than risk interpretability issues.

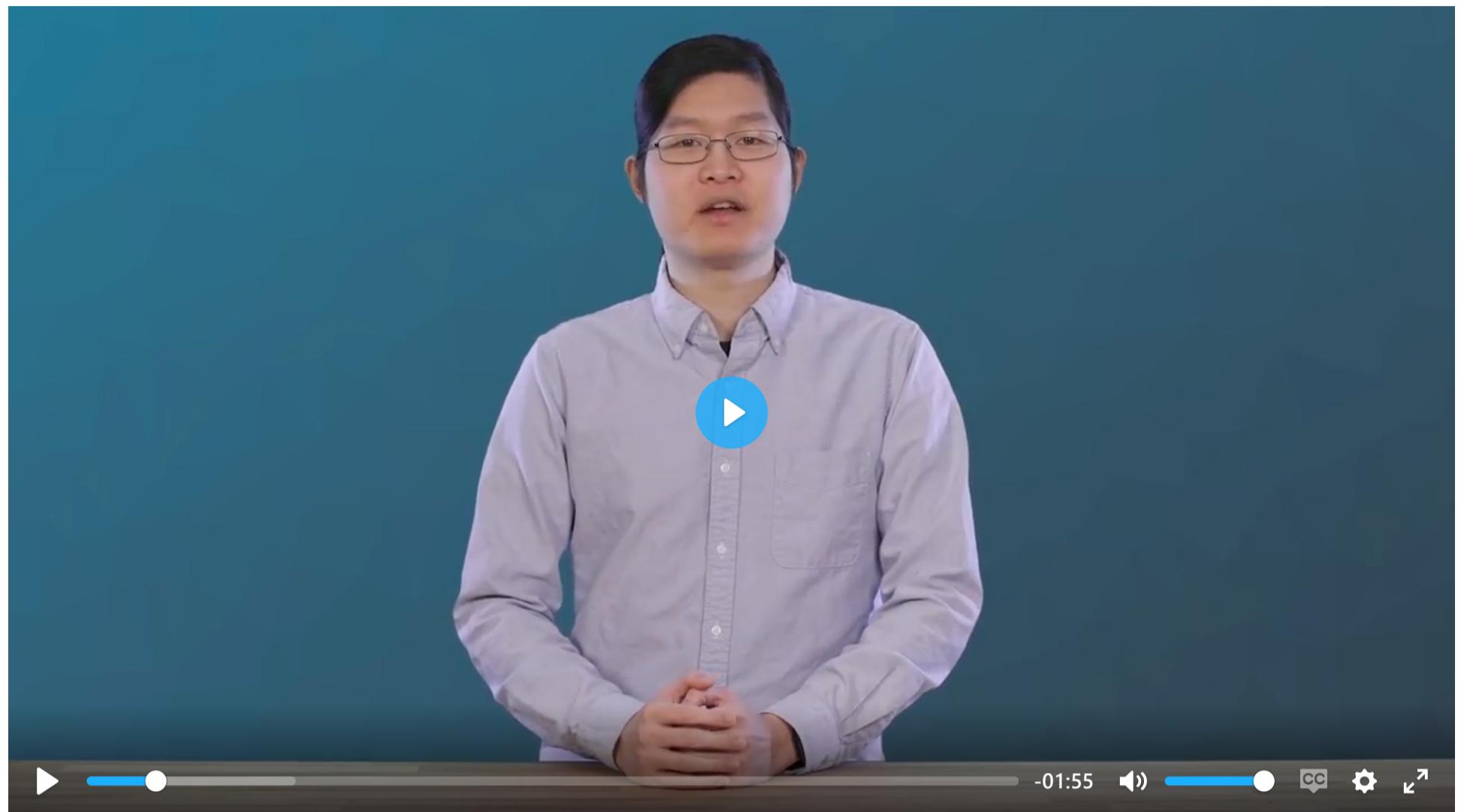
While you might justify plotting discrete numeric data using a bar chart, you'll be less apt to justify the opposite: plotting ordinal data as a histogram. The space between bars in a bar chart helps to remind the reader that values are not contiguous in an 'interval'-type fashion: only that there is an order in levels. With that space removed as in a histogram, it's harder to remember this important bit of interpretation.

---

[Next Concept](#)

## ≡ 12. Descriptive Statistics, Outliers and Axis Limits

L3 111 Descriptive Stats Outliers And Axis Limits V2



DataVis L3 11 V1



Out[2]:

	<b>id</b>	<b>identifier</b>	<b>generation_id</b>	<b>height</b>	<b>weight</b>	<b>base_experience</b>	<b>type_1</b>	<b>type_2</b>	<b>hp</b>	<b>attack</b>	<b>defense</b>	<b>speed</b>	<b>special_attack</b>
<b>0</b>	1	bulbasaur	1	0.7	6.9	64	grass	poison	45	49	49	45	65
<b>1</b>	2	ivysaur	1	1.0	13.0	142	grass	poison	60	62	63	60	80
<b>2</b>	3	venusaur	1	2.0	100.0	236	grass	poison	80	82	83	80	100
<b>3</b>	4	charmander	1	0.6	8.5	62	fire	NaN	39	52	43	65	60
<b>4</b>	5	charmeleon	1	1.1	19.0	142	fire	NaN	58	64	58	80	80
<b>5</b>	6	charizard	1	1.7	90.5	240	fire	flying	78	84	78	100	109
<b>6</b>	7	squirtle	1	0.5	9.0	63	water	NaN	44	48	65	43	50
<b>7</b>	8	wartortle	1	1.0	22.5	142	water	NaN	59	63	80	58	65
<b>8</b>	9	blastoise	1	1.6	85.5	239	water	NaN	79	83	100	78	85
<b>9</b>	10	caterpie	1	0.3	2.9	39	bug	NaN	45	30	35	45	20

```
In []: bins = np.arange(0, pokemon['height'].max()+0.5, 0.5)
plt.hist(data = pokemon, x = 'height', bins = bins);
```

In [ ]:

00:51 CC



## Descriptive Statistics, Outliers, and Axis Limits

As you create your plots and perform your exploration, make sure that you pay attention to what the plots tell you that go beyond just the basic descriptive statistics. Note any aspects of the data like the number of modes and skew, and note the presence of outliers in the data for further investigation.

Related to the latter point, you might need to change the limits or scale of what is plotted to take a closer look at the underlying patterns in the data. Let's see a few examples.

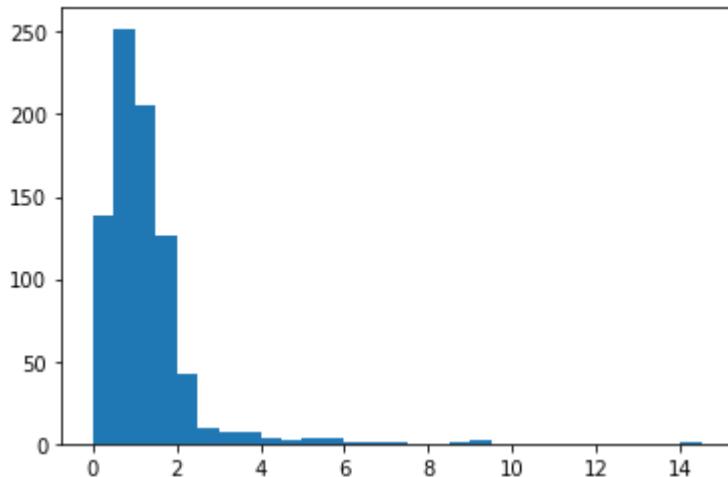
### Example 1. Plot the histogram from the data having a skewed distribution of values

```
TO DO: Necessary import

Load the data, and see the height column
pokemon = pd.read_csv('pokemon.csv')
pokemon.head(10)

Get the ticks for bins between [0-15], at an interval of 0.5
bins = np.arange(0, pokemon['height'].max()+0.5, 0.5)

Plot the histogram for the height column
plt.hist(data=pokemon, x='height', bins=bins);
```



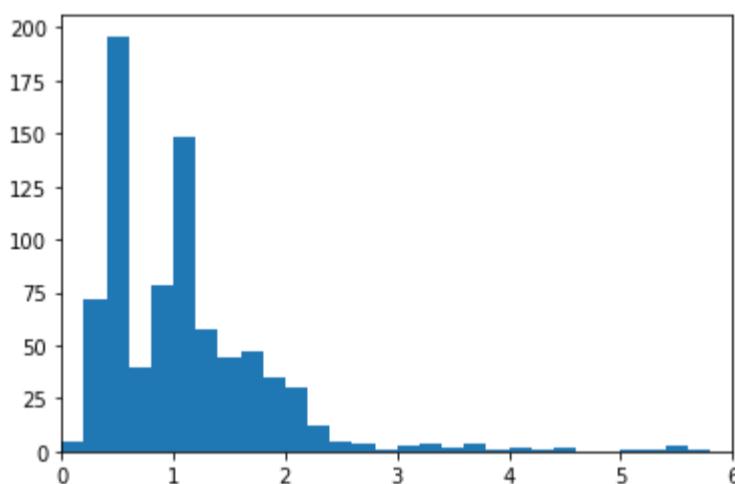
This page covers the topic of axis limits; and the next topic is scales and transformations. In order to change a histogram's axis limits, you can add a Matplotlib `xlim()` call to your code. The function takes a tuple of two numbers specifying the upper and lower bounds of the x-axis range. See the example below.

## Example 2. Plot the histogram with a changed axis limit.

```
Get the ticks for bins between [0-15], at an interval of 0.5
bins = np.arange(0, pokemon['height'].max()+0.2, 0.2)
plt.hist(data=pokemon, x='height', bins=bins);

Set the upper and lower bounds of the bins that are displayed in the plot
Refer here for more information - https://matplotlib.org/3.1.1/api/_as_gen/matplotlib.pyplot.xlim.html
The argument represent a tuple of the new x-axis limits.
plt.xlim((0,6));
```

Alternatively, the `xlim` function can be called with two numeric arguments only, `plt.xlim(0,6)`, to get the same result.



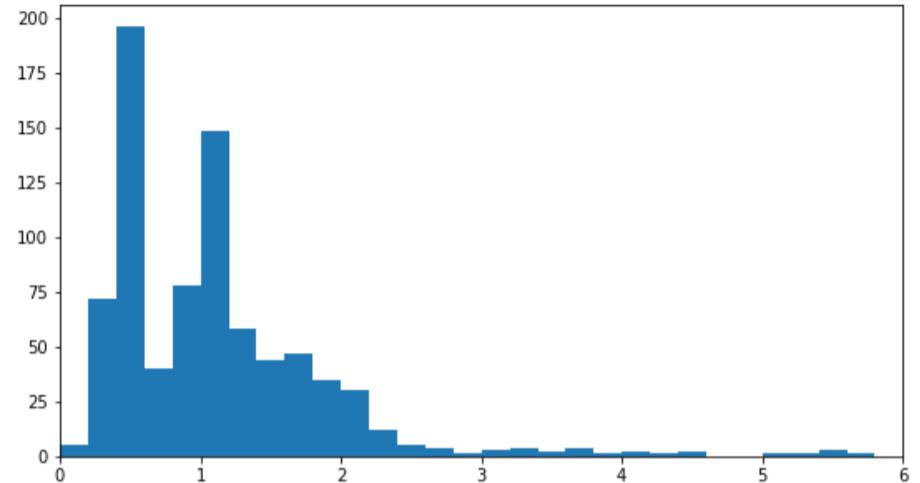
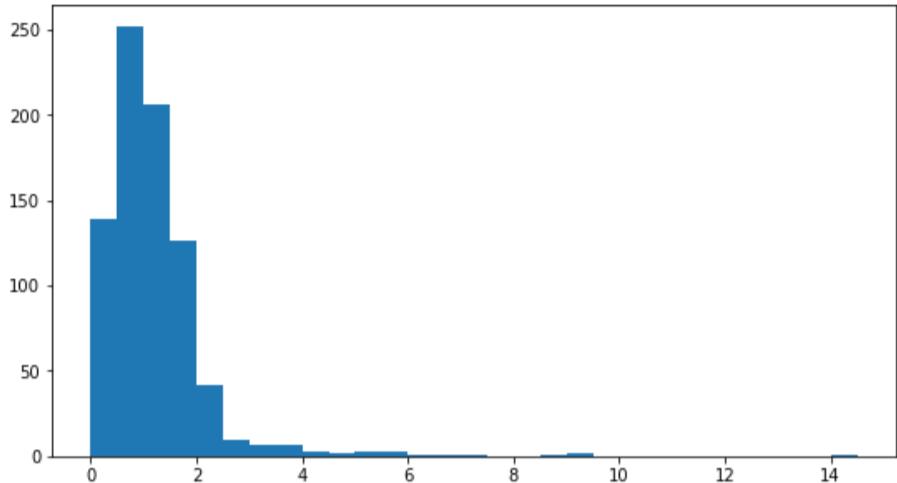
In the generic example above, we might be interested in comparing patterns in other variables between data points that take values less than 6 to those that take values greater than 6. For anything that is concentrated on the bulk of the data in the former group (< 6), the use of axis limits can allow focusing on data points in that range without needing to go through the creation of new DataFrame filtering out the data points in the latter group (> 6).

**TO DO:** Plot the above two graphs in a single figure of size 20 x 5 inches, side-by-side.

**Hint** - Use the steps below:

1. Define the figure size, using `matplotlib.pyplot.figure(figsize = [float, float])`.
2. Add a subplot using `matplotlib.pyplot.subplot(int, int, index)` for the left-graph to the current figure. Then, define the left-graph.
3. Similarly, add a subplot for the right-graph to the current figure. Then, define the right-graph.

The expected output is shown below:



#### Here is one of the possible solutions:

```
Define the figure size
plt.figure(figsize = [20, 5])

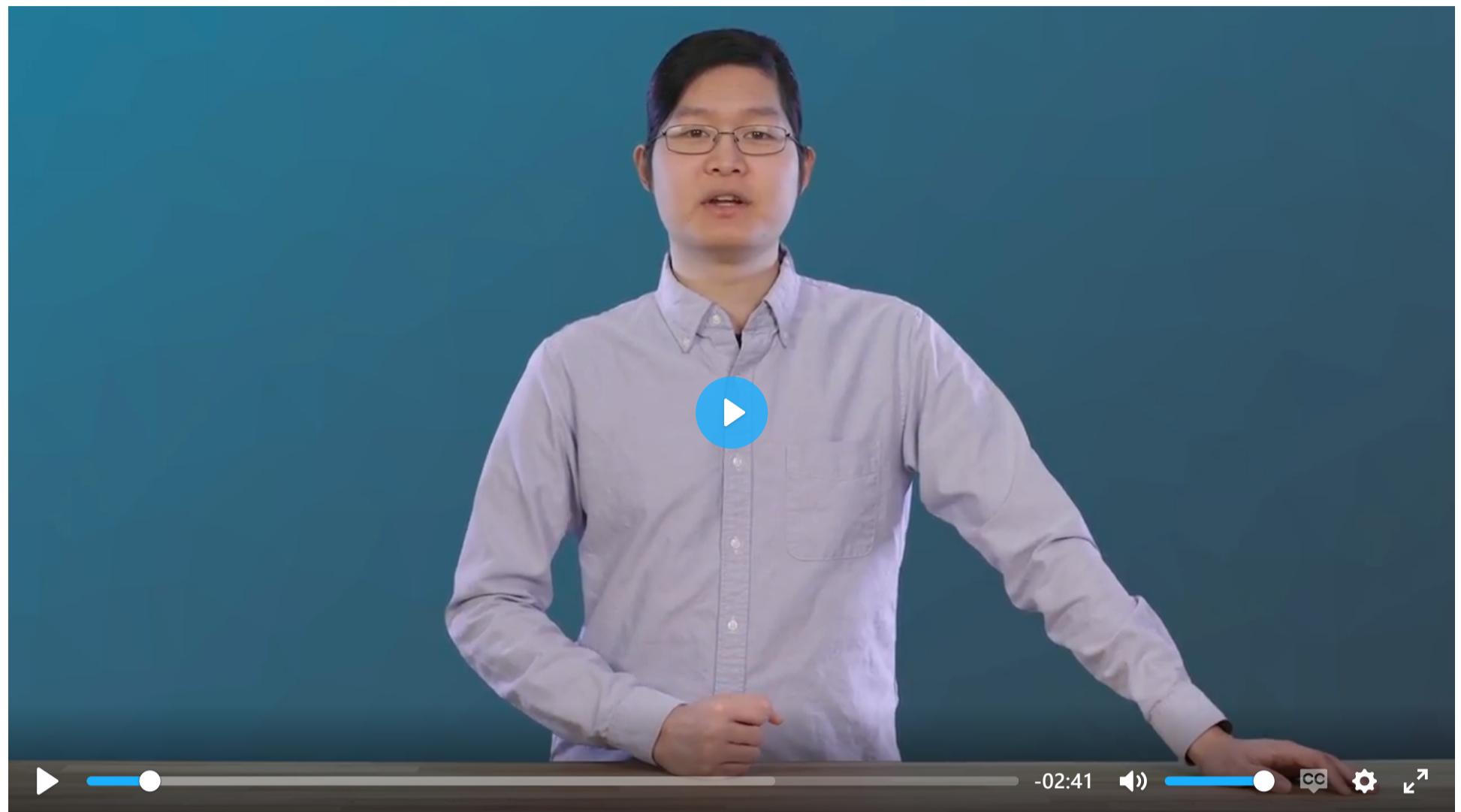
histogram on left: full data
plt.subplot(1, 2, 1)
bin_edges = np.arange(0, pokemon['height'].max()+0.5, 0.5)
plt.hist(data=pokemon, x='height', bins = bin_edges)

histogram on right: focus in on bulk of data < 6
plt.subplot(1, 2, 2)
bin_edges = np.arange(0, pokemon['height'].max()+0.2, 0.2)
plt.hist(data=pokemon, x='height', bins = bin_edges)
plt.xlim(0, 6) # could also be called as plt.xlim((0, 6))
```

[Next Concept](#)

## ≡ 13. Scales and Transformations

L3 121 Scales And Transformations V3

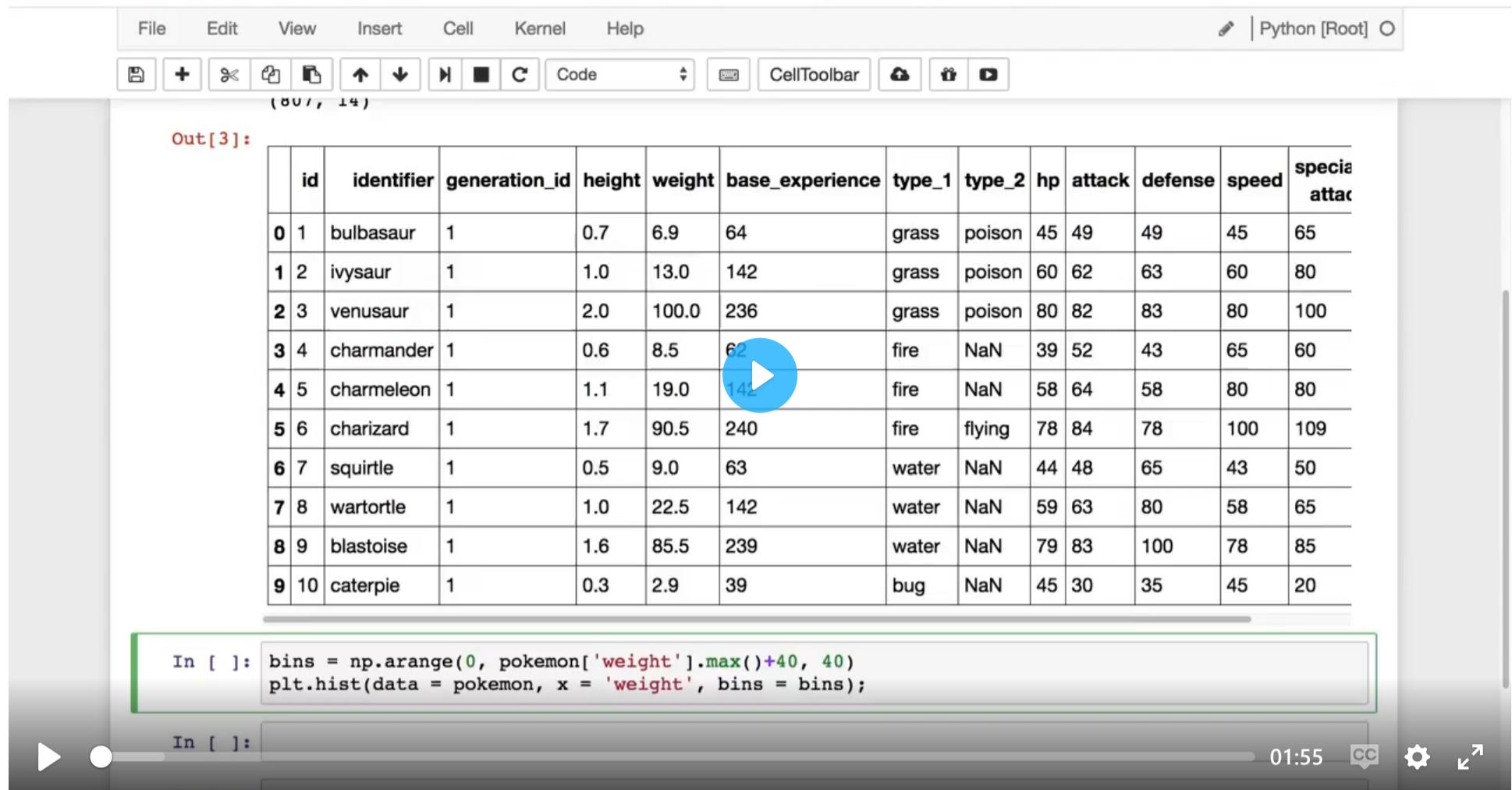


---

---

---

DataVis L3 12 V2



## Scales and Transformations

Certain data distributions will find themselves amenable to scale transformations. The most common example of this is data that follows an approximately [log-normal](#) distribution. This is data that, in their natural units, can look highly skewed: lots of points with low values, with a very long tail of data points with large values. However, after applying a logarithmic transform to the data, the data will follow a normal distribution. (If you need a refresher on the logarithm function, check out [this lesson on Khan Academy](#).)

### Example 1 - Scale the x-axis to log-type

```
Necessary import

pokemon = pd.read_csv('pokemon.csv')
pokemon.head(10)

plt.figure(figsize = [20, 5])

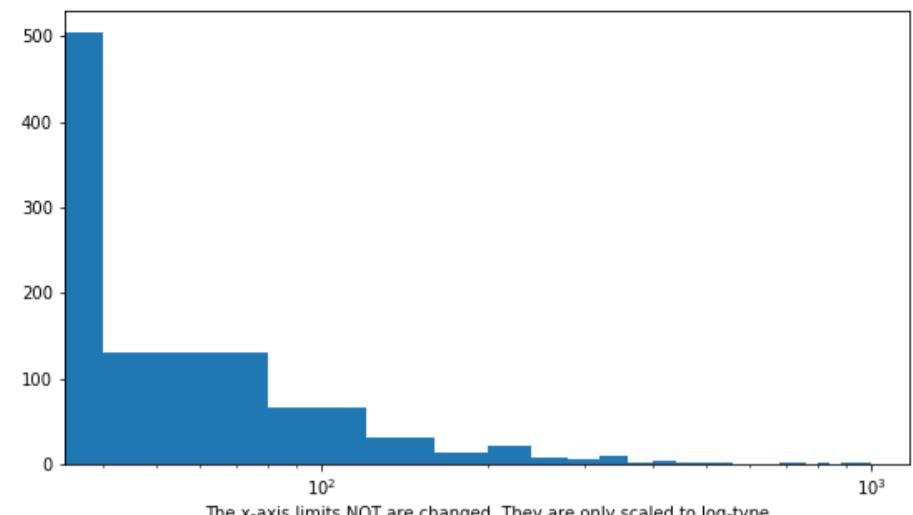
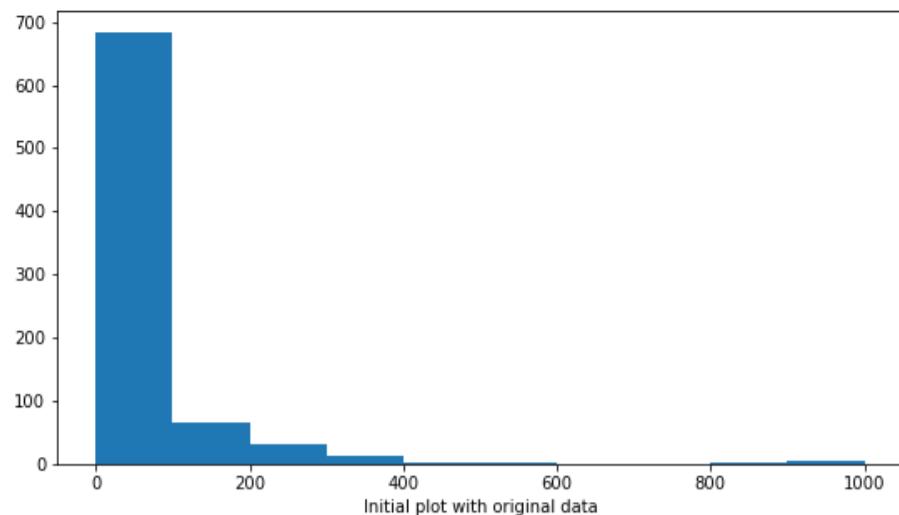
HISTOGRAM ON LEFT: full data without scaling
plt.subplot(1, 2, 1)
plt.hist(data=pokemon, x='weight');
Display a label on the x-axis
plt.xlabel('Initial plot with original data')

HISTOGRAM ON RIGHT
plt.subplot(1, 2, 2)

Get the ticks for bins between [0 - maximum weight]
bins = np.arange(0, pokemon['weight'].max()+40, 40)
plt.hist(data=pokemon, x='weight', bins=bins);

The argument in the xscale() represents the axis scale type to apply.
The possible values are: {"linear", "log", "symlog", "logit", ...}
Refer - https://matplotlib.org/3.1.1/api/_as_gen/matplotlib.pyplot.xscale.html
plt.xscale('log')

plt.xlabel('The x-axis limits NOT are changed. They are only scaled to log-type')
```



```
Describe the data
pokemon['weight'].describe()
```

<b>count</b>	<b>807.000000</b>
<b>mean</b>	<b>61.771128</b>
<b>std</b>	<b>111.519355</b>
<b>min</b>	<b>0.100000</b>
<b>25%</b>	<b>9.000000</b>
<b>50%</b>	<b>27.000000</b>
<b>75%</b>	<b>63.000000</b>
<b>max</b>	<b>999.900000</b>
<b>Name:</b>	<b>weight, dtype: float64</b>

Notice two things about the right histogram of example 1 above, now.

1. Even though the data is on a log scale, the bins are still linearly spaced. This means that they change size from wide on the left to thin on the right, as the values increase multiplicative. Matplotlib's `xscale` function includes a few built-in transformations: we have used the 'log' scale here.
2. Secondly, the default label (x-axis ticks) settings are still somewhat tricky to interpret and are sparse as well.

To address the bin size issue, we just need to change them so that they are evenly-spaced powers of 10. Depending on what you are plotting, a different base power like 2 might be useful instead.

To address the second issue of interpretation of x-axis ticks, the scale transformation is the solution. In a scale transformation, the gaps between values are based on the transformed scale, but you can interpret data in the variable's natural units.

Let's see another example below.

## Example 2 - Scale the x-axis to log-type, and change the axis limit.

```
Transform the describe() to a scale of log10
Documentation: [numpy `log10`](https://docs.scipy.org/doc/numpy/reference/generated/numpy.log10.html)
np.log10(pokemon['weight'].describe())
```

---

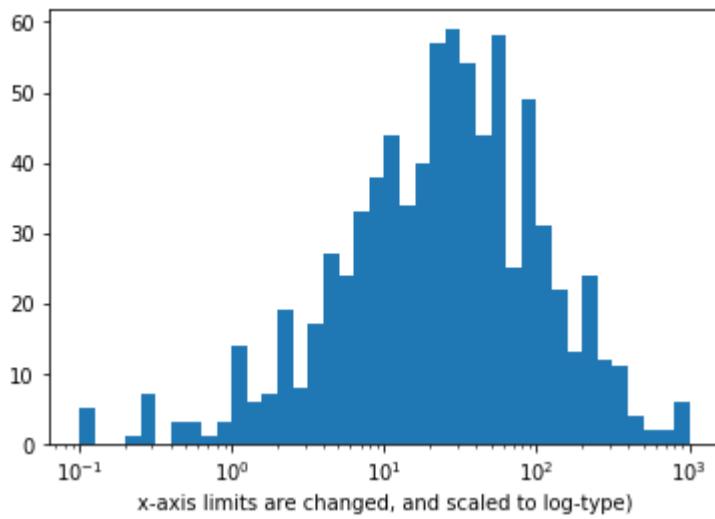
```
count 2.906874
mean 1.790786
std 2.047350
min -1.000000
25% 0.954243
50% 1.431364
75% 1.799341
max 2.999957
Name: weight, dtype: float64
```

---

```
Axis transformation
Bin size
bins = 10 ** np.arange(-1, 3+0.1, 0.1)
plt.hist(data=pokemon, x='weight', bins=bins);

The argument in the xscale() represents the axis scale type to apply.
The possible values are: {"linear", "log", "symlog", "logit", ...}
plt.xscale('log')

Apply x-axis label
Documentation: [matplotlib `xlabel`](https://matplotlib.org/api/_as_gen/matplotlib.pyplot.xlabel.html)
plt.xlabel('x-axis limits are changed, and scaled to log-type')
```



### Example 3 - Scale the x-axis to log-type, change the axis limits, and increase the x-ticks

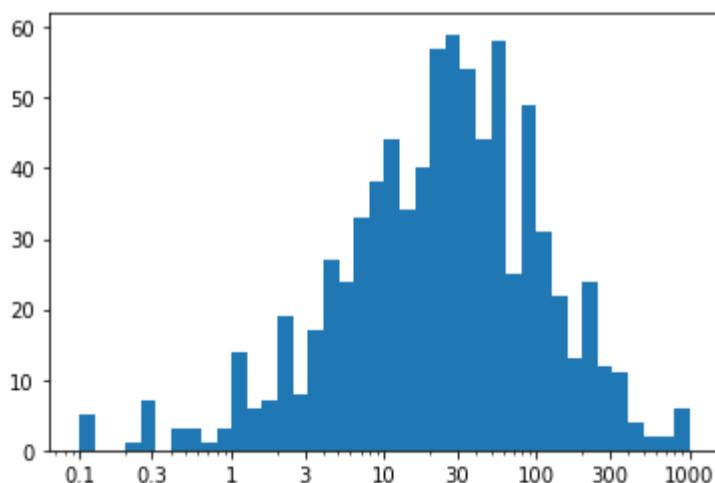
```
Get the ticks for bins between [0 - maximum weight]
bins = 10 ** np.arange(-1, 3+0.1, 0.1)

Generate the x-ticks you want to apply
ticks = [0.1, 0.3, 1, 3, 10, 30, 100, 300, 1000]
Convert ticks into string values, to be displayed along the x-axis
labels = ['{}'.format(v) for v in ticks]

Plot the histogram
plt.hist(data=pokemon, x='weight', bins=bins);

The argument in the xscale() represents the axis scale type to apply.
The possible values are: {"linear", "log", "symlog", "logit", ...}
plt.xscale('log')

Apply x-ticks
plt.xticks(ticks, labels);
```



**Observation** - We've ended up with the same plot as when we performed the direct log transform, but now with a much nicer set of tick marks and labels.

For the ticks, we have used [`xticks\(\)`](#) to specify locations and labels in their natural units. Remember: we aren't changing the values taken by the data, only how they're displayed. Between integer powers of 10, we don't have clean values for even markings, but we can still get close. Setting ticks in cycles of 1-3-10 or 1-2-5-10 are very useful for base-10 log transforms.

It is important that the `xticks` are specified *after* `xscale` since that function has its own built-in tick settings.

## Alternative Approach

Be aware that a logarithmic transformation is not the only one possible. When we perform a logarithmic transformation, our data values have to all be positive; it's impossible to take a log of zero or a negative number. In addition, the transformation implies that additive steps on the log scale will result in multiplicative changes in the natural scale, an important implication when it comes to data modeling. The type of transformation that you choose may be informed by the context for the data. For example, [this Wikipedia section](#) provides a few examples of places where log-normal distributions have been observed.

If you want to use a different transformation that's not available in `xscale`, then you'll have to perform some feature engineering. In cases like this, we want to be systematic by writing a function that applies both the transformation and its inverse. The inverse will be useful in cases where we specify values in their transformed units and need to get the natural units back. For the purposes of demonstration, let's say that we want to try plotting the above data on a square-root transformation. (Perhaps the numbers represent areas, and we think it makes sense to model the data on a rough estimation of radius, length, or some other 1-d dimension.) We can create a visualization on this transformed scale like this:

### Example 4. Custom scaling the given data Series, instead of using the built-in log scale

```
def sqrt_trans(x, inverse = False):
 """ transformation helper function """
 if not inverse:
 return np.sqrt(x)
 else:
 return x ** 2

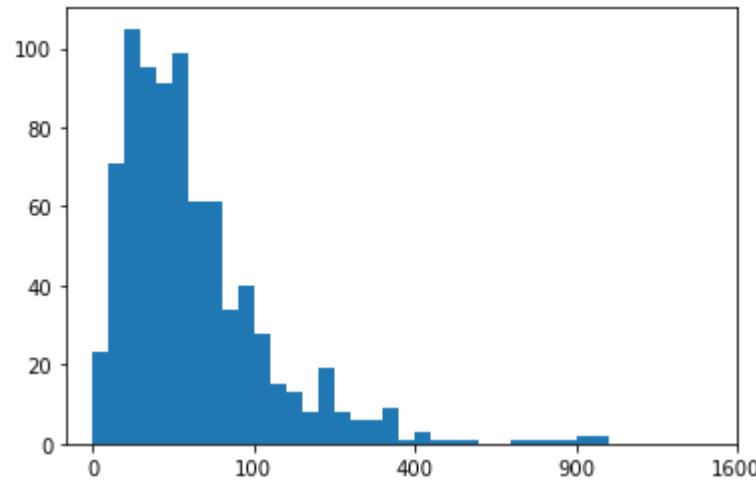
Bin resizing, to transform the x-axis
bin_edges = np.arange(0, sqrt_trans(pokemon['weight'].max())+1, 1)

Plot the scaled data
plt.hist(pokemon['weight'].apply(sqrt_trans), bins = bin_edges)

Identify the tick-locations
tick_locs = np.arange(0, sqrt_trans(pokemon['weight'].max())+10, 10)

Apply x-ticks
plt.xticks(tick_locs, sqrt_trans(tick_locs, inverse = True).astype(int));
```

Note that `data` is a pandas Series, so we can use the `apply` method for the function. If it were a NumPy Array, we would need to apply the function like in the other cases. The tick locations could have also been specified with the natural values, where we apply the standard transformation function on the first argument of `xticks` instead. The output transformed-histogram is shown below:



Histogram based on the custom scaling the given data Series

[Next Concept](#)

## ≡ 16. Extra: Kernel Density Estimation

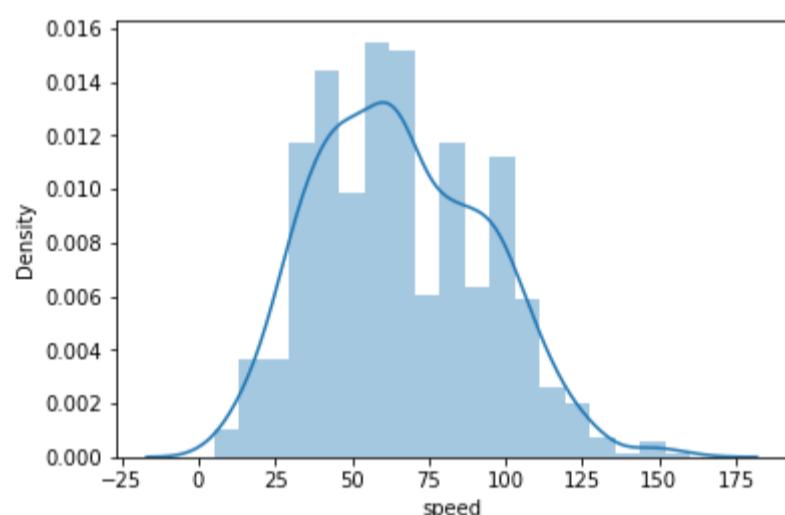
At the end of this lesson and the next, you'll find some extra concepts that didn't really fit into the main lesson flow. These concepts will cover a few additional univariate plots that you might be interested in using or might observe in your own research. While bar charts and histograms should cover your most common needs, the plots in this section might prove useful for both the exploratory and explanatory sides of data visualization.

### Kernel Density Estimation

Earlier in this lesson, you saw an example of [kernel density estimation](#) (KDE) through the use of seaborn's `distplot` function, which plots a KDE on top of a histogram.

#### Example 1. Plot the Kernel Density Estimation (KDE)

```
The pokemon dataset is available to download at the bottom of this page.
sb.distplot(pokemon['speed']);
```

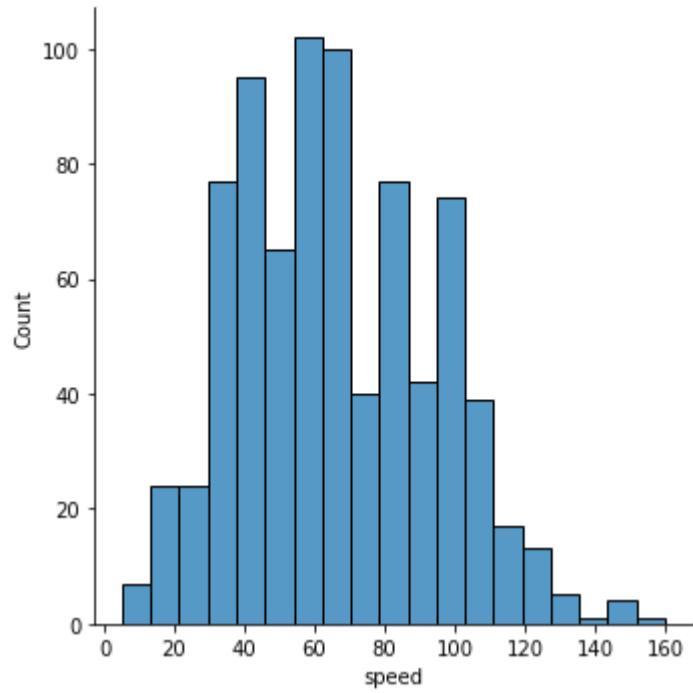


**Note** - The `distplot()` function is deprecated in Seaborn v 0.11.0, and will be removed in a future version. The alternative is either of the following:

1. [displot\(\)](#) - A figure-level function with similar flexibility.
2. [histplot\(\)](#) - An axes-level function for histograms.

See the same example with newer `displot()` function:

```
Use this new function only with Seaborn 0.11.0 and above.
The kind argument can take any one value from {"hist", "kde", "ecdf"}.
sb.displot(pokemon['speed'], kind='hist');
Use the 'kde' kind for kernel density estimation
sb.displot(pokemon['speed'], kind='kde');
```



---

Kernel density estimation is one way of estimating the probability density function of a variable. In a KDE plot, you can think of each observation as replaced by a small 'lump' of area. Stacking these lumps all together produces the final density curve. The default settings use a normal-distribution kernel, but most software that can produce KDE plots also include other kernel function options.

Seaborn's `distplot` function calls another function, `kdeplot`, to generate the KDE. The demonstration code below also uses a third function called by `distplot` for illustration, `rugplot()`. In a rugplot, data points are depicted as dashes on a number line.

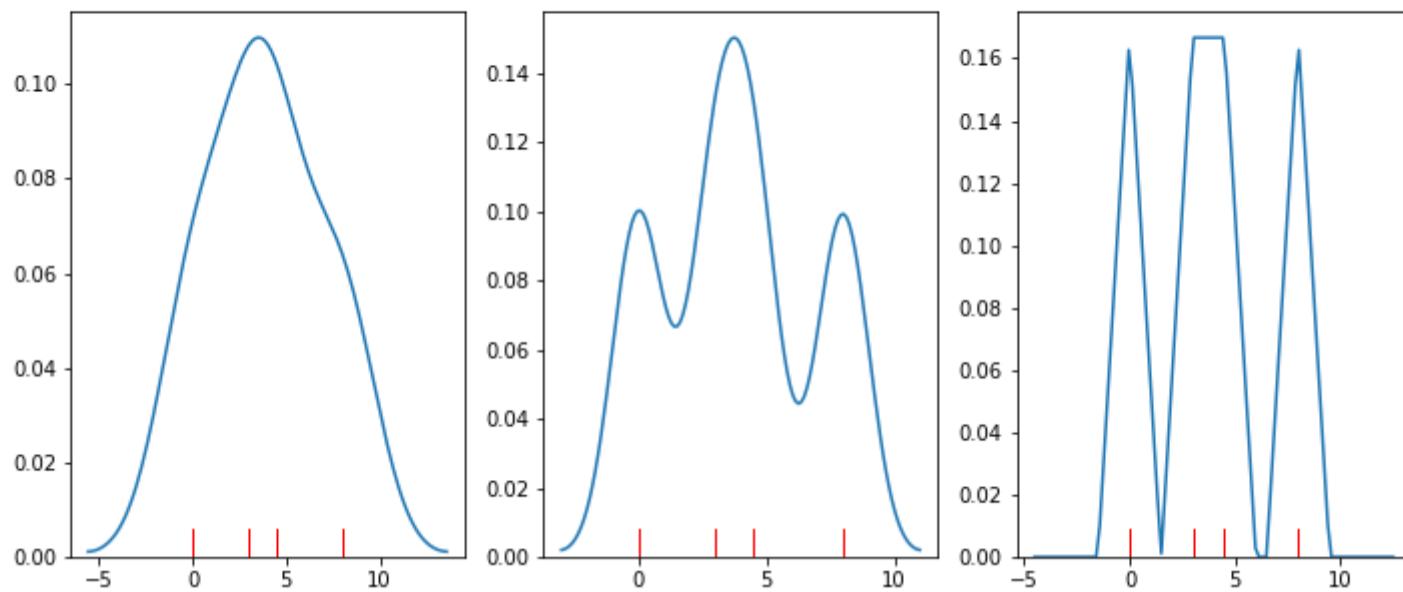
## Example 2. Demonstrating `distplot()` and `rugplot()` to plot the KDE

```
data = [0.0, 3.0, 4.5, 8.0]
plt.figure(figsize = [12, 5])

left plot: showing kde lumps with the default settings
plt.subplot(1, 3, 1)
sb.distplot(data, hist = False, rug = True, rug_kws = {'color' : 'r'})

central plot: kde with narrow bandwidth to show individual probability lumps
plt.subplot(1, 3, 2)
sb.distplot(data, hist = False, rug = True, rug_kws = {'color' : 'r'},
 kde_kws = {'bw' : 1})

right plot: choosing a different, triangular kernel function (lump shape)
plt.subplot(1, 3, 3)
sb.distplot(data, hist = False, rug = True, rug_kws = {'color' : 'r'},
 kde_kws = {'bw' : 1.5, 'kernel' : 'tri'})
```



Interpreting proportions from this plot type is slightly trickier than a standard histogram: the vertical axis indicates a density of data rather than straightforward proportions. Under a KDE plot, the total area between the 0-line and the curve will be 1. The probability of an outcome falling between two values is found by computing the area under the curve that falls between those values. Making area judgments like this without computer assistance is difficult and likely to be inaccurate.

Despite the fact that making specific probability judgments are not as intuitive with KDE plots as histograms, there are still reasons to use kernel density estimation. If there are relatively few data points available, KDE provides a smooth estimate of the overall distribution of data. These ideas may not be so easily conveyed through histograms, in which the large discreteness of jumps may end up misleading.

It should also be noted that there is a bandwidth parameter in KDE that specifies how wide the density lumps are. Similar to bin width for histograms, we need to choose a bandwidth size that best shows the signal in the data. A too-small bandwidth can make the data look noisier than it really is, and a too-large bandwidth can smooth out useful features that we could use to make inferences about the data. It's good to keep this in mind in case the default bandwidths chosen by your visualization software don't look quite right or if you need to perform further investigations.

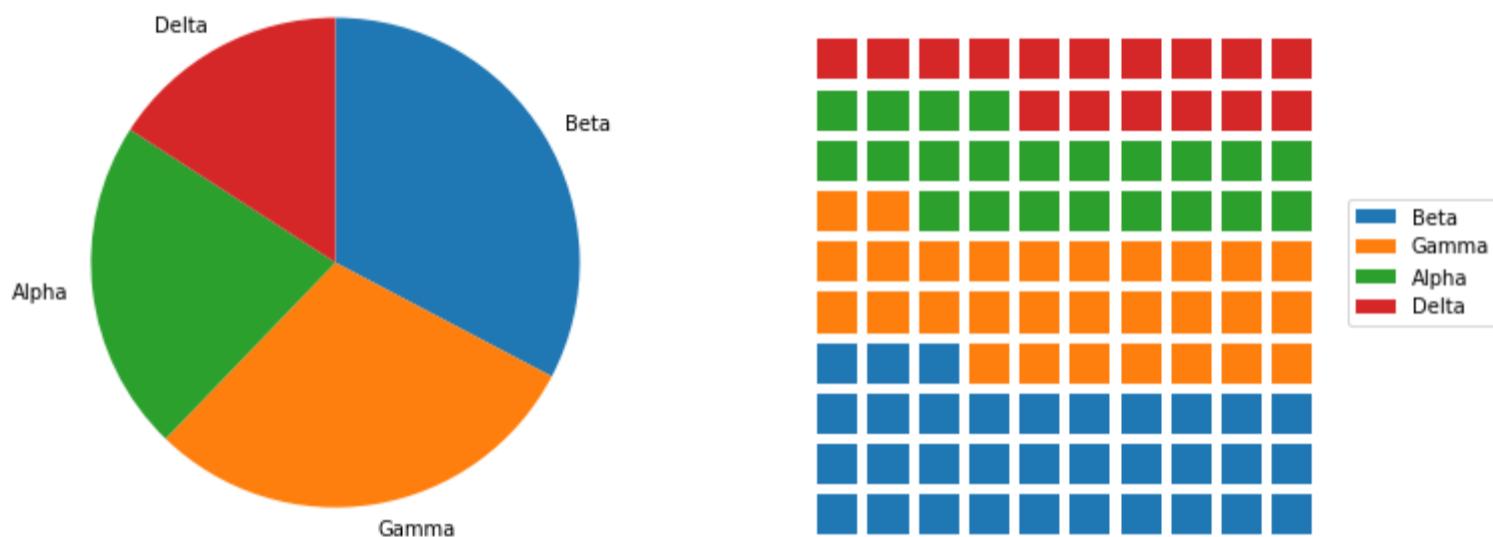
---

[Next Concept](#)

## ≡ 17. Extra: Waffle Plots

### Waffle Plots

One alternative univariate plot type that you might see for categorical data is the **waffle plot**, also known as the square pie chart. While the standard pie chart uses a circle to represent the whole, a waffle plot is plotted onto a square divided into a 10x10 grid. Each small square in the grid represents one percent of the data, and a number of squares are colored by category to indicate total proportions. Compared to a pie chart, it is much easier to make precise assessments of relative frequencies.



You've seen code for the pie chart (left) previously in the lesson. Code for the waffle plot (right) will be walked through below.

There's no built-in function for waffle plots in Matplotlib or Seaborn, so we'll need to take some additional steps in order to build one with the tools available. First, we need to create a function to decide how many blocks to allocate to each category. The function below, `percentage_blocks`, uses a rule where each category gets a number of blocks equal to the number of full percentage points it covers. The remaining blocks to get to the full one hundred are assigned to the categories with the largest fractional parts.

```
def percentage_blocks(df, var):
 """
 Take as input a dataframe and variable, and return a Pandas series with
 approximate percentage values for filling out a waffle plot.
 """
 # compute base quotas
 percentages = 100 * df[var].value_counts() / df.shape[0]
 counts = np.floor(percentages).astype(int) # integer part = minimum quota
 decimal = (percentages - counts).sort_values(ascending = False)

 # add in additional counts to reach 100
 rem = 100 - counts.sum()
 for cat in decimal.index[:rem]:
 counts[cat] += 1

 return counts
```

<code>df['cat_var'].value_counts()</code>	<code>df['cat_var'].value_counts() / df.shape[0]</code>	<code>percentage_blocks(df, 'cat_var')</code>
Beta 103 Gamma 93 Alpha 69 Delta 50 Name: cat_var, dtype: int64	Beta 0.326984 Gamma 0.295238 Alpha 0.219048 Delta 0.158730 Name: cat_var, dtype: float64	Beta 33 Gamma 29 Alpha 22 Delta 16 Name: cat_var, dtype: int64

Note that if we just rounded the proportions (center), we would round all of them up, ending up with a total of 101 blocks.

---

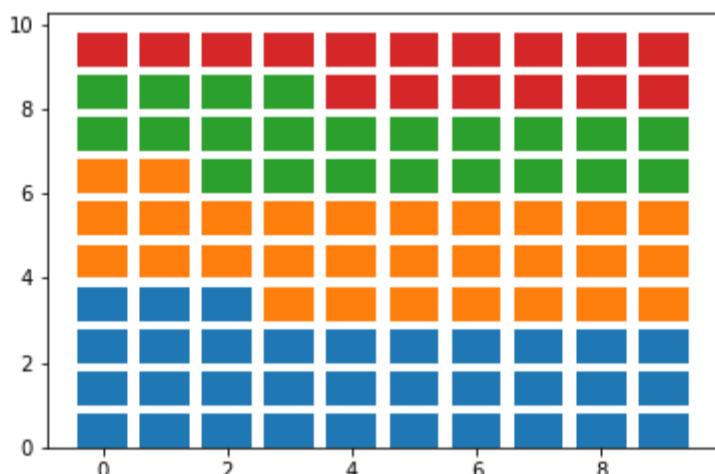
Now it's time to actually plot those counts as boxes in the waffle plot form. To do this, we'll make use of Matplotlib's `bar` function. We could have used this function earlier in the lesson to create our bar charts instead of Seaborn's `countplot`, but it would have required us to aggregate the data first to get the height of each bar. For the case of the waffle plot, we're going to specify the x- and y-coordinates of the boxes, and set their widths and heights to be equal, to create squares. The initial plotting code looks like this:

```
waffle_counts = percentage_blocks(df, 'cat_var')

prev_count = 0
for each category,
for cat in range(waffle_counts.shape[0]):
 # get the block indices
 blocks = np.arange(prev_count, prev_count + waffle_counts[cat])
 # and put a block at each index's location
 x = blocks % 10 # use mod operation to get ones digit
 y = blocks // 10 # use floor division to get tens digit
 plt.bar(x = x, height = 0.8, width = 0.8, bottom = y)
 prev_count += waffle_counts[cat]
```

The blocks are drawn from left to right, bottom to top, using the ones and tens digits for numbers from 0 to 99 to specify the x- and y-positions, respectively. A loop is used to call the `bar` function once for each category; each time it is called, the plotted bars are assigned a different color.

---



The last steps that we need to do involve aesthetic cleaning to polish it up for interpretability. We can take away the plot border and ticks, since they're arbitrary, but we should change the limits so that the boxes are square. We should also add a legend so that the mapping from colors to category levels is clear.

```
waffle_counts = percentage_blocks(df, 'cat_var')

prev_count = 0
for each category,
for cat in range(waffle_counts.shape[0]):
 # get the block indices
 blocks = np.arange(prev_count, prev_count + waffle_counts[cat])
 # and put a block at each index's location
 x = blocks % 10 # use mod operation to get ones digit
 y = blocks // 10 # use floor division to get tens digit
 plt.bar(x = x, height = 0.8, width = 0.8, bottom = y)
 prev_count += waffle_counts[cat]

aesthetic wrangling
plt.legend(waffle_counts.index, bbox_to_anchor = (1, 0.5), loc = 6)
plt.axis('off')
plt.axis('square')
```

The two calls to Matplotlib's `axis` function make use of two convenience strings for arguments: 'off' removes the axis lines, ticks, and labels, while 'square' ensures that the scaling on each axis is equal within a square bounding box. As for the `legend` call, the first argument is a list of categories as obtained from the sorted `waffle_counts` Series variable. This will match each category to each `bar` call, in order. The "bbox\_to\_anchor" argument sets an anchor for the legend to the right side of the plot, and "loc = 6" positions the anchor to the center left of the legend. The final plot is as it looks at the top of the page:



Other variants of the waffle plot exist to extend it beyond just displaying probabilities. By associating each square with an amount rather than a percentage, we can use waffle plots to show absolute frequencies instead. This might cause us to end up with something other than 100 squares.

```

each box represents five full counts
waffle_counts = (df['cat_var'].value_counts() / 5).astype(int)

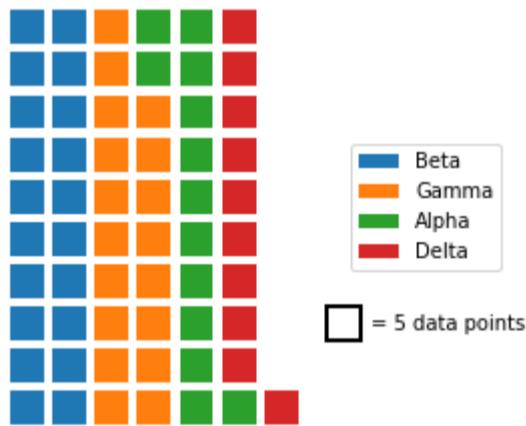
prev_count = 0
for each category,
for cat in range(waffle_counts.shape[0]):
 # get the block indices
 blocks = np.arange(prev_count, prev_count + waffle_counts[cat])
 # and put a block at each index's location
 x = blocks % 10
 y = blocks // 10
 plt.bar(y, 0.8, 0.8, x)
 prev_count += waffle_counts[cat]

box size legend
plt.bar(7.5, 0.8, 0.8, 2, color = 'white', edgecolor = 'black', lw = 2)
plt.text(8.1, 2.4,'= 5 data points', va = 'center')

aesthetic wrangling
plt.legend(waffle_counts.index, bbox_to_anchor = (0.8, 0.5), loc = 6)
plt.axis('off')
plt.axis('square')

```

In the above code, `waffle_counts` has been adjusted so that each box represents 5 data points. Most of the code is the same as before, though it should be noted that the `x` and `y` variables have been swapped in the `bar` function so that the boxes are plotted in columns from left to right. Additional `bar` and `text` calls have been added to the plot to act as an ad hoc legend. The positions of these elements, and the legend, have been adjusted manually through some trial and error to improve the aesthetic appeal. Note that this constitutes more of an explanatory polishing than it is a part of exploration!



As a further extension, there's no restriction against us using icons for each tally, rather than just squares. Infographics often take this approach, by having each icon represent some number of units (e.g. one person icon representing one million people). But while it can be tempting to use icons to represent values as a bit of visual flair, an icon-based plot contains more chart junk than a bar chart that conveys the same information. There's a larger cognitive challenge in having to count a number of icons to understand the scale of a value, compared to just referencing a box's endpoint on a labeled axis.

One other downside of the waffle plot is that it is not commonly supported out of the box for most visualization libraries, including Matplotlib and Seaborn. The length of the demonstration code presented above is a testament to that. The effort required to create a meaningful and useful waffle plot means that it is best employed carefully as a part of explanatory visualizations. During the exploratory phase, you're better off using more traditional plots like the bar chart to more rapidly build your understanding of the data.

## Additional Resources

You don't actually need to go through all of the code wrangling shown above to create waffle plots in Python. The [PyWaffle](#) package can be used with Matplotlib's `figure` function to create waffle plots, with a few options for the orientation and order of icons, but you'll need to install it separately since it's not a major package. One of the main reasons why I didn't use it above is that the syntax for using it is very different from what you've seen and will see in this course. If you want to make use of the library, check out the examples on the linked GitHub page.

---