

## Module 7: 1D data

Let's 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 [3]:

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

In [5]:

```
cars = data.cars()
cars.head()
```

Out[5]:

	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

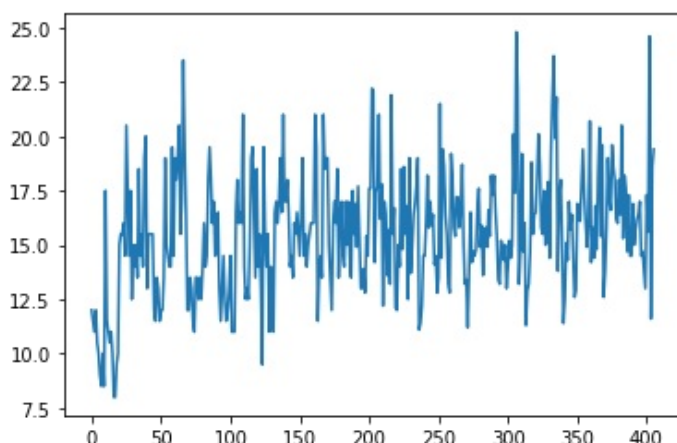
Let's consider the `Acceleration` column as our 1D data. If we ask pandas to plot this series, it'll produce a line graph where the index becomes the horizontal axis.

In [ ]:

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

Out[ ]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f4c58d23cd0>



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

definitely not the plot we want!

It's actually not trivial to use pandas to create an 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 [ ]:

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

Out[ ]:

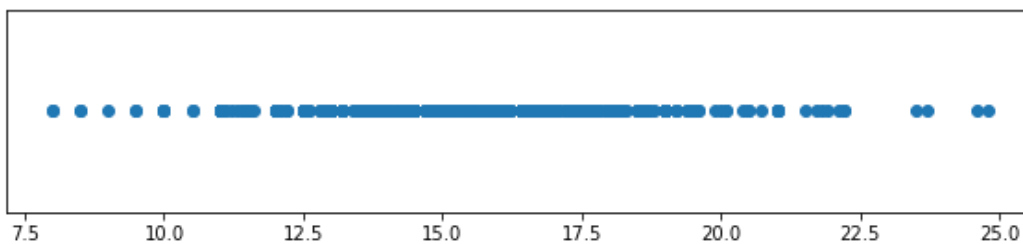
```
array([0, 0, 0])
```

In [32]:

**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 [42]:

```
# TODO: put your code here
cars_np = np.zeros_like([cars.Acceleration])
plt.figure(figsize=(10, 2))
plt.scatter(cars.Acceleration, cars_np)
plt.gca().axes.yaxis.set_visible(False)
```

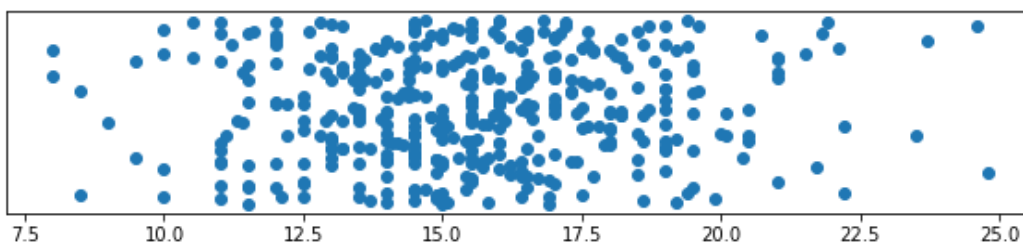


As you can see, there are lots of occlusions. So this plot cannot show the distribution properly and we would like to fix it. How about adding some jitters? You can use `numpy`'s `random.rand()` function to generate random numbers, instead of using an array with zeros.

**Q: create a jittered 1D scatter plot.**

In [54]:

```
# TODO: put your code here
# jittered_y = ...
# plt ...
jittered_y = np.random.rand(cars.Acceleration.size)
plt.figure(figsize=(10, 2))
plt.scatter(cars.Acceleration, jittered_y)
plt.gca().axes.yaxis.set_visible(False)
```

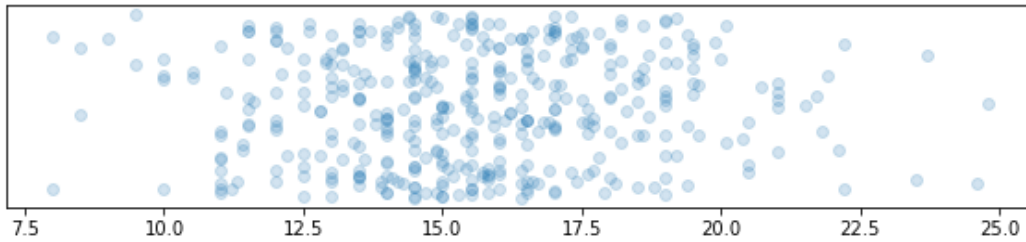


We can further improve this by adding transparency to the symbols. The transparency option for `scatter` function is called `alpha`. Set it to be 0.2.

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

In [61]:

```
# TODO: put your code here
jittered_y = np.random.rand(cars.Acceleration.size)
plt.figure(figsize=(10, 2))
plt.scatter(cars.Acceleration, jittered_y, alpha=0.2)
plt.gca().axes.yaxis.set_visible(False)
```

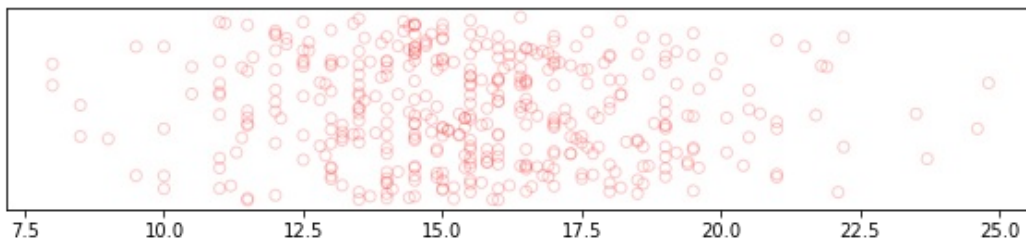


Another strategy is using empty symbols. The option is `facecolors`. You can also change the stroke color (`edgecolors`).

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

In [6]:

```
# TODO: put your code here
jittered_y = np.random.rand(cars.Acceleration.size)
plt.figure(figsize=(10, 2))
plt.scatter(cars.Acceleration, jittered_y, alpha=0.2, facecolor = 'none', edgecolors='r'
)
plt.gca().axes.yaxis.set_visible(False)
```



In [ ]:

```
# TODO: put your code here
jittered_y = np.random.rand(cars.Acceleration.size)
plt.figure(figsize=(10, 2))
plt.scatter(cars.Acceleration, jittered_y, alpha=0.2, facecolor = 'none', edgecolors='r'
)
plt.gca().axes.yaxis.set_visible(False)
```

## What happens if you have lots and lots of points?

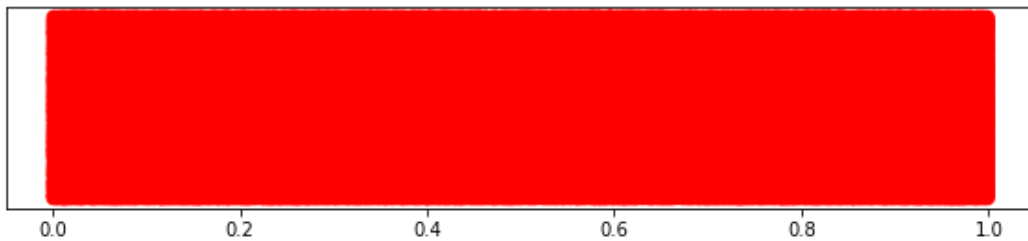
Whatever strategy that you use, it's almost useless if you have too many data points. Let's play with different number of data points and see how it looks.

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

In [22]:

```
# TODO: play with N and see what happens.
x = np.random.rand(900000)
y = np.random.rand(900000)

# TODO: 1D scatter plot code here
plt.figure(figsize=(10, 2))
plt.scatter(x, y, alpha=0.2, facecolor = 'none', edgecolors='r')
plt.gca().axes.yaxis.set_visible(False)
```



## Histogram and boxplot

When you have lots of data points, you can't no longer 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. Histogram is about aggregating and counting the data while boxplot is about summarizing the data. Let's first draw some histograms.

### Histogram

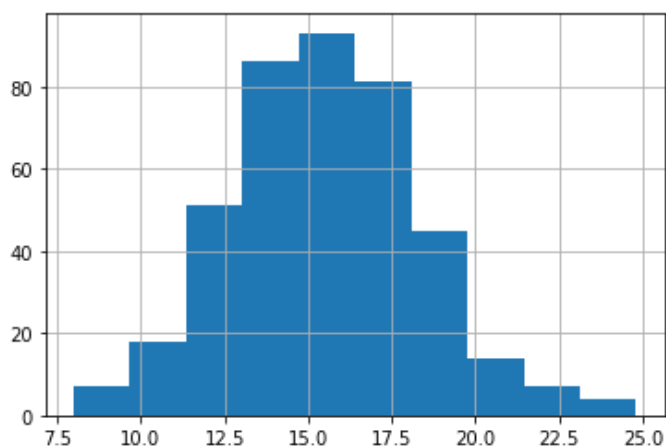
It's very easy to draw a histogram with pandas.

In [23]:

```
cars.Acceleration.hist()
```

Out[23]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f083b1a10d0>



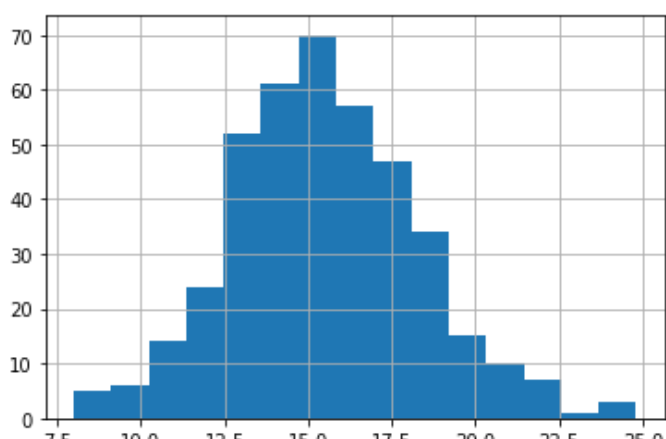
You can adjust the bin size, which is the main parameter of the histogram.

In [24]:

```
cars.Acceleration.hist(bins=15)
```

Out[24]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f083b199550>



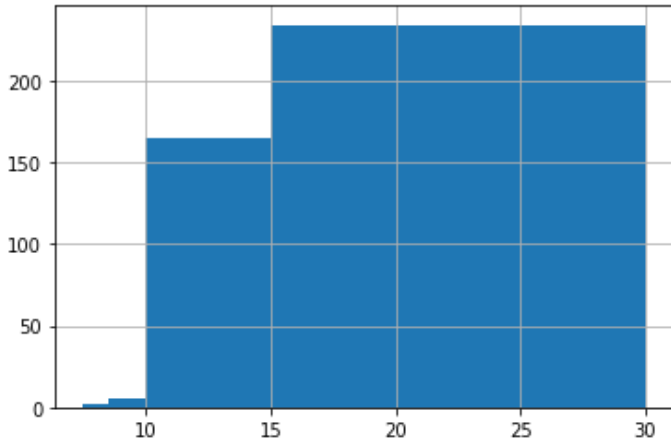
**You can even specify the actual bins.**

In [25]:

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

Out[25]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f083a970ed0>



**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.**

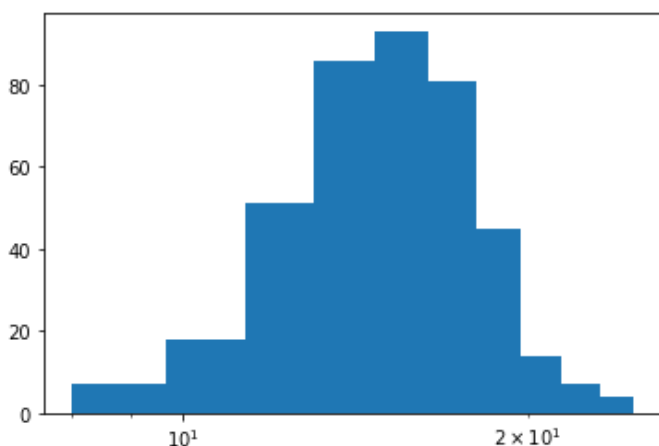
(hints: do you remember what we discussed regarding histogram? Also [pandas documentation](#) does not show the option that you should use. You should take a look at the `matplotlib`'s documentation.

In [29]:

```
# TODO: put your code here
counts, bins = np.histogram(cars.Acceleration)
plt.gca().set_xscale("log")
plt.hist(bins[:-1], bins, weights=counts)
```

Out[29]:

```
(array([ 7., 18., 51., 86., 93., 81., 45., 14.,  7.,  4.]),
 array([ 8. ,  9.68, 11.36, 13.04, 14.72, 16.4 , 18.08, 19.76, 21.44,
        23.12, 24.8 ]),
 <a list of 10 Patch objects>)
```



## 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

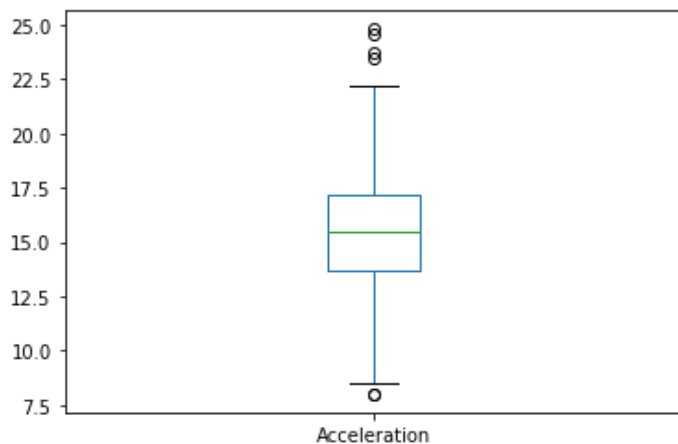
In [36]:

```
cars_np = pd.DataFrame(cars.Acceleration)

cars_np.boxplot(grid = False)
```

Out[36]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f0839a6c8d0>



## 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. In other words, they are very cool. 😊

I'd like to 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 [37]:

```
import seaborn as sns
import altair as alt

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

In [38]:

```
cars.head()
```

Out[38]:

	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

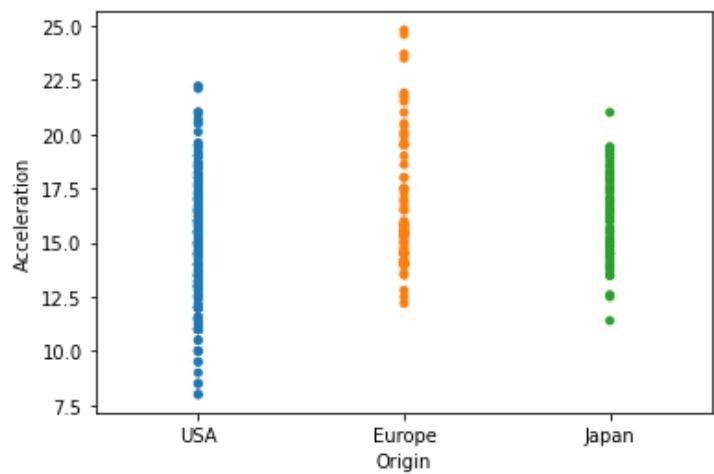
Seaborn has a built-in function to create 1D scatter plots with multiple categories, and it adds jittering by default.

In [39]:

```
sns.stripplot(x='Origin', y='Acceleration', data=cars, jitter=False)
```

Out[39]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f082a6170d0>



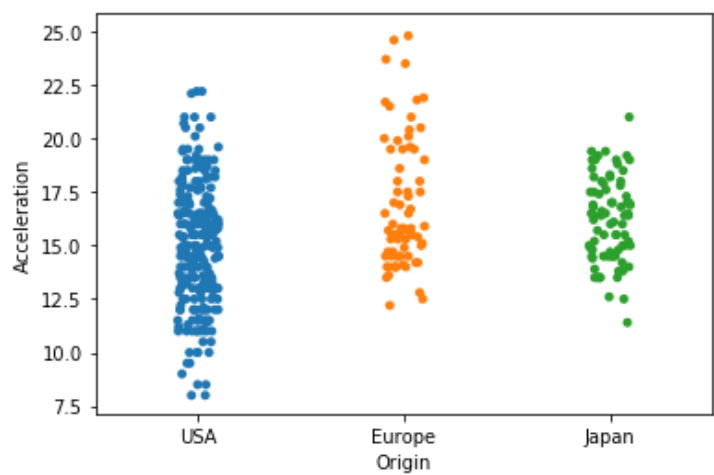
And you can easily add jitters or even create a beeswarm plot.

In [ ]:

```
sns.stripplot(x='Origin', y='Acceleration', data=cars)
```

Out[ ]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x1a1ada8050>



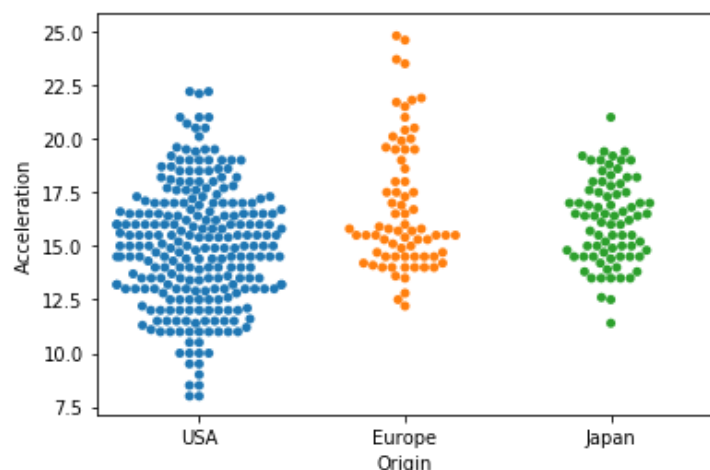
Seems like European cars tend to have good acceleration. 😊 Let's look at the beeswarm plot, which is a pretty nice option for fairly small datasets.

In [ ]:

```
sns.swarmplot(x='Origin', y='Acceleration', data=cars)
```

Out[ ]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x1a1a9ad0d0>



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` ?**

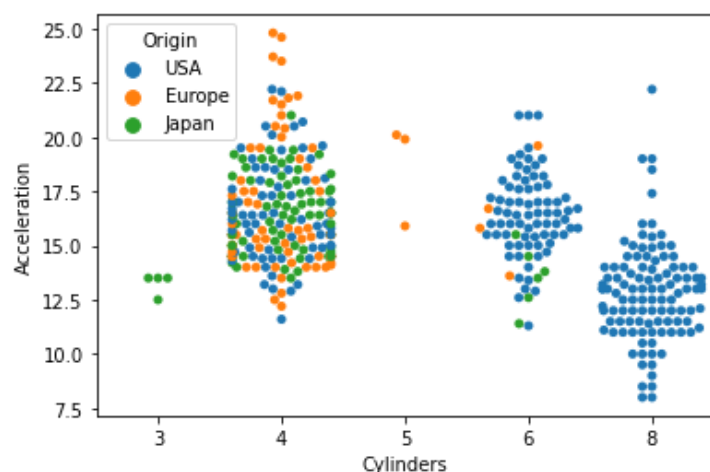
In [43]:

```
# TODO: put your code here
sns.swarmplot(x='Cylinders', y='Acceleration', data=cars, hue='Origin')
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/categorical.py:1296: UserWarning: 25.1% of
the points cannot be placed; you may want to decrease the size of the markers or use stri
pplot.
  warnings.warn(msg, UserWarning)
/usr/local/lib/python3.7/dist-packages/seaborn/categorical.py:1296: UserWarning: 6.5% of
the points cannot be placed; you may want to decrease the size of the markers or use stri
pplot.
  warnings.warn(msg, UserWarning)
```

Out[43]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f0821bcc310>



And of course you can create box plots too.



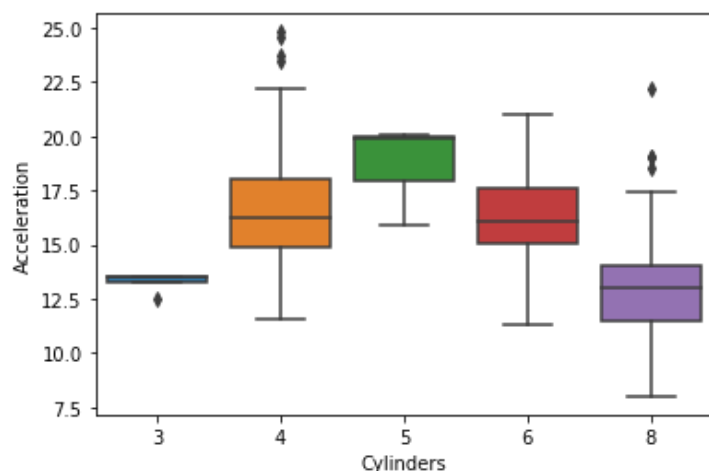
**Q: Create boxplots to show the relationships between Cylinders and Acceleration .**

In [45]:

```
# TODO: put your code here
sns.boxplot(x='Cylinders', y='Acceleration', data=cars)
```

Out[45]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f0821afde10>



## 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 [56]:

```
alt.Chart(cars)
```

```
-----
SchemaValidationError                                Traceback (most recent call last)
/usr/local/lib/python3.7/dist-packages/altair/vegalite/v4/api.py in to_dict(self, *args,
**kwargs)
    380         if dct is None:
    381             kwargs["validate"] = "deep"
--> 382         dct = super(TopLevelMixin, copy).to_dict(*args, **kwargs)
    383
    384         # TODO: following entries are added after validation. Should they be vali
dated?

/usr/local/lib/python3.7/dist-packages/altair/utils/schemapi.py in to_dict(self, validate
, ignore, context)
    337         self.validate(result)
    338         except jsonschema.ValidationError as err:
--> 339             raise SchemaValidationError(self, err)
    340         return result
    341
```

SchemaValidationError: Invalid specification

```
altair.vegalite.v4.api.Chart, validating 'required'
'mark' is a required property
```

Out[56]:

```
alt.Chart(...)
```

In [55]:

```
alt.Chart(cars).mark_point().encode(
```

```
alt.Chart(cars).mark_point().encode(
    x='Acceleration',
)
alt.figure
```

Out[55]:

**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 [57]:

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

Out[57]:

```
RendererRegistry.enable('default')
```

In [53]:

```
alt.Chart(cars).mark_point()
```

Out[53]:

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 [61]:

```
alt.Chart(cars).mark_point().encode(
    x='Acceleration',
).properties(width=400, height=300)
```

Out[61]:

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 [62]:

```
alt.Chart(cars).mark_point().encode(
    x=alt.X('Acceleration')
).properties(width=400, height=300)
```

Out[62]:

There is another nice mark called `tick`:

In [63]:

```
alt.Chart(cars).mark_tick().encode(
    x='Acceleration',
).properties(width=400, height=300)
```

Out[63]:

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 [66]:

```
alt.Chart(cars).mark_bar().encode(
```

```
x=alt.X('Acceleration', bin=True),
y='count()'
).properties(width=400, height=300)
```

Out[66]:

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

In [84]:

```
# TODO: put your code here

alt.Chart(cars).mark_point().encode(
    x='Acceleration',
    y='Horsepower',
    color='Origin'
)
chart
```

Out[84]:

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](https://altair-viz.github.io/getting_started/starting.html#publishing-your-visualization)

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

In [86]:

```
# TODO: your code here.
chart = alt.Chart(cars).mark_point().encode(
    x='Acceleration',
    y='Horsepower',
    color='Origin'
)
chart
chart.save('m07.html')
```