

Visualizing the distribution of a dataset

- ▶ When dealing with a set of data, often the first thing you'll want to do is get a sense for how the variables are distributed.
- ▶ This chapter of the tutorial will give a brief introduction to some of the tools in seaborn for examining univariate and bivariate distributions.
- ▶ You may also want to look at the categorical plots chapter for examples of functions that make it easy to compare the distribution of a variable across levels of other variables.

```
%matplotlib inline
import numpy as np
import pandas as pd
from scipy import stats, integrate
import matplotlib.pyplot as plt
import seaborn as sns
sns.set(color_codes=True)
np.random.seed(sum(map(ord,
                        "distributions"))))
```

The most convenient way to take a quick look at a univariate distribution in seaborn is the `distplot()` function. By default, this will draw a histogram and fit a kernel density estimate (KDE).

```
x = np.random.normal(size=100)
sns.distplot(x);
```

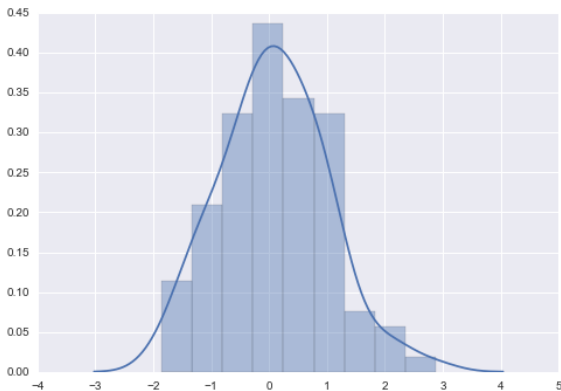


Figure:

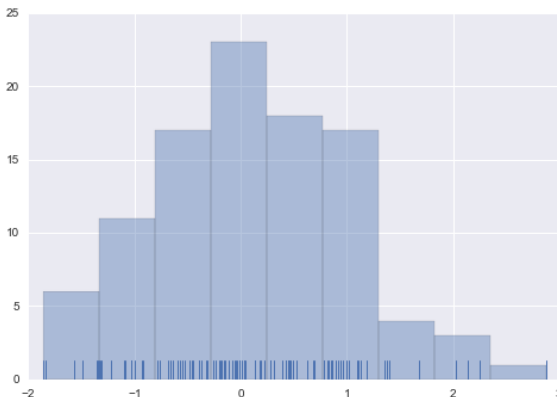
Histograms

- ▶ Histograms are likely familiar, and a `hist` function already exists in `matplotlib`.
- ▶ A histogram represents the distribution of data by forming bins along the range of the data and then drawing bars to show the number of observations that fall in each bin.

Seaborn Workshop

- ▶ To illustrate this, let's remove the density curve and add a rug plot, which draws a small vertical tick at each observation.
- ▶ You can make the rug plot itself with the `rugplot()` function, but it is also available in `distplot()`:

```
sns.distplot(x, kde=False, rug=True);
```

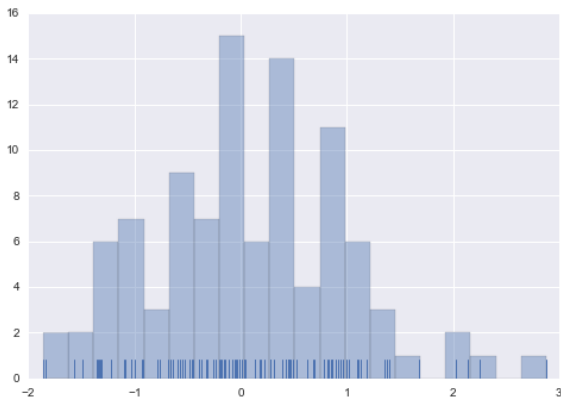


Seaborn Workshop

- ▶ When drawing histograms, the main choice you have is the number of bins to use and where to place them.
- ▶ `distplot()` uses a simple rule to make a good guess for what the right number is by default, but trying more or fewer bins might reveal other features in the data:

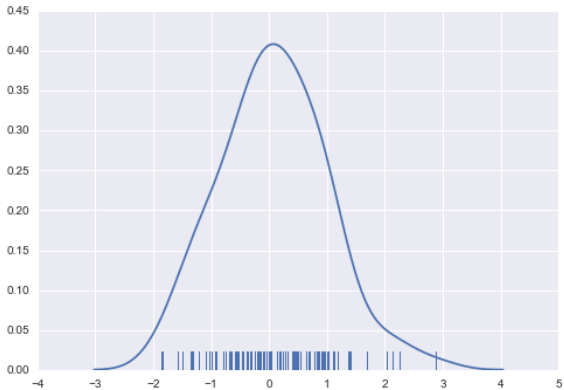
Seaborn Workshop

```
sns.distplot(x, bins=20, kde=False, rug=True);
```



- ▶ The kernel density estimate may be less familiar, but it can be a useful tool for plotting the shape of a distribution.
- ▶ Like the histogram, the KDE plots encodes the density of observations on one axis with height along the other axis:


```
sns.distplot(x, hist=False, rug=True);
```



- ▶ Drawing a KDE is more computationally involved than drawing a histogram.
- ▶ What happens is that each observation is first replaced with a normal (Gaussian) curve centered at that value:

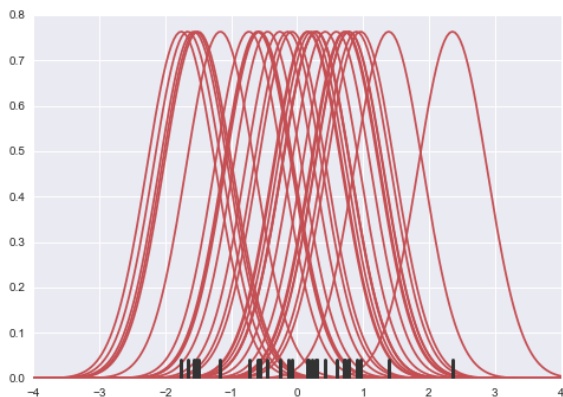
```
x = np.random.normal(0, 1, size=30)
bandwidth = 1.06 * x.std() * x.size ** (-1 / 5.)
support = np.linspace(-4, 4, 200)

kernels = []
for x_i in x:

    kernel = stats.norm(x_i, bandwidth).pdf(support)
    kernels.append(kernel)
plt.plot(support, kernel, color="r")
```

Seaborn Workshop

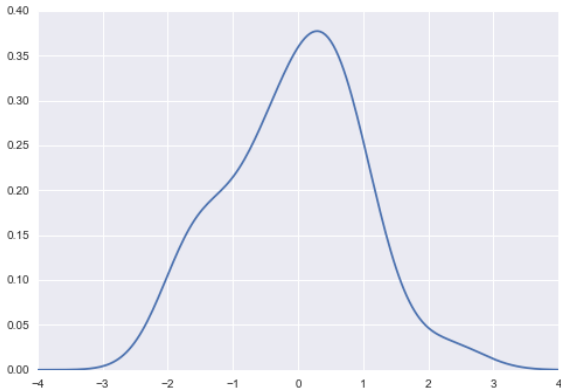
```
sns.rugplot(x, color=".2", linewidth=3);
```



Next, these curves are summed to compute the value of the density at each point in the support grid. The resulting curve is then normalized so that the area under it is equal to 1:

Seaborn Workshop

```
density = np.sum(kernels, axis=0)  
density /= integrate.trapz(density, support)  
plt.plot(support, density);
```

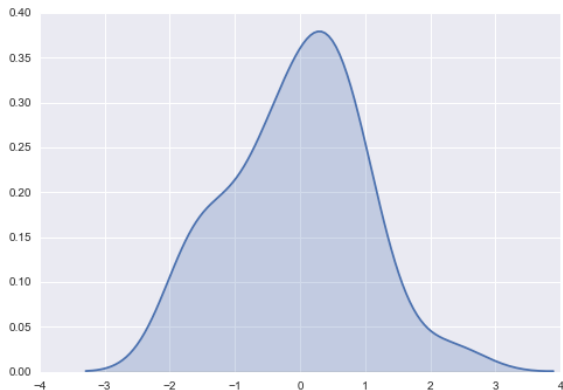


Seaborn Workshop

- ▶ We can see that if we use the `kdeplot()` function in `seaborn`, we get the same curve.
- ▶ This function is used by `distplot()`, but it provides a more direct interface with easier access to other options when you just want the density estimate:

Seaborn Workshop

```
sns.kdeplot(x, shade=True);
```

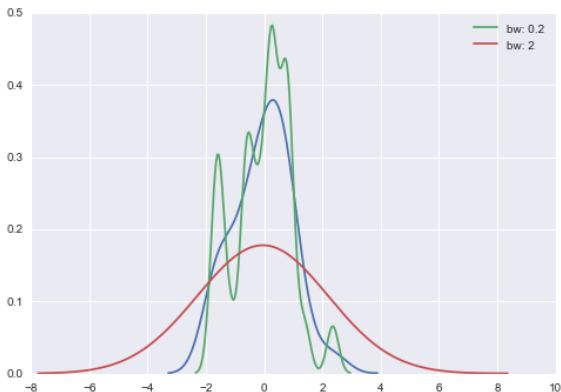


Seaborn Workshop

- ▶ The bandwidth (bw) parameter of the KDE controls how tightly the estimation is fit to the data, much like the bin size in a histogram.
- ▶ It corresponds to the width of the kernels we plotted above.
- ▶ The default behavior tries to guess a good value using a common reference rule, but it may be helpful to try larger or smaller values:

Seaborn Workshop

```
{sns.kdeplot(x)  
sns.kdeplot(x, bw=.2, label="bw: 0.2")  
sns.kdeplot(x, bw=2, label="bw: 2")  
plt.legend();...}
```

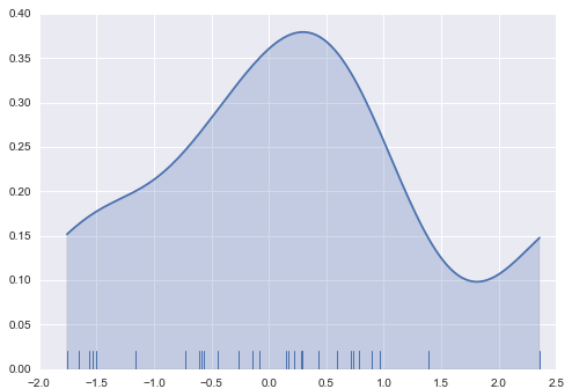


Seaborn Workshop

- ▶ As you can see above, the nature of the Gaussian KDE process means that estimation extends past the largest and smallest values in the dataset.
- ▶ Its possible to control how far past the extreme values the curve is drawn with the cut parameter; however, this only influences how the curve is drawn and not how it is fit:

Seaborn Workshop

```
sns.kdeplot(x, shade=True, cut=0)  
sns.rugplot(x);
```

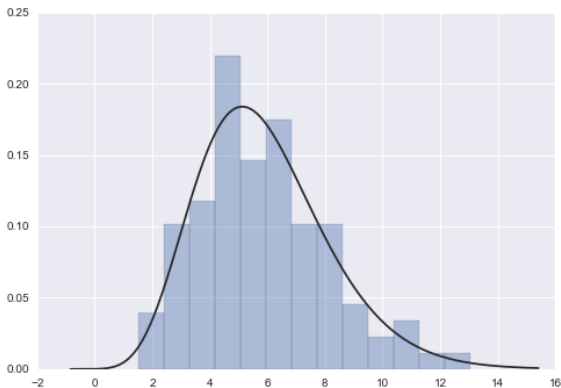


Fitting parametric distributions

You can also use `distplot()` to fit a parametric distribution to a dataset and visually evaluate how closely it corresponds to the observed data:

Seaborn Workshop

```
x = np.random.gamma(6, size=200)
sns.distplot(x, kde=False, fit=stats.gamma);
```



Plotting Bivariate Bistributions

- ▶ It can also be useful to visualize a bivariate distribution of two variables.
- ▶ The easiest way to do this in seaborn is to just the `jointplot()` function, which creates a multi-panel figure that shows both the bivariate (or joint) relationship between two variables along with the univariate (or marginal) distribution of each on separate axes.

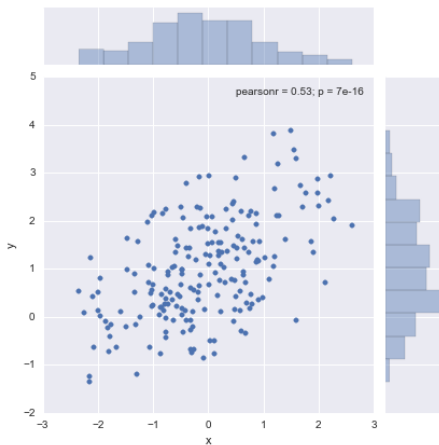
Seaborn Workshop

```
mean, cov = [0, 1], [(1, .5), (.5, 1)]  
  
data = np.random.multivariate_normal(mean,  
    cov, 200)  
  
df = pd.DataFrame(data, columns=["x", "y"])
```

- ▶ The most familiar way to visualize a bivariate distribution is a scatterplot, where each observation is shown with point at the x and y values.
- ▶ This is analogous to a rug plot on two dimensions.
- ▶ You can draw a scatterplot with the matplotlib `plt.scatter` function, and it is also the default kind of plot shown by the `jointplot()` function:

Seaborn Workshop

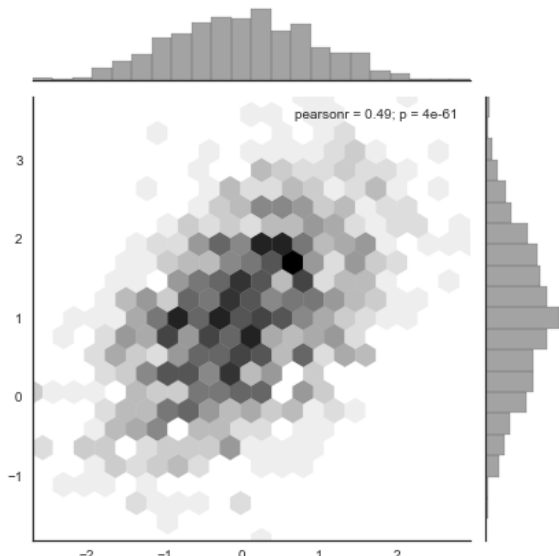
```
sns.jointplot(x="x", y="y", data=df);
```



Hexbin plots

- ▶ The bivariate analogue of a histogram is known as a hexbin plot, because it shows the counts of observations that fall within hexagonal bins.
- ▶ This plot works best with relatively large datasets.
- ▶ Its available through the matplotlib `plt.hexbin` function and as a style in `jointplot()`.
- ▶ It looks best with a white background:

```
x, y = np.random.multivariate_normal(mean, cov, 1000).T  
with sns.axes_style("white"):  
sns.jointplot(x=x, y=y, kind="hex", color="k");
```

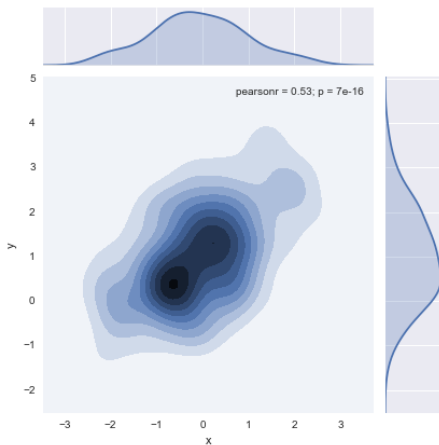


Kernel density estimation

- ▶ It is also possible to use the kernel density estimation procedure described above to visualize a bivariate distribution.
- ▶ In seaborn, this kind of plot is shown with a contour plot and is available as a style in `jointplot()`:

Seaborn Workshop

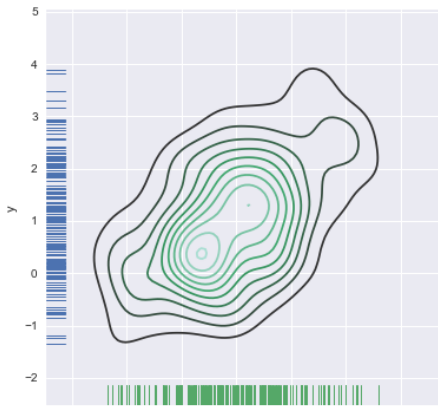
```
sns.jointplot(x="x", y="y", data=df, kind="kde");
```



Seaborn Workshop

- ▶ You can also draw a two-dimensional kernel density plot with the `kdeplot()` function.
- ▶ This allows you to draw this kind of plot onto a specific (and possibly already existing) matplotlib axes, whereas the `jointplot()` function manages its own figure:

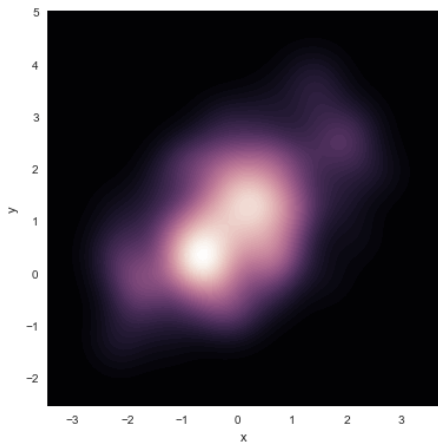
```
f, ax = plt.subplots(figsize=(6, 6))  
sns.kdeplot(df.x, df.y, ax=ax)  
sns.rugplot(df.x, color="g", ax=ax)  
sns.rugplot(df.y, vertical=True, ax=ax);
```



If you wish to show the bivariate density more continuously, you can simply increase the number of contour levels:

```
f, ax = plt.subplots(figsize=(6, 6))  
cmap = sns.cubehelix_palette(as_cmap=True, dark=0, light=1)  
sns.kdeplot(df.x, df.y, cmap=cmap, n_levels=60, shade=True)
```


Hexbin Plots

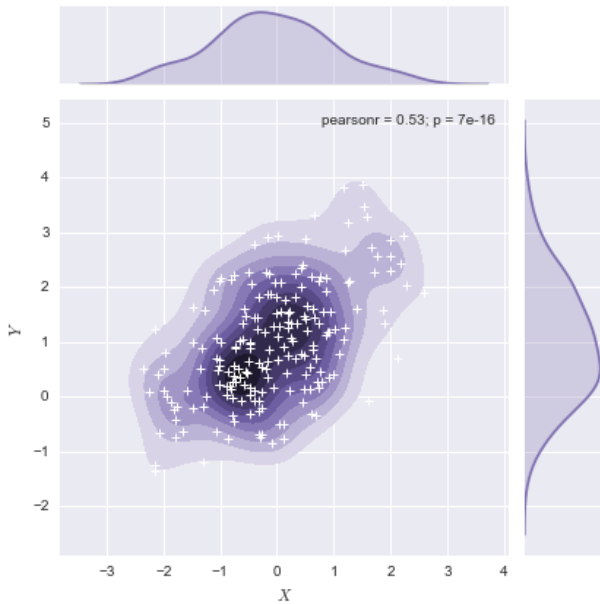


Seaborn Workshop

- ▶ The `jointplot()` function uses a `JointGrid` to manage the figure.
- ▶ For more flexibility, you may want to draw your figure by using `JointGrid` directly.
- ▶ `jointplot()` returns the `JointGrid` object after plotting, which you can use to add more layers or to tweak other aspects of the visualization:

Seaborn Workshop

```
g = sns.jointplot(x="x", y="y",  
data=df, kind="kde", color="m")  
  
g.plot_joint(plt.scatter, c="w",  
s=30, linewidth=1, marker="+")  
  
g.ax_joint.collections[0].set_alpha(0)  
  
g.set_axis_labels("$X$", "$Y$");
```



- ▶ To plot multiple pairwise bivariate distributions in a dataset, you can use the `pairplot()` function.
- ▶ This creates a matrix of axes and shows the relationship for each pair of columns in a `DataFrame`.
- ▶ By default, it also draws the univariate distribution of each variable on the diagonal Axes:

```
iris = sns.load_dataset("iris")  
sns.pairplot(iris);
```

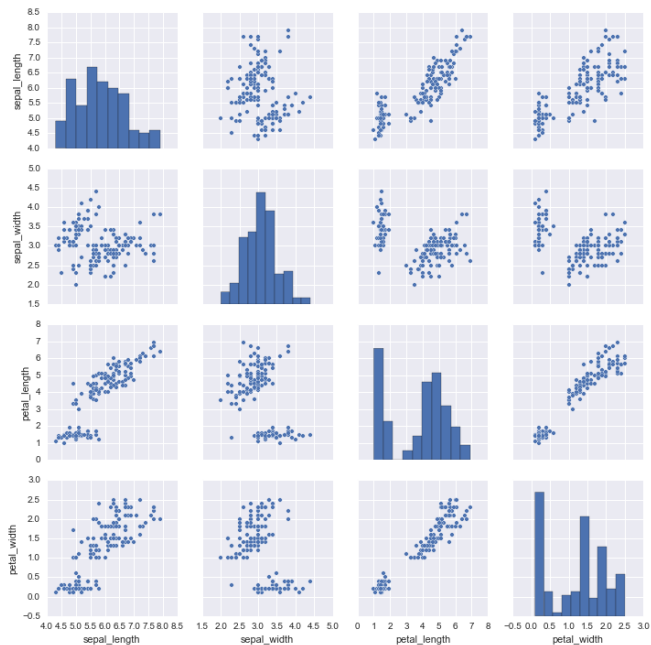


Figure:

Seaborn Workshop

Much like the relationship between `jointplot()` and `JointGrid`, the `pairplot()` function is built on top of a `PairGrid` object, which can be used directly for more flexibility:

```
g = sns.PairGrid(iris)
g.map_diag(sns.kdeplot)
g.map_offdiag(sns.kdeplot, cmap="Blues_d",
n_levels=6);
```

