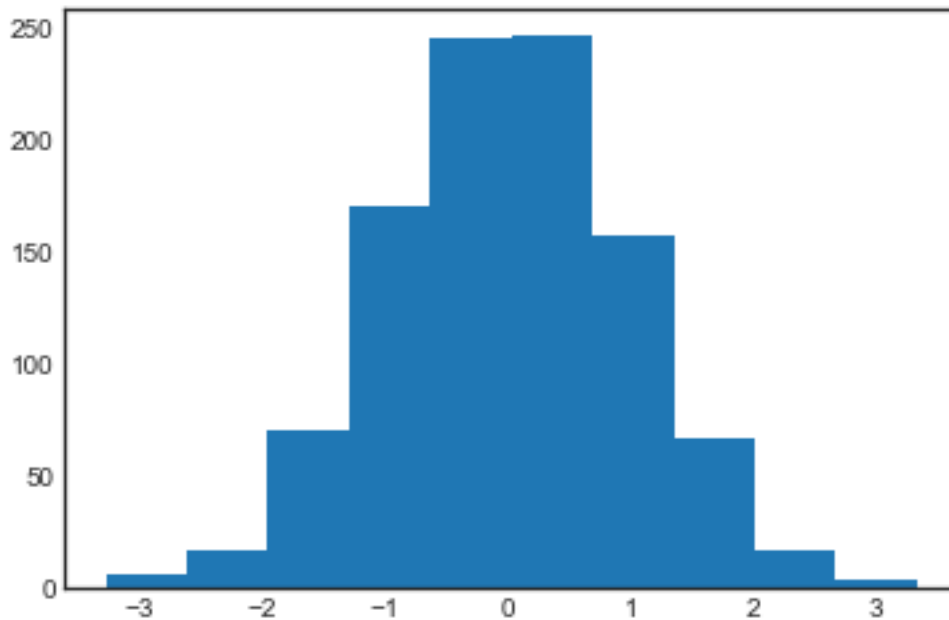# 04.05-Histograms-and-Binnings

November 10, 2024

## 1  Histograms, Binnings, and Density

A simple histogram can be a great first step in understanding a dataset. Earlier, we saw a preview of Matplotlib's histogram function (discussed in Comparisons, Masks, and Boolean Logic), which creates a basic histogram in one line, once the normal boilerplate imports are done (see the following figure):

```
[1]: %matplotlib inline
     import numpy as np
     import matplotlib.pyplot as plt
     plt.style.use('seaborn-white')

     rng = np.random.default_rng(1701)
     data = rng.normal(size=1000)
```
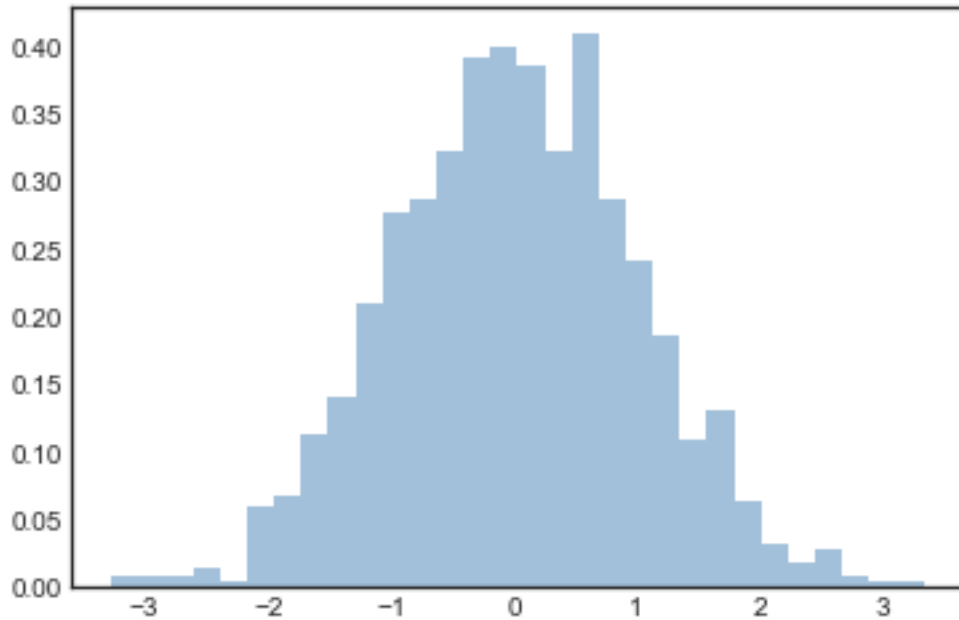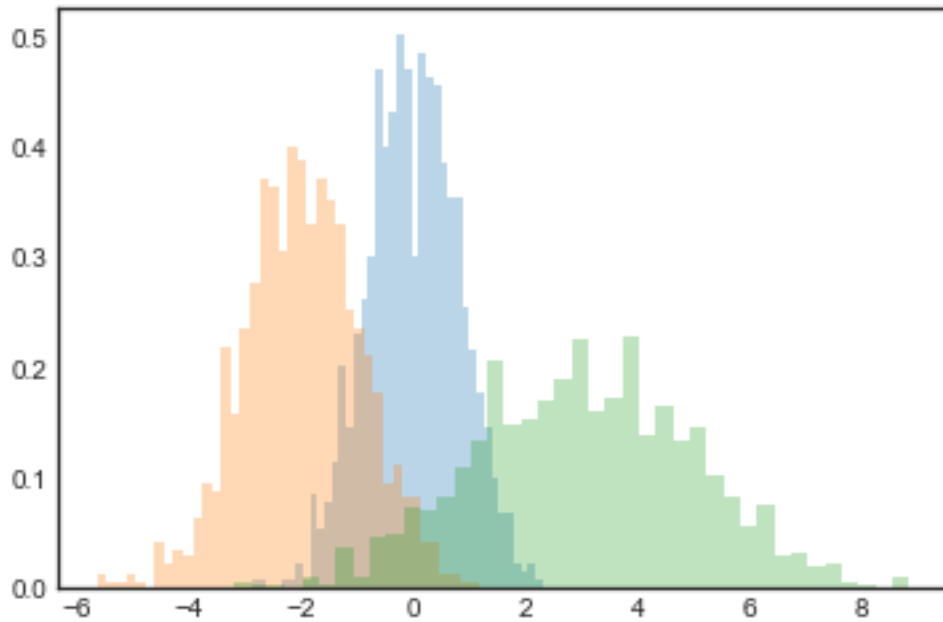
```
[2]: plt.hist(data);
```

The `hist` function has many options to tune both the calculation and the display; here's an example of a more customized histogram, shown in the following figure:

```
[3]: plt.hist(data, bins=30, density=True, alpha=0.5,
              histtype='stepfilled', color='steelblue',
              edgecolor='none');
```



The `plt.hist` docstring has more information on other available customization options. I find this combination of `histtype='stepfilled'` along with some transparency `alpha` to be helpful when comparing histograms of several distributions (see the following figure):

```
[4]: x1 = rng.normal(0, 0.8, 1000)
     x2 = rng.normal(-2, 1, 1000)
     x3 = rng.normal(3, 2, 1000)

     kwargs = dict(histtype='stepfilled', alpha=0.3, density=True, bins=40)

     plt.hist(x1, **kwargs)
     plt.hist(x2, **kwargs)
     plt.hist(x3, **kwargs);
```

If you are interested in computing, but not displaying, the histogram (that is, counting the number of points in a given bin), you can use the `np.histogram` function:

```
[5]: counts, bin_edges = np.histogram(data, bins=5)
     print(counts)
```

```
[ 23 241 491 224  21]
```

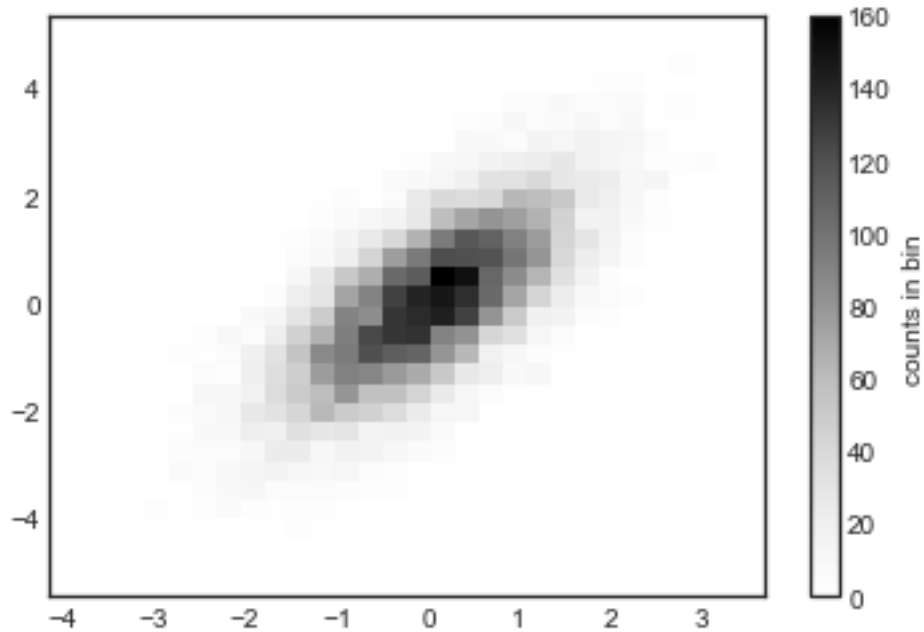## 1.1 Two-Dimensional Histograms and Binnings

Just as we create histograms in one dimension by dividing the number line into bins, we can also create histograms in two dimensions by dividing points among two-dimensional bins. We'll take a brief look at several ways to do this here. We'll start by defining some data—an `x` and `y` array drawn from a multivariate Gaussian distribution:

```
[6]: mean = [0, 0]
     cov = [[1, 1], [1, 2]]
     x, y = rng.multivariate_normal(mean, cov, 10000).T
```

### 1.1.1 plt.hist2d: Two-dimensional histogram

One straightforward way to plot a two-dimensional histogram is to use Matplotlib's `plt.hist2d` function (see the following figure):

```
[7]: plt.hist2d(x, y, bins=30)
     cb = plt.colorbar()
     cb.set_label('counts in bin')
```

Just like `plt.hist`, `plt.hist2d` has a number of extra options to fine-tune the plot and the binning, which are nicely outlined in the function docstring. Further, just as `plt.hist` has a counterpart in `np.histogram`, `plt.hist2d` has a counterpart in `np.histogram2d`:

```
[8]: counts, xedges, yedges = np.histogram2d(x, y, bins=30)
     print(counts.shape)
```
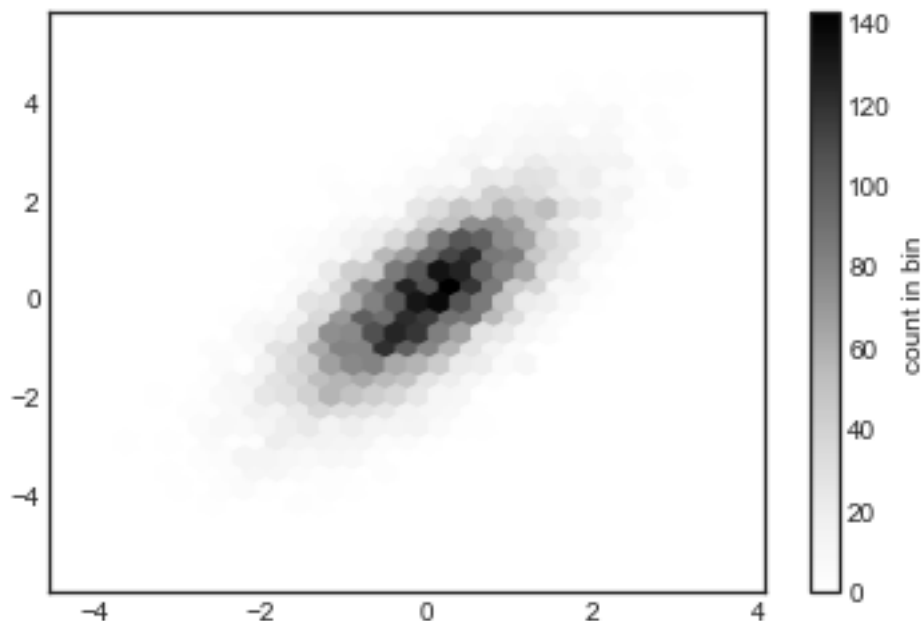
```
(30, 30)
```

For the generalization of this histogram binning when there are more than two dimensions, see the `np.histogramdd` function.

### 1.1.2  plt.hexbin: Hexagonal binnings

The two-dimensional histogram creates a tesselation of squares across the axes. Another natural shape for such a tesselation is the regular hexagon. For this purpose, Matplotlib provides the `plt.hexbin` routine, which represents a two-dimensional dataset binned within a grid of hexagons (see the following figure):

```
[9]: plt.hexbin(x, y, gridsize=30)
     cb = plt.colorbar(label='count in bin')
```

`plt.hexbin` has a number of additional options, including the ability to specify weights for each point and to change the output in each bin to any NumPy aggregate (mean of weights, standard deviation of weights, etc.).

### 1.1.3 Kernel density estimation

Another common method for estimating and representing densities in multiple dimensions is *kernel density estimation* (KDE). This will be discussed more fully in In-Depth: Kernel Density Estimation, but for now I'll simply mention that KDE can be thought of as a way to "smear out" the points in space and add up the result to obtain a smooth function. One extremely quick and simple KDE implementation exists in the `scipy.stats` package. Here is a quick example of using KDE (see the following figure):

```python
from scipy.stats import gaussian_kde

# fit an array of size [Ndim, Nsamples]
data = np.vstack([x, y])
kde = gaussian_kde(data)

# evaluate on a regular grid
xgrid = np.linspace(-3.5, 3.5, 40)
ygrid = np.linspace(-6, 6, 40)
Xgrid, Ygrid = np.meshgrid(xgrid, ygrid)
Z = kde.evaluate(np.vstack([Xgrid.ravel(), Ygrid.ravel()]))

# Plot the result as an image
plt.imshow(Z.reshape(Xgrid.shape),
```
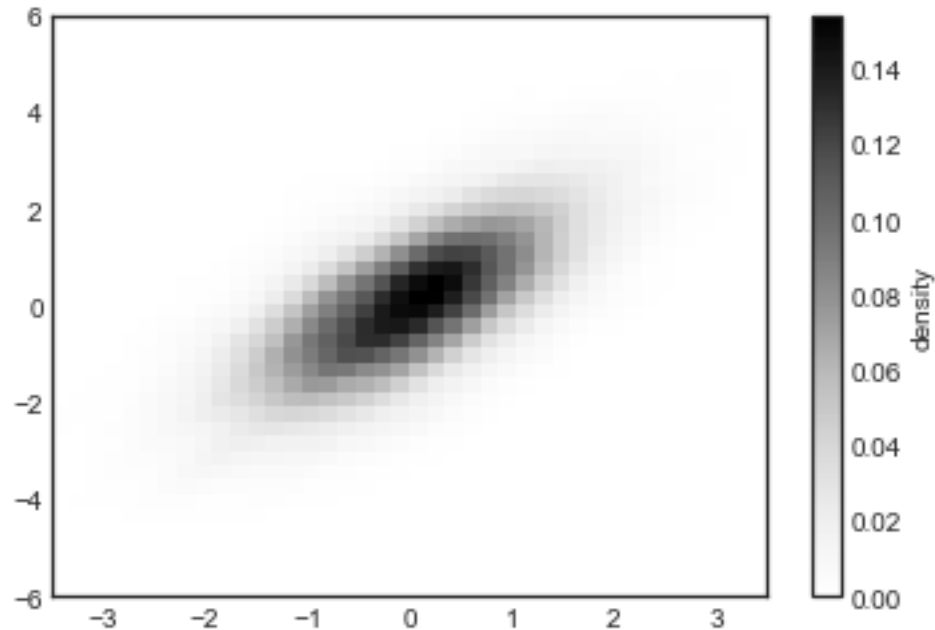
```
            origin='lower', aspect='auto',
            extent=[-3.5, 3.5, -6, 6])
cb = plt.colorbar()
cb.set_label("density")
```



KDE has a smoothing length that effectively slides the knob between detail and smoothness (one example of the ubiquitous bias–variance trade-off). The literature on choosing an appropriate smoothing length is vast; `gaussian_kde` uses a rule of thumb to attempt to find a nearly optimal smoothing length for the input data.

Other KDE implementations are available within the SciPy ecosystem, each with its own strengths and weaknesses; see, for example, `sklearn.neighbors.KernelDensity` and `statsmodels.nonparametric.KDEMultivariate`. For visualizations based on KDE, using Matplotlib tends to be overly verbose. The Seaborn library, discussed in Visualization With Seaborn, provides a much more compact API for creating KDE-based visualizations.