# Sheffield Weather Station

**The following analysis was excerpted from the Data Camp course "Statistical Thinking in Python"**

## Visualizing bootstrap samples

In this exercise, you will generate bootstrap samples from the set of annual rainfall data measured at the Sheffield Weather Station in the UK from 1883 to 2015. The data are stored in the NumPy array `rainfall` in units of millimeters (mm). By graphically displaying the bootstrap samples with an ECDF, you can get a feel for how bootstrap sampling allows probabilistic descriptions of data.

In [6]:

```python
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

In [7]:

```python
rainfall = np.array([ 875.5,  648.2,  788.1,  940.3,  491.1,  743.5,  730.1,  686.5,
        878.8,  865.6,  654.9,  831.5,  798.1,  681.8,  743.8,  689.1,
        752.1,  837.2,  710.6,  749.2,  967.1,  701.2,  619. ,  747.6,
        803.4,  645.6,  804.1,  787.4,  646.8,  997.1,  774. ,  734.5,
        835. ,  840.7,  659.6,  828.3,  909.7,  856.9,  578.3,  904.2,
        883.9,  740.1,  773.9,  741.4,  866.8,  871.1,  712.5,  919.2,
        927.9,  809.4,  633.8,  626.8,  871.3,  774.3,  898.8,  789.6,
        936.3,  765.4,  882.1,  681.1,  661.3,  847.9,  683.9,  985.7,
        771.1,  736.6,  713.2,  774.5,  937.7,  694.5,  598.2,  983.8,
        700.2,  901.3,  733.5,  964.4,  609.3, 1035.2,  718. ,  688.6,
        736.8,  643.3, 1038.5,  969. ,  802.7,  876.6,  944.7,  786.6,
        770.4,  808.6,  761.3,  774.2,  559.3,  674.2,  883.6,  823.9,
        960.4,  877.8,  940.6,  831.8,  906.2,  866.5,  674.1,  998.1,
        789.3,  915. ,  737.1,  763. ,  666.7,  824.5,  913.8,  905.1,
        667.8,  747.4,  784.7,  925.4,  880.2, 1086.9,  764.4, 1050.1,
        595.2,  855.2,  726.9,  785.2,  948.8,  970.6,  896. ,  618.4,
        572.4, 1146.4,  728.2,  864.2,  793. ])
```

In [8]:

```python
def ecdf(data):
    """Compute ECDF for a one-dimensional array of measurements."""
    # Number of data points: n
    n = len(data)

    # x-data for the ECDF: x
    x = np.sort(data)

    # y-data for the ECDF: y
    y = np.arange(1, n+1) / n

    return x, y
```

In [9]:

```python
for _ in range(50):
    # Generate bootstrap sample: bs_sample
```
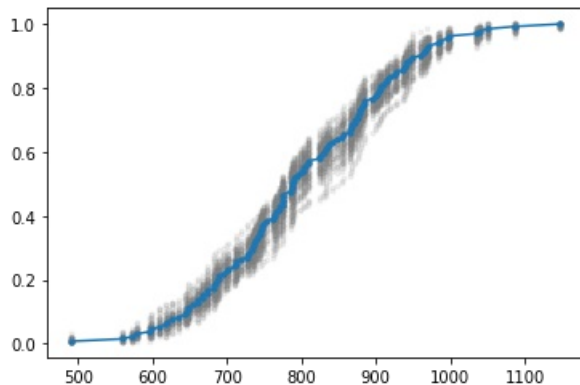
```
    # Generate bootstrap sample: bs_sample
    bs_sample = np.random.choice(rainfall, size=len(rainfall))

    # Compute and plot ECDF from bootstrap sample
    x, y = ecdf(bs_sample)
    _ = plt.plot(x, y, marker='.', linestyle='none',
                 color='gray', alpha=0.1)

# Compute and plot ECDF from original data
x, y = ecdf(rainfall)
_ = plt.plot(x, y, marker='.')
```



Good job! Notice how the bootstrap samples give an idea of how the distribution of rainfalls is spread.

In [10]:

```
def bootstrap_replicate_1d(data, func):
    """Generate bootstrap replicate of 1D data."""
    bs_sample = np.random.choice(data, len(data))
    return func(bs_sample)
```

In [11]:

```
def draw_bs_reps(data, func, size=1):
    """Draw bootstrap replicates."""

    # Initialize array of replicates: bs_replicates
    bs_replicates = np.empty(shape=size)

    # Generate replicates
    for i in range(size):
        bs_replicates[i] = bootstrap_replicate_1d(data, func)

    return bs_replicates
```

## Bootstrap replicates of the mean and the SEM

In this exercise, you will compute a bootstrap estimate of the probability density function of the mean annual rainfall at the Sheffield Weather Station. Remember, we are estimating the mean annual rainfall we would get if the Sheffield Weather Station could repeat all of the measurements from 1883 to 2015 over and over again. This is a *probabilistic* estimate of the mean. You will plot the PDF as a histogram, and you will see that it is Normal.

In fact, it can be shown theoretically that under not-too-restrictive conditions, the value of the mean will always be Normally distributed. (This does not hold in general, just for the mean and a few other statistics.) The standard deviation of this distribution, called the **standard error of the mean**, or SEM, is given by the standard deviation of the data divided by the square root of the number of data points. I.e., for a data set,

`sem = np.std(data) / np.sqrt(len(data))`. Using hacker statistics, you get this same result without the need to derive it, but you will verify this result from your bootstrap replicates.

The dataset has been pre-loaded for you into an array called `rainfall`.

In [12]:

```
# Take 10,000 bootstrap replicates of the mean: bs_replicates
bs_replicates = draw_bs_reps(rainfall,np.mean,10000)

# Compute and print SEM
sem = np.std(rainfall) / np.sqrt(len(rainfall))
print(sem)

# Compute and print standard deviation of bootstrap replicates
bs_std = np.std(bs_replicates)
print(bs_std)

# Make a histogram of the results

#C:\Users\smpet\Anaconda3\lib\site-packages\ipykernel_launcher.py:13:
MatplotlibDeprecationWarning:
#The 'normed' kwarg was deprecated in Matplotlib 2.1 and will be removed in 3.1. Use 'density' ins
tead.
#  del sys.path[0]

_ = plt.hist(bs_replicates, bins=50, density=True)
_ = plt.xlabel('mean annual rainfall (mm)')
_ = plt.ylabel('PDF')

# Show the plot
plt.show()
```
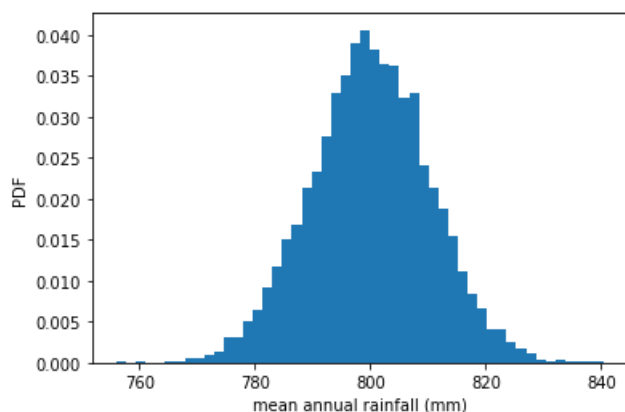
10.510549150506188
10.364395840934675

## Confidence intervals of rainfall data

A *confidence interval* gives upper and lower bounds on the range of parameter values you might
expect to get if we repeat our measurements. For named distributions, you can compute them
analytically or look them up, but one of the many beautiful properties of the bootstrap method is that
you can take percentiles of your bootstrap replicates to get your confidence interval. Conveniently,
you can use the `np.percentile()` function.

Use the bootstrap replicates you just generated to compute the 95% confidence interval. That is, give
the 2.5th and 97.5th percentile of your bootstrap replicates stored as `bs_replicates`. What is the
95% confidence interval?

In [13]:

```python
np.percentile(bs_replicates,[2.5,97.5])
```

Out[13]:

```
array([779.42353383, 820.23086466])
```

## Bootstrap replicates of other statistics

We saw in a previous exercise that the mean is Normally distributed. This does
not necessarily hold for other statistics, but no worry: as hackers, we can
always take bootstrap replicates! In this exercise, you'll generate bootstrap
replicates for the variance of the annual rainfall at the Sheffield Weather
Station and plot the histogram of the replicates.

In [14]:

```python
# Generate 10,000 bootstrap replicates of the variance: bs_replicates
bs_replicates = draw_bs_reps(rainfall,np.var,10000)

# Put the variance in units of square centimeters
bs_replicates = bs_replicates/100

# Make a histogram of the results
_ = plt.hist(bs_replicates, bins=50, density=True)
_ = plt.xlabel('variance of annual rainfall (sq. cm)')
_ = plt.ylabel('PDF')

# Show the plot
```
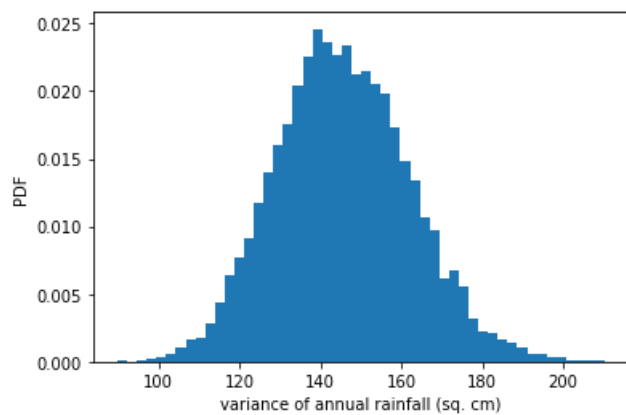
```
plt.show()
```



Great work! This is not normally distributed, as it has a longer tail to the right. Note that you can also compute a confidence interval on the variance, or any other statistic, using `np.percentile()` with your bootstrap replicates.

## Generating a permutation sample

In the video, you learned that permutation sampling is a great way to simulate the hypothesis that two variables have identical probability distributions. This is often a hypothesis you want to test, so in this exercise, you will write a function to generate a permutation sample from two data sets.

Remember, a permutation sample of two arrays having respectively `n1` and `n2` entries is constructed by concatenating the arrays together, scrambling the contents of the concatenated array, and then taking the first `n1` entries as the permutation sample of the first array and the last `n2` entries as the permutation sample of the second array.

In [15]:

```python
def permutation_sample(data1, data2):
    """Generate a permutation sample from two data sets."""

    # Concatenate the data sets: data
    data = np.concatenate((data1,data2))

    # Permute the concatenated array: permuted_data
    permuted_data = np.random.permutation(data)

    # Split the permuted array into two: perm_sample_1, perm_sample_2
    perm_sample_1 = permuted_data[:len(data1)]
    perm_sample_2 = permuted_data[len(data1):]

    return perm_sample_1, perm_sample_2
```

# Visualizing permutation sampling

To help see how permutation sampling works, in this exercise you will generate permutation samples and look at them graphically.

We will use the Sheffield Weather Station data again, this time considering the monthly rainfall in June (a dry month) and November (a wet month). We expect these might be differently distributed, so we will take permutation samples to see how their ECDFs *would look if* they were identically distributed.

The data are stored in the Numpy arrays `rain_june` and `rain_november`.

In [16]:

```
rain_june = np.array([ 66.2,   39.7,   76.4,   26.5,   11.2,   61.8,    6.1,   48.4,   89.2,
       104. ,   34. ,   60.6,   57.1,   79.1,   90.9,   32.3,   63.8,   78.2,
        27.5,   43.4,   30.1,   17.3,   77.5,   44.9,   92.2,   39.6,   79.4,
        66.1,   53.5,   98.5,   20.8,   55.5,   39.6,   56. ,   65.1,   14.8,
        13.2,   88.1,    8.4,   32.1,   19.6,   40.4,    2.2,   77.5,  105.4,
        77.2,   38. ,   27.1,  111.8,   17.2,   26.7,   23.3,   77.2,   87.2,
        27.7,   50.6,   60.3,   15.1,    6. ,   29.4,   39.3,   56.3,   80.4,
        85.3,   68.4,   72.5,   13.3,   28.4,   14.7,   37.4,   49.5,   57.2,
        85.9,   82.1,   31.8,  126.6,   30.7,   41.4,   33.9,   13.5,   99.1,
        70.2,   91.8,   61.3,   13.7,   54.9,   62.5,   24.2,   69.4,   83.1,
        44. ,   48.5,   11.9,   16.6,   66.4,   90. ,   34.9,  132.8,   33.4,
       225. ,    7.6,   40.9,   76.5,   48. ,  140. ,   55.9,   54.1,   46.4,
        68.6,   52.2,  108.3,   14.6,   11.3,   29.8,  130.9,  152.4,   61. ,
        46.6,   43.9,   30.9,  111.1,   68.5,   42.2,    9.8,  285.6,   56.7,
       168.2,   41.2,   47.8,  166.6,   37.8,   45.4,   43.2])
```

In [17]:

```
rain_november = np.array([ 83.6,   30.9,   62.2,   37. ,   41. ,  160.2,   18.2,  122.4,   71.3,
        44.2,   49.1,   37.6,  114.5,   28.8,   82.5,   71.9,   50.7,   67.7,
       112. ,   63.6,   42.8,   57.2,   99.1,   86.4,   84.4,   38.1,   17.7,
       102.2,  101.3,   58. ,   82. ,  101.4,   81.4,  100.1,   54.6,   39.6,
        57.5,   29.2,   48.8,   37.3,  115.4,   55.6,   62. ,   95. ,   84.2,
       118.1,  153.2,   83.4,  104.7,   59. ,   46.4,   50. ,  147.6,   76.8,
        59.9,  101.8,  136.6,  173. ,   92.5,   37. ,   59.8,  142.1,    9.9,
       158.2,   72.6,   28. ,  112.9,  119.3,  199.2,   50.7,   44. ,  170.7,
        67.2,   21.4,   61.3,   15.6,  106. ,  116.2,   42.3,   38.5,  132.5,
        40.8,  147.5,   93.9,   71.4,   87.3,  163.7,  141.4,   62.6,   84.9,
        28.8,  121.1,   28.6,   32.4,  112. ,   50. ,   96.9,   81.8,   70.4,
       117.5,   41.2,  124.9,   78.2,   93. ,   53.5,   50.5,   42.6,   47.9,
        73.1,  129.1,   56.9,  103.3,   60.5,  134.3,   93.1,   49.5,   48.2,
       167.9,   27. ,  111.1,   55.4,   36.2,   57.4,   66.8,   58.3,   60. ,
       161.6,  112.7,   37.4,  110.6,   56.6,   95.8,  126.8])
```
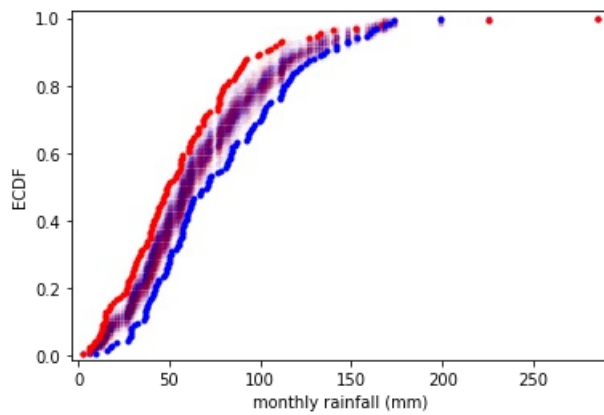
In [18]:

```
for _ in range(50):
    # Generate permutation samples
    perm_sample_1, perm_sample_2 = permutation_sample(rain_june, rain_november)


    # Compute ECDFs
    x_1, y_1 = ecdf(perm_sample_1)
    x_2, y_2 = ecdf(perm_sample_2)

    # Plot ECDFs of permutation sample
    _ = plt.plot(x_1, y_1, marker='.', linestyle='none',
                 color='red', alpha=0.02)
    _ = plt.plot(x_2, y_2, marker='.', linestyle='none',
                 color='blue', alpha=0.02)
```

```
# Create and plot ECDFs from original data
x_1, y_1 = ecdf(rain_june)
x_2, y_2 = ecdf(rain_november)
_ = plt.plot(x_1, y_1, marker='.', linestyle='none', color='red')
_ = plt.plot(x_2, y_2, marker='.', linestyle='none', color='blue')

# Label axes, set margin, and show plot
plt.margins(0.02)
_ = plt.xlabel('monthly rainfall (mm)')
_ = plt.ylabel('ECDF')
plt.show()
```



Great work! Notice that the permutation samples ECDFs overlap and give a purple haze. None of the ECDFs from the permutation samples overlap with the observed data, suggesting that the hypothesis is not commensurate with the data. June and November rainfall are not identically distributed.