## 4.2   Estimation of Sample Mean and Variance

### 4.2.1   Estimation and Sample Statistics

The mean and variance definitions of section 3.7 and section 3.8 apply only to a random variable for which we have a theoretical probability distribution. Typically, it is not until after having performed many measurements of a random variable that we can assign a good distribution model. Until then, measurements can help us *estimate* aspects of the data. We usually start by estimating basic parameters such as *mean* and *variance* before estimating a probability distribution.

There are two key aspects to randomness in the measurement of a random variable. First, of course, there is the underlying randomness with its probability distribution, mean, standard deviation, etc., which we call the *population statistics*. Second, there is the *statistical variability* that is due to the fact that we are *estimating* the random variable's statistics—called its *sample statistics*—from some sample. Statistical variability is decreased with greater sample size and number of samples, whereas the underlying randomness of the random variable does not decrease. Instead, our estimates of its probability distribution and statistics improve.

### 4.2.2   Sample Mean, Variance, and Standard Deviation

The *arithmetic mean* or **sample mean** of a measurand with sample size $N$, represented by random variable $X$, is defined as

$$\overline{x} = \frac{1}{N} \sum_{i=1}^{N} x_i.$$

If the sample size is large, $\overline{x} \to m_X$ (the sample mean approaches the mean). The **population mean** is another name for the mean $m_X$, which is equal to

$$m_X = \lim_{N \to \infty} \frac{1}{N} \sum_{i=1}^{N} x_i.$$

Recall that the *definition* of the mean is $m_X = \mathrm{E}[x]$.

The **sample variance** of a measurand represented by random variable $X$ is defined as

$$S_X^2 = \frac{1}{N-1} \sum_{i=1}^{N} (x_i - \overline{x})^2.$$

If the sample size is large, $S_X^2 \to \sigma_X^2$ (the sample variance approaches the variance). The **population variance** is another term for the variance $\sigma_X^2$, and can be expressed

as

$$\sigma_X^2 = \lim_{N \to \infty} \frac{1}{N-1} \sum_{i=1}^{N} (x_i - \overline{x})^2.$$

Recall that the *definition* of the variance is $\sigma_X^2 = \mathrm{E}\left[(X - m_X)^2\right]$.

The *sample standard deviation* of a measurand represented by random variable $X$ is defined as

$$S_X = \sqrt{S_X^2}.$$

If the sample size is large, $S_X \to \sigma_X$ (the sample standard deviation approaches the standard deviation). The *population standard deviation* is another term for the standard deviation $\sigma_X$, and can be expressed as

$$\sigma_X = \lim_{N \to \infty} \sqrt{S_X^2}.$$

Recall that the *definition* of the standard deviation is $\sigma_X = \sqrt{\sigma_X^2}$.

### 4.2.3   Sample Statistics as Random Variables

There is an ambiguity in our usage of the term "sample." It can mean just one measurement or it can mean a collection of measurements gathered together. Hopefully, it is clear from context.

In the latter sense, often we collect multiple samples, each of which has its own sample mean $\overline{X}_i$ and standard deviation $S_{X_i}$. In this situation, $\overline{X}_i$ and $S_{X_i}$ are themselves random variables (meta af, I know). This means they have their own sample means $\overline{\overline{X}}_i$ and $\overline{S_{X_i}}$ and standard deviations $S_{\overline{X}_i}$ and $S_{S_{X_i}}$.

The **mean of means** $\overline{\overline{X}}_i$ is equivalent to a mean with a larger sample size and is therefore our best estimate of the mean of the underlying random process. The **mean of standard deviations** $\overline{S_{X_i}}$ is our best estimate of the standard deviation of the underlying random process. The **standard deviation of means** $S_{\overline{X}_i}$ is a measure of the spread in our estimates of the mean. It is our best estimate of the standard deviation of the statistical variation and should therefore tend to zero as sample size and number of samples increases. The **standard deviation of standard deviations** $S_{S_{X_i}}$ is a measure of the spread in our estimates of the standard deviation of the underlying process. It should also tend to zero as sample size and number of samples increases.

Let $N$ be the size of each sample. It can be shown that the standard deviation of the means $S_{\overline{X}_i}$ can be estimated from a single sample standard deviation:

$$S_{\overline{X}_i} \approx \frac{S_{X_i}}{\sqrt{N}}.$$

This shows that as the sample size $N$ increases, the statistical variability of the mean decreases (and in the limit approaches zero).

### 4.2.4 Nonstationary Signal Statistics

The sample mean, variance, and standard deviation definitions, above, assume the random process is *stationary*—that is, its population mean does not vary with time. However, a great many measurement signals have populations that *do* vary with time, i.e. they are *nonstationary*. Sometimes the nonstationarity arises from a "drift" in the dc value of a signal or some other slowly changing variable. But dynamic signals can also change in a recognizable and predictable manner, as when, say, the temperature of a room changes when a window is opened or when a water level changes with the tide.

Typically, we would like to minimize the effect of nonstationarity on the signal statistics. In certain cases, such as drift, the variation is a nuissance only, but other times it is the point of the measurement.

Two common techniques are used, depending on the overall type of nonstationarity. If it is periodic with some known or estimated period, the measurement data series can be "folded" or "reshaped" such that the $i$th measurement of each period corresponds to the $i$th measurement of all other periods. In this case, somewhat counterintuitively, we can consider the $i$th measurements to correspond to a sample of size $N$, where $N$ is the number of periods over which measurements are made.

When the signal is aperiodic, we often simply divide it into "small" (relative to its overall trend) intervals over which statistics are computed, separately.

Note that in this discussion, we have assumed that the nonstationarity of the signal is due to a variable that is deterministic (not random).

### Example 4.2

Consider the measurement of the temperature inside a desktop computer chassis via an inexpensive *thermistor*, a resistor that changes resistance with temperature. The processor and power supply heat the chassis in a manner that depends on processing demand. For the test protocol, the processors are cycled sinusoidally through processing power levels at a frequency of 50 mHz for $n_T = 12$ periods and sampled at 1 Hz. Assume a temperature fluctuation between about 20 and 50 C and gaussian noise with standard deviation 4 C. Consider a *sample* to be the multiple measurements of a certain instant in the period.

1. Generate and plot simulated temperature data as a time series and as a histogram or frequency distribution. Comment on why the frequency distribution sucks.

2. Compute the sample mean and standard deviation *for each sample in the cycle*.
3. Subtract the mean from each sample in the period such that each sample distribution is centered at zero. Plot the composite frequency distribution of all samples, together. This represents our best estimate of the frequency distribution of the underlying process.
4. Plot a comparison of the theoretical mean, which is 35, and the sample mean of means with an error bar. Vary the number of samples $n_T$ and comment on its effect on the estimate.
5. Plot a comparison of the theoretical standard deviation and the sample mean of sample standard deviations with an error bar. Vary the number of samples $n_T$ and comment on its effect on the estimate.
6. Plot the sample means over a single period with error bars of ± one sample standard deviation of the means. This represents our best estimate of the sinusoidal heating temperature. Vary the number of samples $n_T$ and comment on the estimate.

We proceed in Python. First, load packages:

```python
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import norm
```

**Generate the Temperature Data**   The temperature data can be generated by constructing an array that is passed to a sinusoid, then "randomized" by gaussian random numbers.

Set a random seed for reproducible pseudorandom numbers.

```python
np.random.seed(43)
```

Define constants with

```python
f = 50e-3   # Hz
a = 15      # C
dc = 35     # C
fs = 1      # Hz
nT = 12     # number of sinusoid periods
s = 4       # C
np_ = int(fs / f + 1)  # number of samples per period
n = nT * np_ + 1  # total number of samples
```

Generate the temperature data.

```
t_a = np.linspace(0, nT / f, n)
sin_a = dc + a * np.sin(2 * np.pi * f * t_a)
noise_a = s * np.random.randn(n)
signal_a = sin_a + noise_a
```

Plot temperature over time

```
fig, ax = plt.subplots()
ax.plot(t_a, signal_a, 'o-', color='0.8', markerfacecolor='b', markersize=3)
plt.xlabel('time (s)')
plt.ylabel('temperature (C)')
plt.draw()
```
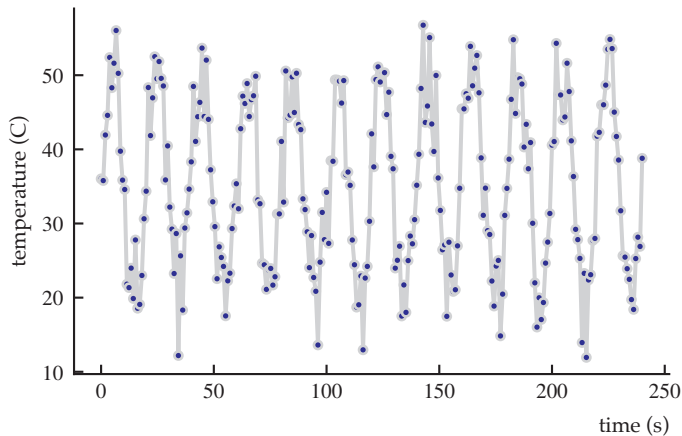


Figure 4.1. Raw temperature data over time.

This is something like what we might see for continuous measurement data. Now, the histogram.

```
fig, ax = plt.subplots()
ax.hist(signal_a, bins=30, density=True, alpha=0.5)
plt.xlabel('temperature (C)')
plt.ylabel('probability')
plt.draw()
```
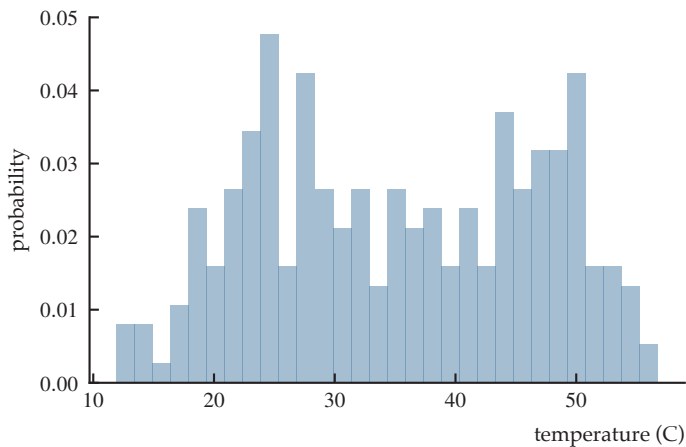
Figure 4.2. Raw temperature data histogram.

This sucks because we plot a frequency distribution to tell us about the random variation, but this data includes the sinusoid.

**Sample Mean, Variance, and Standard Deviation**   To compute the sample mean $\mu$ and standard deviation $s$ for each sample in the period, we must "pick out" the nT data points that correspond to each other. Currently, they're in one long $1 \times n$ array `signal_a`. It is helpful to *reshape* the data so it is in an nT $\times$ np array, which each row corresponding to a new period. This leaves the correct points aligned in columns. It is important to note that we can do this "folding" operation only when we know rather precisely the period of the underlying sinusoid. It is given in the problem that it is a controlled experiment variable. If we did not know it, we would have to estimate it, too, from the data.

Reshape data for sample mean, variance, and standard deviation calculations with

```
signal_ar = signal_a[:-1].reshape((nT, np_))
```

Compute sample mean, variance, and standard deviations with

```
mu_a = np.array([np.mean(col) for col in signal_ar.T])
var_a = np.array([np.var(col) for col in signal_ar.T])
s_a = np.array([np.std(col) for col in signal_ar.T])
```

**Composite Frequency Distribution** The columns represent samples. We want to subtract the mean from each column. We use `repmat` to reproduce `mu_a` in `nT` rows so it can be easily subtracted.

```
signal_arz = signal_ar - mu_a[np.newaxis,:]
x_a = np.linspace(-15, 15, 100)
pdfit_a = norm.pdf(x_a, loc=0, scale=s)
pdf_a = norm.pdf(x_a, loc=0, scale=s)
```

Now that all samples have the same mean, we can lump them into one big bin for the frequency distribution.

Plot composite frequency distribution with a probability distribution fit and the original probability distribution used to generate the data.

```
fig,ax = plt.subplots()
ax.hist(signal_arz.ravel(), bins=int(s * np.sqrt(nT)), density=True, alpha=0.5)
ax.plot(x_a, pdfit_a, 'b-', linewidth=2, label='pdf est.')
ax.plot(x_a, pdf_a, 'g--', linewidth=2, label='pdf')
plt.xlabel('Zero-mean temperature (C)')
plt.ylabel('Probability mass/density')
plt.legend()
plt.draw()
```
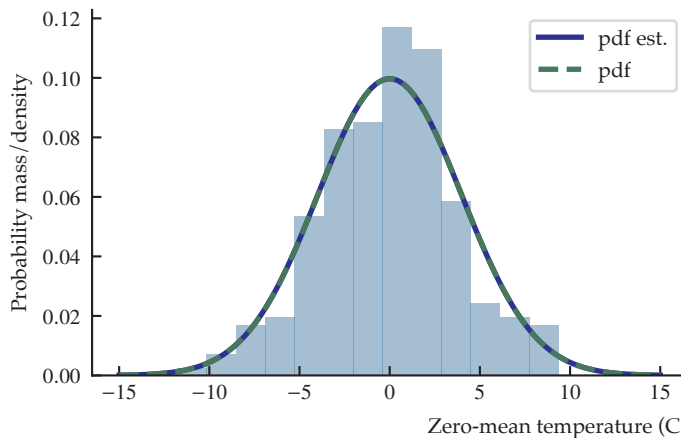


Figure 4.3. Composite frequency distribution of zero-mean temperature data.

**Means Comparison** The sample mean of means is simply the following:

```
mu_mu = np.mean(mu_a)
```

The standard deviation that works as an error bar, which should reflect how well we can estimate the point plotted, is the standard deviation of the means. It is difficult to compute this directly for a nonstationary process. We use the estimate given above and improve upon it by using the mean of standard deviations instead of a single sample's standard deviation.

```
s_mu = np.mean(s_a) / np.sqrt(nT)
```

Plot sample mean of means with an error bar as follows:

```
fig,ax = plt.subplots()
ax.bar(['$\overline{\overline{X}}$'], [mu_mu], yerr=s_mu, color='b', capsize=5)
plt.xlabel('Sample Mean of Means')
plt.draw()
```
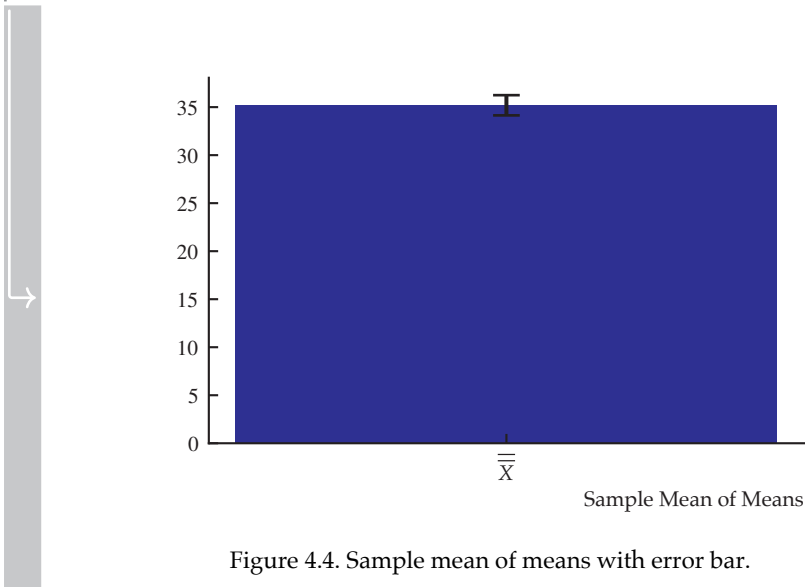


Figure 4.4. Sample mean of means with error bar.

**Standard Deviations Comparison**   The sample mean of standard deviations is simply the following:

```
mu_s = np.mean(s_a)
```

The standard deviation that works as an error bar, which should reflect how well we can estimate the point plotted, is the standard deviation of the standard deviations. We can compute this directly.

```
s_s = np.std(s_a)
```

Plot sample mean of standard deviations with error bar as follows:

```python
fig,ax = plt.subplots()
ax.bar(['$\overline{S_X}$'], [mu_s], yerr=s_s, color='b', capsize=5)
plt.xlabel('Sample Mean of Sample Standard Deviations')
plt.draw()
```
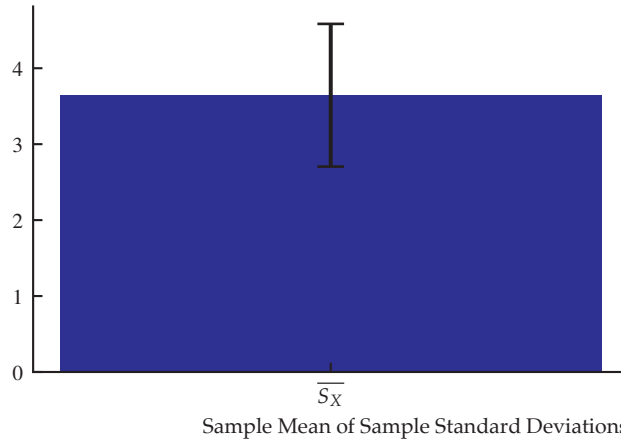


Figure 4.5. Sample mean of sample standard deviations with error bar.

**Plot a Period with Error Bars**   Plotting the data with error bars is fairly straight-forward. The main question is "which standard deviation?" Since we're plotting the means, it makes sense to plot the error bars as a single sample standard deviation of the means.

Plot sample means over a single period with error bars as follows:

```python
fig,ax = plt.subplots()
ax.errorbar(t_a[:np_], mu_a, yerr=s_a, fmt='o-', capsize=2, label='sample mean', color
t_a2 = np.linspace(0, 1 / f, 101)
ax.plot(t_a2, dc + a * np.sin(2 * np.pi * f * t_a2), 'r-', label='population mean')
plt.xlim([t_a[0], t_a[np_ - 1]])
plt.xlabel('Folded time (s)')
plt.ylabel('Temperature (C)')
plt.legend()
plt.show()  # Show all the plots
```
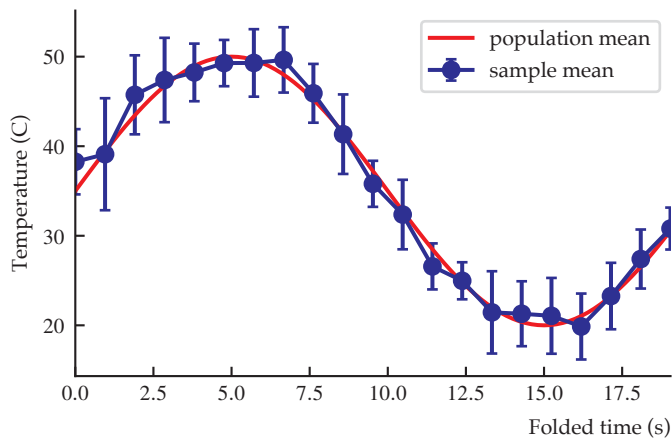
Figure 4.6. Sample means over a single period with error bars.

## 4.3   Confidence

One really ought to have it to give a lecture named it, but we'll give it a try anyway. **Confidence** is used in the common sense, although we do endow it with a mathematical definition to scare business majors, who aren't actually impressed, but indifferent. Approximately: if, under some reasonable assumptions (probabilistic model), we estimate the probability of some event to be $P\%$, we say we have $P\%$ confidence in it. I mean, business majors are all, "Supply and demand? Let's call that a 'law,'" so I think we're even.

So we're back to computing probability from distributions—probability density functions (PDFs) and probability mass functions (PMFs). Usually we care most about estimating the mean of our distribution. Recall from the previous lecture that when several samples are taken, each with its own mean, the mean is itself a random variable—with a mean, of course. Meanception.

But the mean has a probability distribution of its own. The **central limit theorem** has as one of its implications that, as the sample size $N$ gets large, *regardless of the sample distributions, this distribution of means approaches the Gaussian distribution*.

But sometimes I always worry I'm being lied to, so let's check.