

Physics 115/242

A summary of some important, and often poorly understood, results concerning the mean of the distribution, μ , the mean of a sample of N data points, \bar{x} , the standard deviation of the distribution, σ , the standard deviation of the data points, s , and the error bar on the mean, σ_μ .

Peter Young

(Dated: April 27, 2011)

Suppose we have a set of experimental data, x_i , ($i = 1, \dots, N$), which has some random noise. We shall often refer to this as a *sample* of data. The values of the x_i are governed by a distribution $P(x)$, *which we don't know*. This distribution has a mean $\mu \equiv \langle x \rangle$, and a variance σ^2 . (The term “standard deviation” is used for σ , the square root of the variance.) We denote an average *over the exact distribution* by angular brackets, e.g.

$$\mu \equiv \langle x \rangle = \int x P(x) dx. \quad (1a)$$

$$\sigma^2 \equiv \langle (x - \langle x \rangle)^2 \rangle = \langle x^2 \rangle - \langle x \rangle^2 = \int x^2 P(x) dx - \left(\int x P(x) dx \right)^2. \quad (1b)$$

Our goal is to determine $\langle x \rangle$, and the uncertainty in our estimate of it, from the N data points x_i . In order to do this we will assume that the data are uncorrelated with each other. This is a crucial assumption, without which it is very difficult to proceed. However, it is not always clear if the data points are truly independent of each other; some correlations may be present but not immediately obvious. Here, we take the usual approach of assuming that even if there are some correlations, they are sufficiently weak so as not to significantly perturb the results of the analysis.

The information *from the data* is usefully encoded in two parameters, the sample mean \bar{x} and the sample standard deviation s which are defined by¹

¹ The sample variance is often defined with a factor of $N - 1$ rather than N in Eq. (2b). However, to me this seems unnatural, and I prefer to define the sample variance in the naive way as the variance over the data. The reason that $N - 1$ is often put here, is so the factor of $(N - 1)/N$ in Eqs. (6d) and (7) become unity, and the factor of $N - 1$ in Eqs. (8) and (10) becomes N . The final result for the error bar in the mean, σ_μ , is unaffected by how the intermediate quantity s is defined. Furthermore, the difference between N and $N - 1$ is negligible for large N .

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i, \quad (2a)$$

$$s^2 \equiv \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 \quad (2b)$$

$$= \overline{(x - \bar{x})^2} = \overline{x^2} - (\bar{x})^2 = \frac{1}{N} \sum_{i=1}^N x_i^2 - \left(\frac{1}{N} \sum_{i=1}^N x_i \right)^2. \quad (2c)$$

In statistics, notation is often confusing but crucial to understand. Here, an average indicated by an over-bar, $\overline{\dots}$, is an average over the *sample of N data points*. This is to be distinguished from an exact average over the distribution $\langle \dots \rangle$, as in Eqs. (1a) and (1b). The latter is, however, just a theoretical construct since we *don't know* the distribution $P(x)$, only the set of N data points x_i which have been sampled from it.

Now we describe an important thought experiment. Let's *suppose* that we could repeat the set of N measurements *very many* many times, each time obtaining a value of the sample average \bar{x} . From these results we could construct a distribution, $\tilde{P}(\bar{x})$, for the sample average as shown in Fig. 1.

If we do enough repetitions we are effectively averaging over the exact distribution. Hence the average of the sample mean, \bar{x} , over very many repetitions of the data, is given by

$$\langle \bar{x} \rangle = \frac{1}{N} \sum_{i=1}^N \langle x_i \rangle = \langle x \rangle \equiv \mu, \quad (3)$$

i.e. it is the exact average over the distribution of x , as one would intuitively expect, see Fig. 1.

In fact, though, we have only the *one* set of data, so we can not determine μ exactly. However, Eq. (3) shows that

$$\boxed{\text{the best estimate of } \mu \text{ is } \bar{x}}, \quad (4)$$

i.e. the sample mean, since averaging the sample mean over many repetitions of the N data points gives the true mean of the distribution, μ . An estimate like this, which gives the exact result if averaged over many repetitions of the experiment, is said to be unbiased.

We would also like an estimate of the uncertainty or “error bar” in our estimate of \bar{x} for the exact average μ . This would be useful, for example, if we have a theoretical prediction for its

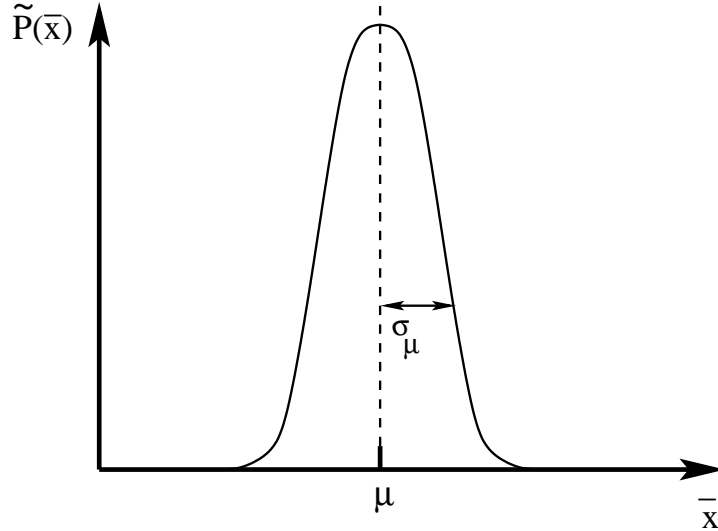


FIG. 1: The distribution of results for the sample mean \bar{x} obtained by repeating the measurements of the N data points x_i many times. The average of this distribution is μ , the exact average value of x . The width of this distribution is the uncertainty in the mean as determined from a single set of data, σ_μ .

value and would like to know if the experiment agrees with it. We can't tell unless we know the uncertainty in the experimental estimate.

We take σ_μ , the standard deviation in \bar{x} (obtained if one did many repetitions of the N measurements), to be the uncertainty, or error bar, in \bar{x} . This is the width of the distribution $\tilde{P}(\bar{x})$ shown in Fig. 1. A *single* estimate \bar{x} typically differs from the exact result μ by an amount of order σ_μ .

We shall now show that the variance of the *mean* of a set of N random variables is the variance of *one* variable divided by N . To see this, we have

$$\sigma_\mu^2 \equiv \langle \bar{x}^2 \rangle - \langle \bar{x} \rangle^2 = \left\langle \left(\frac{1}{N} \sum_{i=1}^N x_i \right)^2 \right\rangle - \left\langle \left(\frac{1}{N} \sum_{i=1}^N x_i \right) \right\rangle^2 \quad (5a)$$

$$= \frac{1}{N^2} \sum_{i=1}^N (\langle x_i x_j \rangle - \langle x_i \rangle \langle x_j \rangle) \quad (5b)$$

$$= \frac{1}{N^2} \sum_{i=1}^N (\langle x_i^2 \rangle - \langle x_i \rangle^2) \quad (5c)$$

$$= \frac{1}{N} (\langle x^2 \rangle - \langle x \rangle^2) \quad (5d)$$

$$\boxed{= \frac{\sigma^2}{N}}. \quad (5e)$$

To get from Eq. (5b) to Eq. (5c) we note that, for $i \neq j$, $\langle x_i x_j \rangle = \langle x_i \rangle \langle x_j \rangle$ since x_i and x_j are assumed to be statistically independent. (This is where the statistical independence of the data is

needed.)

The problem with Eq. (5e) is that **we don't know** σ^2 since it is a function of the exact distribution $P(x)$. We do, however, know the *sample* variance s^2 , see Eq. (2b), and the average of this over many repetitions of the N data points, is related to σ^2 since

$$\langle s^2 \rangle = \frac{1}{N} \sum_{i=1}^N \langle x_i^2 \rangle - \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \langle x_i x_j \rangle \quad (6a)$$

$$= \langle x^2 \rangle - \frac{1}{N^2} [N(N-1)\langle x \rangle^2 + N\langle x^2 \rangle] \quad (6b)$$

$$= \frac{N-1}{N} [\langle x^2 \rangle - \langle x \rangle^2] \quad (6c)$$

$$= \frac{N-1}{N} \sigma^2. \quad (6d)$$

To get from Eq. (6a) to Eq. (6b), we have separated the terms with $i = j$ in the last term of Eq. (6a) from those with $i \neq j$, and used the facts that each of the x_i is chosen from the same distribution and is statistically independent of the others. It follows from Eq. (6d) that

$$\boxed{\text{the best estimate of } \sigma^2 \text{ is } \frac{N}{N-1} s^2,} \quad (7)$$

since averaging $[N/(N-1)]s^2$ over many repetitions of N data points gives σ^2 . The estimate for σ^2 in Eq. (7) is therefore unbiased.

Combining Eqs. (5e) and (7) gives

$$\boxed{\text{the best estimate of } \sigma_\mu^2 \text{ is } \frac{s^2}{N-1}.} \quad (8)$$

This estimate is also unbiased. We have now obtained, using only information from from the data, that the mean is given by

$$\boxed{\mu = \bar{x} \pm \sigma_\mu.} \quad (9)$$

where we estimate

$$\boxed{\sigma_\mu = \frac{s}{\sqrt{N-1}}.} \quad (10)$$

Remember that \bar{x} and s are the mean and standard deviation of the (one set) of data that is available to us, see Eqs. (2a) and (2b).

As an example, suppose $N = 5$ and the data points are

$$x_i = 10, 11, 12, 13, 14, \quad (11)$$

(not very random looking data it must be admitted!). Then, from Eq. (2a) we have $\bar{x} = 12$, and from Eq. (2b)

$$s^2 = \frac{1}{5} [(-2)^2 + (-1)^2 + 0^2 + 1^2 + 2^2] = \frac{10}{5} = 2. \quad (12)$$

Hence, from Eq. (10),

$$\sigma_{\mu} = \frac{1}{\sqrt{4}} \sqrt{2} = \frac{1}{\sqrt{2}}. \quad (13)$$

so

$$\mu = 12 \pm \frac{1}{\sqrt{2}}. \quad (14)$$

For large N we can neglect factors of -1 compared with N , so, from Eq. (6d), s is equal to σ , and hence, from Eq. (10), we have

the error bar in the mean goes down like $1/\sqrt{N}$.

Hence, to reduce the error bar by a factor of 10 one needs 100 times as much data. This is discouraging, but is a fact of life when dealing with random noise.

For Eq. (10) to be really useful we need to know the probability that the true answer μ lies more than σ_{μ} away from our estimate \bar{x} . Fortunately, for large N the central limit theorem tells us (for distributions where the first two moments are finite) that the distribution of \bar{x} is a Gaussian, even if the distribution of the individual data points is not Gaussian. For the Gaussian distribution we know that the probability of finding a result more than one standard deviation away from the mean is 32%, more than two standard deviations is 4.5% and more than three standard deviations is 0.3%. Hence we expect that most of the time \bar{x} will be within σ_{μ} of the correct result μ , and only occasionally will be more than two times $\sigma_{\bar{x}}$ from it. Even if N is not very large, so there are some deviations from the Gaussian form, the above numbers may be a reasonable guide. Hence, if N is large and the theoretical prediction differs from the experimental value of \bar{x} by several times σ_{μ} , or more, there is likely to be either some systematic error in the experiment, or the theory does not apply.

Note, though, that distributions found in nature tend to have a much greater weight “in the wings” of the distribution than the Gaussian, which falls off very fast at large $(x - \mu)/\sigma$. Hence, if N is not large enough for the central limit theorem to apply, the probabilities of getting a certain number of standard deviations from the mean will be larger than the figures quoted above.

As an aside, although scientists usually quote σ_{μ} as the statistical uncertainty in \bar{x} , by convention, surveys of voters in elections use $2\sigma_{\mu}$ as a measure of the statistical uncertainty.