

1 Some Statistical Basics.

Statistics treats random errors. (There are also *systematic* errors; e.g., if your watch is 5 minutes fast, you will always get the wrong time, but it won't be random.) The two most important equations which describe the distribution of measured values in the presence of random error are the Gaussian distribution and the Poisson distribution. The first treats measured quantities which have a continuous range of values, while the second applies to cases where only discrete results can occur (e.g., the number of radioactive nuclei that decay in a given interval of time: there might be 41 decays or 42 decays but never 41.3744 ...).

1.1 The Poisson Equation

Let us consider a process which produces discrete events, but with a fixed *average* rate R per unit time. That is, if we count the events for a very long time, then

$$\lim_{t \rightarrow \infty} \left\{ \frac{\text{number of events}}{\text{length of time}} \right\} \rightarrow R \quad (1)$$

If we then observe for a (short) time t , then the *expected* number of events will be $N = R t$. But if the individual events are random, the result may not be N , it may be more or less than N . In fact, there is no reason why N should be an integer, while the number n we actually observe in time interval t *must* be integral. Given an expected number N , the Poisson distribution, $P_n(N)$, is the probability that our measurement will find a particular number of events n .

The Poisson equation is given by the expression:

$$P_n(N) = \frac{N^n}{n!} e^{-N} \quad (2)$$

Any observation is certain to yield *some* number of events, either none ($n = 0$), or $n = 1, 2, 3, \dots$. So if we add up all the $P_n(N)$ for a given N , we must obtain unity. Let's check that:

$$\sum_{n=0}^{\infty} P_n(N) = \sum_{n=0}^{\infty} \frac{N^n}{n!} e^{-N} = e^{-N} \left[1 + N + \frac{1}{2!}N^2 + \frac{1}{3!}N^3 + \frac{1}{4!}N^4 + \dots \right] \quad (3)$$

But the term in $[\dots]$ is just the series expansion of e^N , so

$$\sum_{n=0}^{\infty} P_n(N) = e^{-N} e^N = 1 \quad (4)$$

Now suppose we took a whole series of observations, each for a length of time t , so that the expected value each time is N . The actual number n seen in each observation will vary, but we could average the results of the series to find some mean number \bar{n} . If the Poisson distribution is obeyed, then we can compute the mean number we should see: We multiply the number of events observed, n , by the probability $P_n(N)$ that that number occurs, and sum over all possible values of n :

$$\bar{n} = \sum_{n=0}^{\infty} n P_n(N) = \sum_{n=0}^{\infty} \frac{N^n}{(n-1)!} e^{-N} = N e^{-N} \sum_{n=1}^{\infty} \frac{N^{n-1}}{(n-1)!} \quad (5)$$

Now the last sum is just e^N again, as we see if we let $n - 1 = k$:

$$\bar{n} = N e^{-N} \sum_{k=0}^{\infty} \frac{N^k}{k!} = N e^{-N} e^N = N . \quad (6)$$

This is just a check that equation (2) makes sense, since we already said that N represents the expected (average) number of events observed.

Lets look at some examples of this distribution:

N	$P_0(N)$	$P_1(N)$	$P_2(N)$	$P_3(N)$	$P_4(N)$	$P_5(N)$	$P_6(N)$	$P_7(N)$	$P_8(N)$
1	0.367879	0.367879	0.183940	0.061313	0.015328	0.003066	0.000511	0.000073	0.000009
1.5	0.22313	0.334695	0.251021	0.125511	0.047067	0.014120	0.003530	0.000756	0.000142
2	0.135335	0.270671	0.270671	0.180447	0.090224	0.036089	0.012030	0.003437	0.000859
2.5	0.082085	0.205212	0.256516	0.213763	0.133602	0.066801	0.027834	0.009941	0.003106
3	0.049787	0.149361	0.224042	0.224042	0.168031	0.100819	0.050409	0.021604	0.008102
5	0.006738	0.033690	0.084224	0.140374	0.175467	0.175467	0.146223	0.104445	0.065278
10	0.000045	0.000454	0.002270	0.007567	0.018917	0.037833	0.063056	0.090079	0.112599

After the mean value \bar{n} , the most important quantity which describes the Poisson distribution is some measure of how much, on the average, a particular value n will differ from the mean. You might think that you could just take the quantity of interest, $n - \bar{n}$, multiply by $P_n(N)$, and sum over all n . That would tell us nothing, because the result would be zero. To see why, consider the values of the $\Delta n = n - \bar{n}$: there will be as many negative as positive values. (Try it and see.) The point is that we don't so much care about the sign of Δn as its magnitude. We could compute the mean value of the absolute values $|\Delta n|$. But absolute values are not so easy to work with. It turns out that the best quantity to consider is the mean of the *squares* of the Δn , which are of course always positive. This is called the **variance**, V , of the distribution:

$$V = \sum_{n=0}^{\infty} (\Delta n)^2 P_n(N) = \sum_{n=0}^{\infty} (n - \bar{n})^2 P_n(N) = \sum_{n=0}^{\infty} (n - N)^2 P_n(N) . \quad (7)$$

Let's try to evaluate this. Expand the square to get

$$V = \sum_{n=0}^{\infty} n^2 P_n(N) - 2N \sum_{n=0}^{\infty} n P_n(N) + N^2 \sum_{n=0}^{\infty} P_n(N) \quad (8)$$

From equations (4), (5), and (6), we already have the values of the second two sums. Thus the last terms become $-2N \times N + N^2 = -N^2$, and we have

$$V = \sum_{n=0}^{\infty} n^2 P_n(N) - N^2 . \quad (9)$$

Now we must evaluate mean of the squares, $\overline{n^2}$. This is just

$$\overline{n^2} = \sum_{n=0}^{\infty} n^2 P_n(N) = \sum_{n=1}^{\infty} n^2 \frac{N^n}{n!} e^{-N} , \quad (10)$$

since the $n = 0$ term of the series is zero. We need to manipulate this a bit:

$$\overline{n^2} = N e^{-N} \sum_{n=1}^{\infty} \frac{n N^{n-1}}{(n-1)!} = N e^{-N} \sum_{k=0}^{\infty} \frac{(1+k)N^k}{k!} = N e^{-N} \left[\sum_{k=0}^{\infty} \frac{N^k}{k!} + \sum_{k=0}^{\infty} \frac{k N^k}{k!} \right]$$

$$\overline{n^2} = Ne^{-N} \left[e^N + N \sum_{k=1}^{\infty} \frac{N^{k-1}}{(k-1)!} \right] = Ne^{-N} \left[e^N + N \sum_{n=0}^{\infty} \frac{N^n}{n!} \right] = N + N^2$$

Going back to equation (9), we thus obtain our result

$$V = \overline{n^2} - N^2 = [N + N^2] - N^2 = N \quad (11)$$

Since the variance involves the square of Δn , we often take the square root of the variance, which is called the *standard deviation* or *root mean square deviation*, and is denoted by the Greek letter σ :

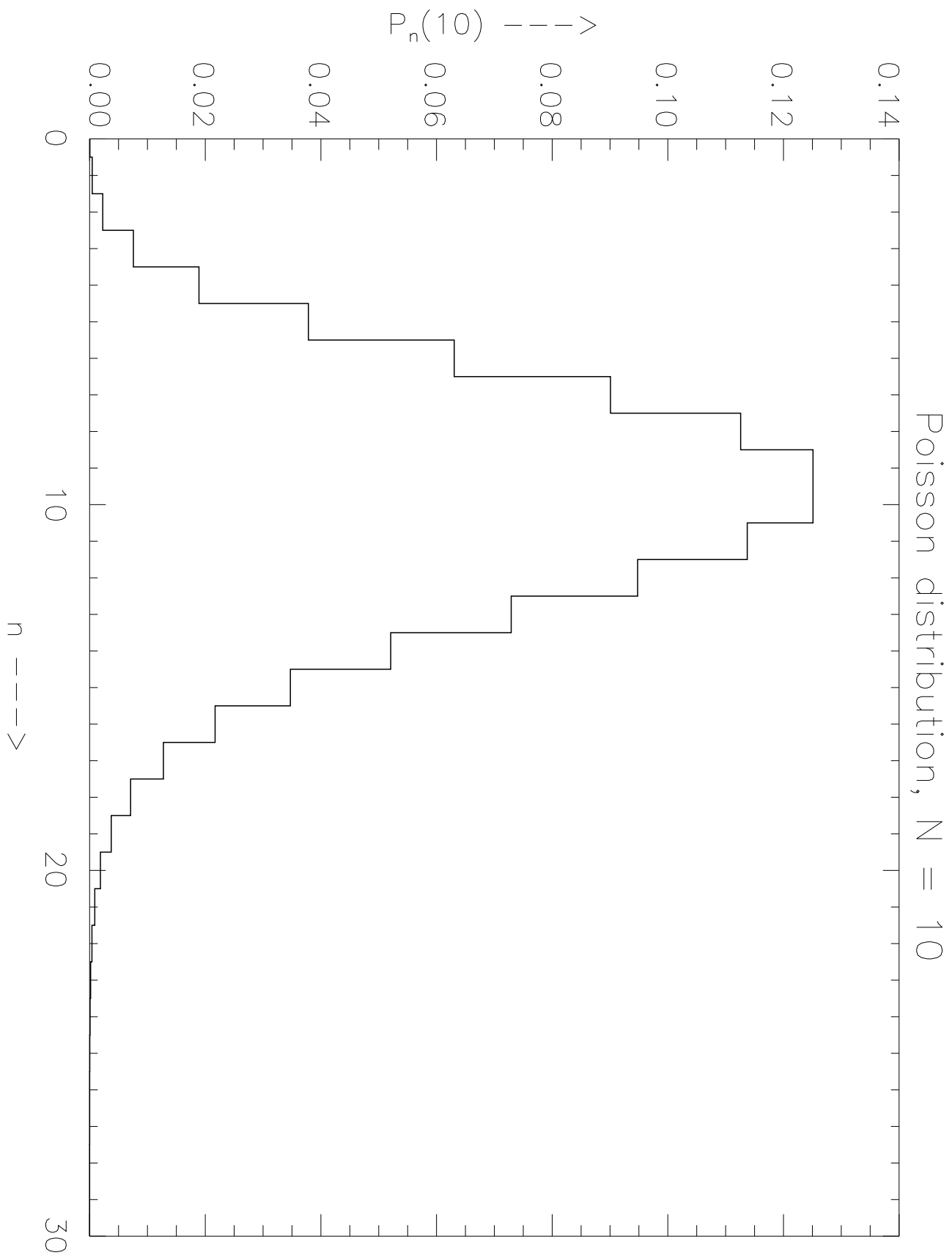
For the Poisson distribution: $\sigma \equiv \sqrt{V} = \sqrt{N}$

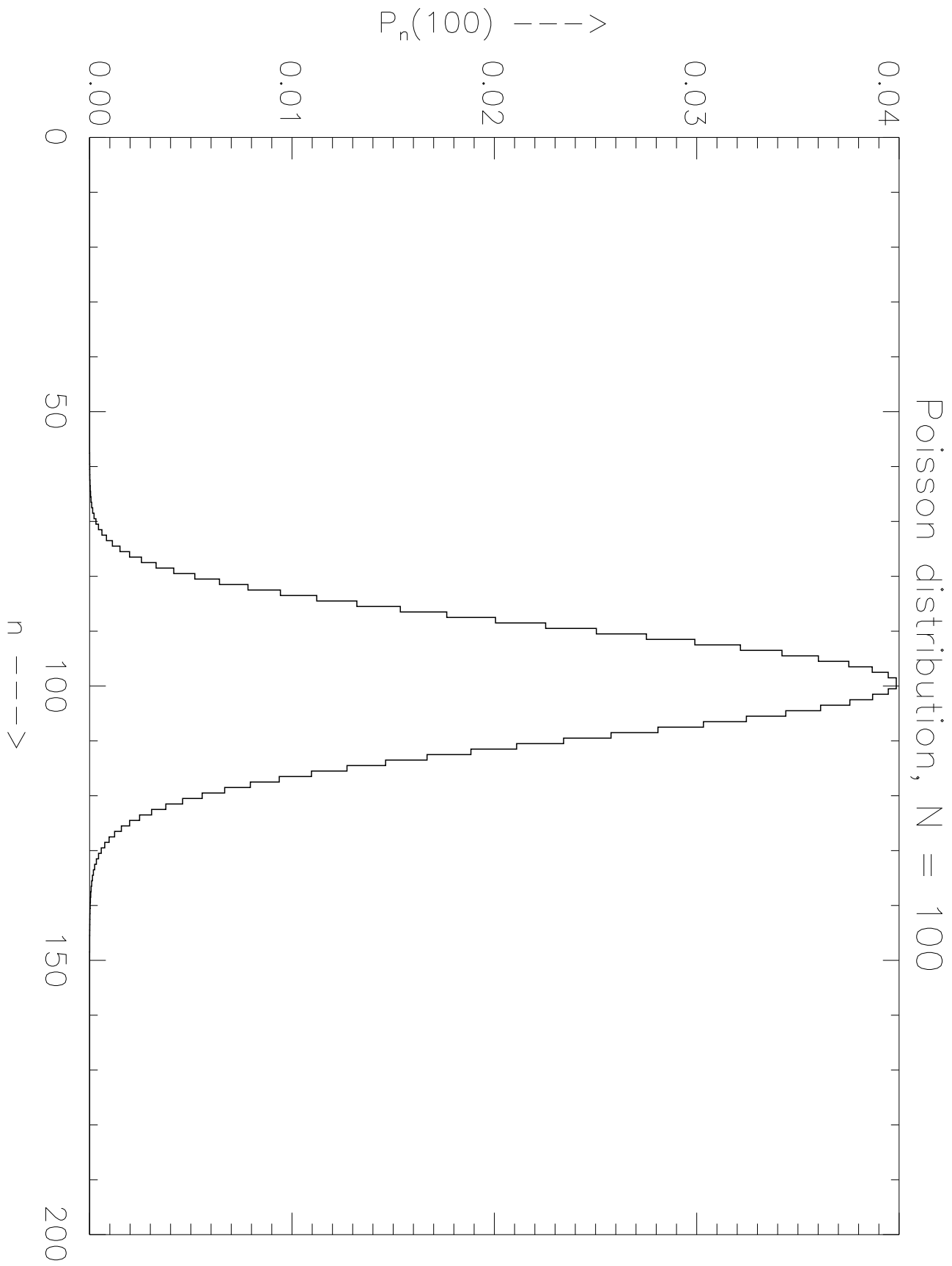
(12)

There are two important points: First, the spread in the values of n about \bar{n} increases as \bar{n} increases. Second, the spread (σ) does not increase as fast as \bar{n} . Thus, the *relative scatter* in the measured values *decreases* as $1/\sqrt{N}$:

N	$\sigma (= \sqrt{N})$	σ/N
10	3.16228	0.316228
100	10	0.100
1000	31.6228	0.0316228
10000	100	0.01000

To see how this applies to astronomical measurements, consider the case where we have a photo-multiplier tube attached to a telescope, and we are counting the photons from a faint star. What we really want to measure is the average rate of photons, R , as this is proportional to the star's brightness. If we count for some length of time t , then we can find R from equation (1). If the number we counted was N , then $R = N/t$ would be the right answer. But what we actually measure is some number n drawn from the Poisson distribution. How much error we make depends on how far n is likely to be from N . Now since N increases in proportion to t , the longer we count, the smaller σ/N will be, and the closer we are likely to be to the true brightness. But the accuracy does not increase very fast, only as \sqrt{t} . E.g., if you get 100 counts in one minute, which implies an accuracy of 10%, it will take 100 minutes to get 1% accuracy!





1.2 The Gaussian Distribution

The Gaussian distribution applies to continuous variables rather than the discrete events described by the Poisson distribution. The two are in fact related, and the histogram of a Poisson distribution for large N looks very much like a Gaussian distribution. Suppose we have a measured quantity x which can take on a continuous range of values $[-\infty < x < +\infty]$. Then the formula for the Gaussian distribution is:

$$P_\sigma(x) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(x-\bar{x})^2}{2\sigma^2}} \quad (13)$$

The distribution is a function of two parameters: the expected or *mean value* \bar{x} (like \bar{n} in the Poisson distribution), and the *standard deviation* σ . We interpret this formula in the following way: $P_\sigma(x) dx$ is the probability that our measurement will produce a value of x in the interval $[x, (x + dx)]$. Since any measurement must yield *some* value of x , integrating over all possible x must give unity:

$$\int_{-\infty}^{\infty} P_\sigma(x) dx = 1 . \quad (14)$$

The Gaussian distribution has its maximum value at $x = \bar{x}$, and that value is $P_\sigma(\bar{x}) = 1/\sqrt{2\pi}\sigma$. If we consider values one σ removed from \bar{x} , $x = \bar{x} \pm \sigma$, then from equation (13) we see that the value of $P_\sigma(x)$ will drop by a factor of $e^{-\frac{1}{2}} = 0.606531 \dots$. If we move out to $x = \bar{x} \pm 2\sigma$, the values will drop by $e^{-2} = 0.135335 \dots$.

In astronomy, you will often see a measure of the width of the Gaussian distribution which is a bit easier to visualize. This is called the *full width at half-maximum*, or “FWHM”. It is the width of the distribution at the point at which it has dropped to half its maximum value. This will happen for a value of x which satisfies the relation $P_\sigma(x) = \frac{1}{2}P_\sigma(\bar{x})$. From the definition (13), we see that this implies

$$e^{-\frac{(x-\bar{x})^2}{2\sigma^2}} = \frac{1}{2}$$

thus

$$\frac{(x-\bar{x})^2}{2\sigma^2} = \ln(2)$$

and

$$(x-\bar{x}) = \pm \sqrt{2 \ln(2)} \sigma = 1.17741 \dots \sigma$$

Since the *full-width* is just twice this value, we see that the relation between the FWHM and the standard deviation is:

$$\text{FWHM} = 2\sqrt{2 \ln(2)} \sigma = 2.35482 \sigma \quad (15)$$

We also might like to know the probability that the measured x will fall within a certain range. Thus we can look at the probability that x will be within one σ of the mean:

$$P[1\sigma] = \int_{\bar{x}-\sigma}^{\bar{x}+\sigma} P_\sigma(x) dx = \frac{1}{\sqrt{2\pi}} \int_{-1}^1 e^{-\frac{1}{2}y^2} dy = 0.6828 \dots \quad (16)$$

Likewise, the probability of falling within 2σ of the mean is

$$P[2\sigma] = \frac{1}{\sqrt{2\pi}} \int_{-2}^2 e^{-\frac{1}{2}y^2} dy = 0.9544 \dots \quad (17)$$

The corresponding probability of falling further away than 2σ is thus $1 - P[2\sigma] = 0.0456$. With increasing distance from \bar{x} , the probabilities decrease very rapidly; e.g., $P[3\sigma] = 0.9974$ so that $1 - P[3\sigma] = 0.0026$.

The standard deviation (if we can estimate it!) is sometimes used to put “error bars” on measured points on plots. Notice that from (16) you would expect almost 1/3 of the points to differ from the true value by *more than the error bars*.

Another quantity sometimes employed is the *probable error* (p.e.). This is the distance from the mean such that one half the values of x should be found within $\bar{x} \pm \text{p.e.}$ It turns out that the p.e. = 0.674σ .

So far we have been talking about \bar{x} and σ as if we knew their values. But for measured quantities, we would need an infinite number of measurements to find them exactly. We can only *estimate* them. It can be shown \dots — that the best estimate of \bar{x} from N measurements is

$$\bar{x}_N = \frac{1}{N} \sum_{i=1}^N x_i \quad (18)$$

Furthermore, the best estimate of σ is

$$\sigma_N = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x}_N)^2} \quad (19)$$

This quantity σ_N is the standard deviation of a single measurement. How good is our estimate \bar{x}_N ? It turns out that

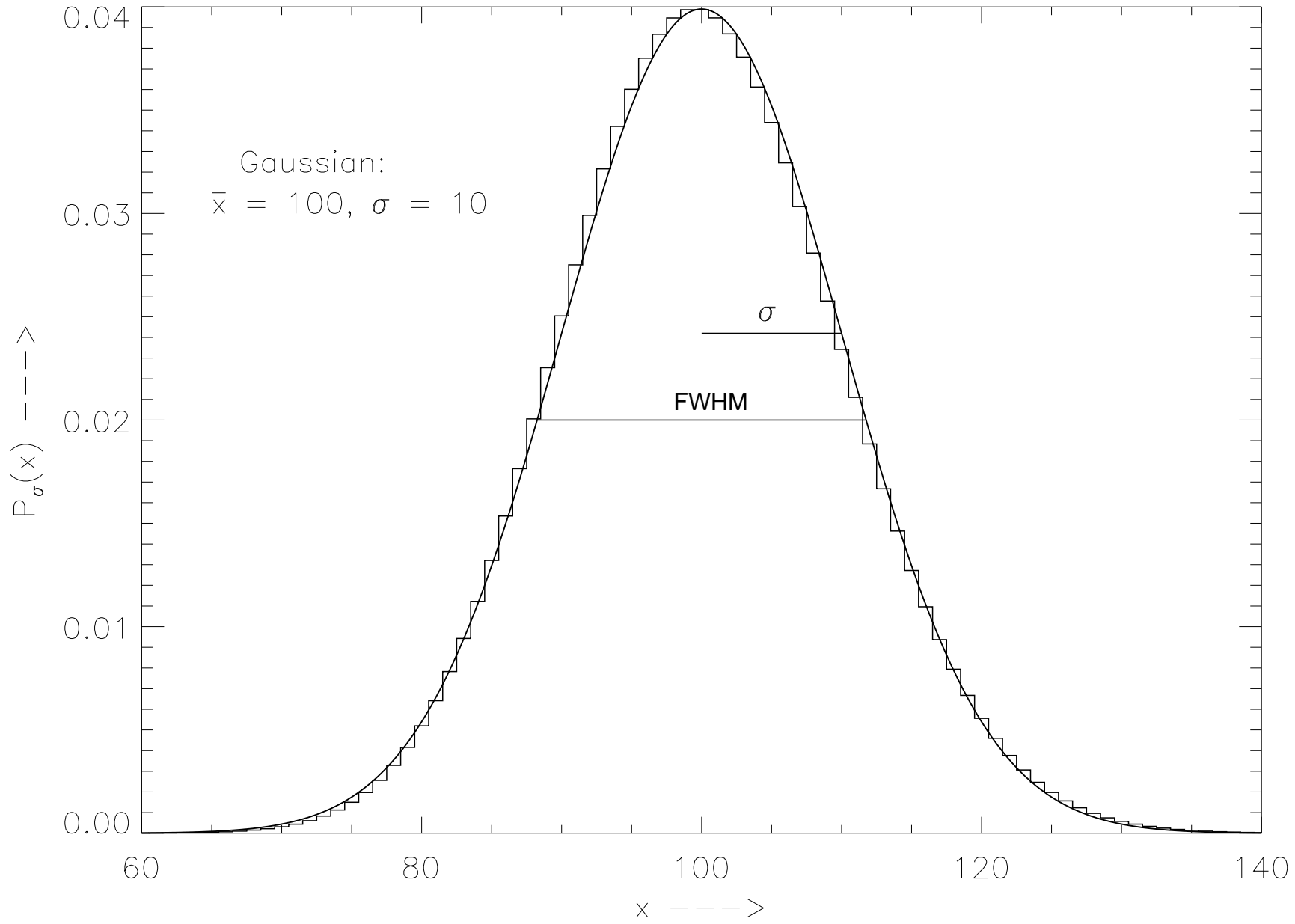
$$\bar{x} = \bar{x}_N \pm \sigma_{\bar{x}} = \bar{x}_N \pm \frac{\sigma_N}{\sqrt{N-1}} \quad (20)$$

The quantity

$$\sigma_{\bar{x}} = \frac{\sigma_N}{\sqrt{N-1}} \quad (21)$$

is called the *standard deviation of the mean*. Note that \bar{x}_N approaches the true value of the mean as $\sim 1/\sqrt{N}$ (just like the error in the Poisson distribution).

Gaussian & Poisson distributions compared.



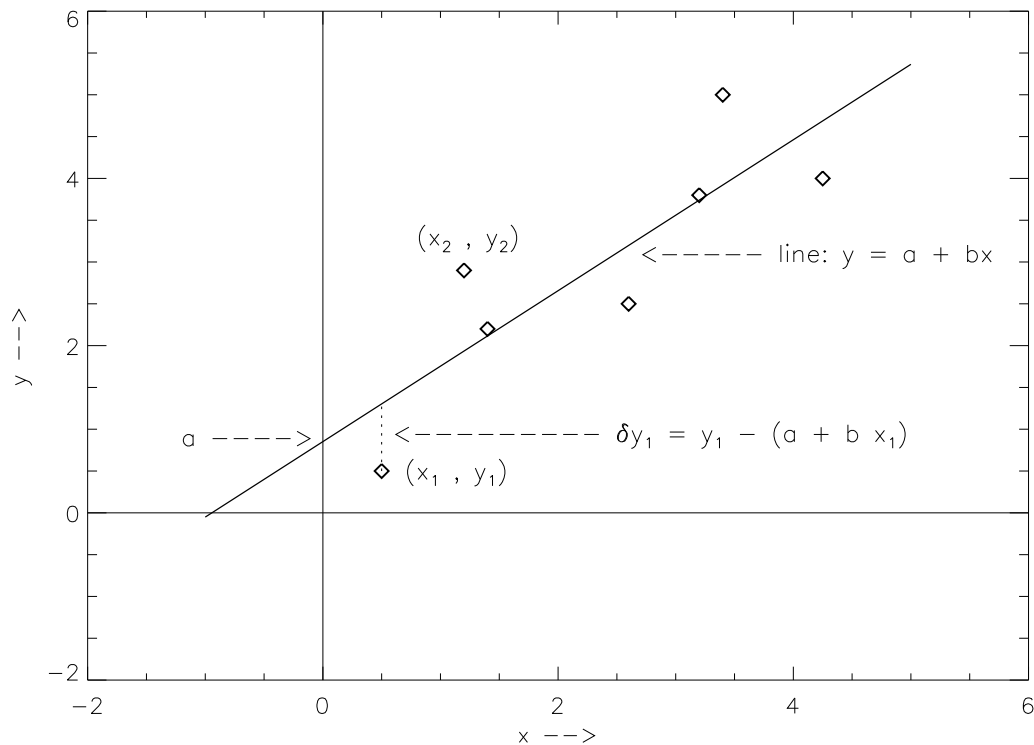
1.3 Least-Squares Fit to a Straight Line

We want to find the line which is the “best fit” to a set of points (x_i, y_i) . The line will be defined in terms of its slope b and y-intercept a :

$$y = a + b x \quad (22)$$

Then the error in y of the first point with respect to the line is

$$\delta y_1 = y_1 - (a + b x_1) \quad (23)$$



Let us choose as our measure of error the sum of the squares of all the individual errors:

$$V = \sum_{i=1}^N (\delta y_i)^2 = \sum_{i=1}^N [y_i - (a + b x_i)]^2 \quad (24)$$

We want to find the line for which V is a minimum. Thus we want to adjust the parameters a and b to obtain this minimum, and this will happen when

$$\frac{\partial V}{\partial a} = 0 \quad \text{and} \quad \frac{\partial V}{\partial b} = 0 \quad . \quad (25)$$

Expanding the square in equation (24), we have

$$V = \sum_{i=1}^N \left\{ y_i^2 - 2y_i a - 2y_i b x_i + a^2 + 2a b x_i + b^2 x_i^2 \right\} \quad (26)$$

We then see that the partial derivatives are

$$\frac{\partial V}{\partial a} = \sum_{i=1}^N \{ -2y_i + 2a + 2bx_i \} \quad (27)$$

$$\frac{\partial V}{\partial b} = \sum_{i=1}^N \{ -2x_i y_i + 2ax_i + 2bx_i^2 \} \quad (28)$$

Canceling the factors of 2, we see that these equations are just

$$N a + b \sum_{i=1}^N x_i = \sum_{i=1}^N y_i \quad (29)$$

$$a \sum_{i=1}^N x_i + b \sum_{i=1}^N x_i^2 = \sum_{i=1}^N x_i y_i \quad (30)$$

It saves the trouble of writing the summation signs if we use the over-bar notation:

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i, \quad \overline{xy} = \frac{1}{N} \sum_{i=1}^N x_i y_i, \quad \text{etc.} \quad (31)$$

Then equations (29) and (30) become

$$a + b\bar{x} = \bar{y} \quad (32)$$

$$a\bar{x} + b\overline{x^2} = \overline{xy} \quad (33)$$

We may use equation (32) to eliminate a from equation (33) to obtain

$$(\bar{y} - b\bar{x})\bar{x} + b\overline{x^2} = \overline{xy} \quad (34)$$

Rearranging, we have

$$b(\overline{x^2} - \bar{x} \cdot \bar{x}) = \overline{xy} - \bar{x} \cdot \bar{y} \quad (35)$$

Thus, the slope, b , of the *least-squares fit* line is

$$b = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{\overline{x^2} - \bar{x}^2} = \frac{N \sum_{i=1}^N x_i y_i - \sum_{i=1}^N x_i \sum_{i=1}^N y_i}{N \sum_{i=1}^N x_i^2 - \left(\sum_{i=1}^N x_i \right)^2} \quad (36)$$

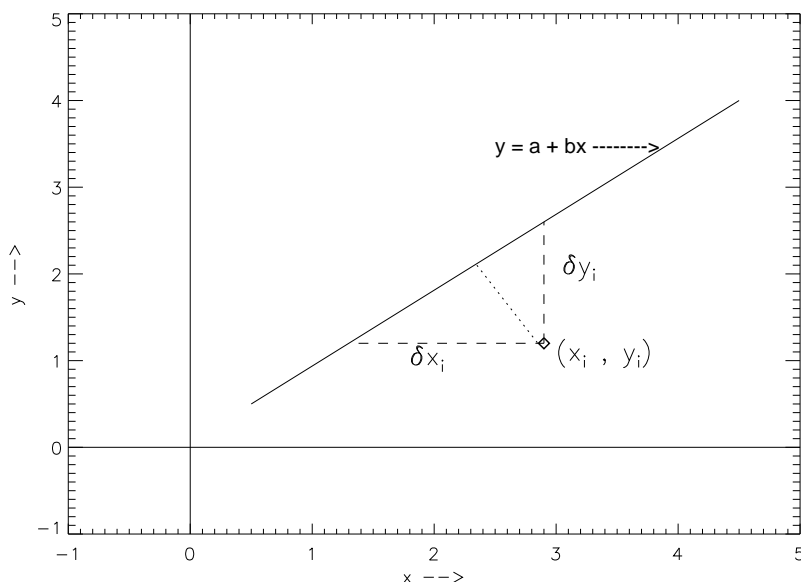
We now substitute this expression for b into equation (32) and solve for the intercept a :

$$a = \bar{y} - \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{\overline{x^2} - \bar{x}^2} \bar{x} \quad (37)$$

which reduces to

$$a = \frac{\bar{y} \cdot \overline{x^2} - \overline{xy} \cdot \bar{x}}{\overline{x^2} - \bar{x}^2} = \frac{\sum_{i=1}^N x_i^2 \sum_{i=1}^N y_i - \sum_{i=1}^N x_i \sum_{i=1}^N x_i y_i}{N \sum_{i=1}^N x_i^2 - \left(\sum_{i=1}^N x_i \right)^2} \quad (38)$$

Note that this derivation assumes that *all the error is in y!* This may or may not be a good assumption. Perhaps the errors are in x . Would we get the same result if we were to minimize the δx_i ? The answer is no.



Consider the point on the line which has the same y -value as some point (x_i, y_i) : The value of x on the line at $y = y_i$ is $x = (y_i - a)/b$. Thus the expression for δx_i will be

$$\delta x_i = x_i - \frac{y_i - a}{b} \quad (39)$$

We could go through the same process and minimize

$$V' = \sum_{i=1}^N (\delta x_i)^2 = \sum_{i=1}^N \left[x_i - \frac{1}{b}(y_i - a) \right]^2 \quad (40)$$

However, it is simpler to just write the line equation (22) in the form

$$x = \left(-\frac{a}{b} \right) + \left(\frac{1}{b} \right) y = A + B y \quad (41)$$

Then we see from (36) and (38) that

$$B = \frac{\overline{x^2} - \bar{x}^2}{\overline{xy} - \bar{x} \cdot \bar{y}} \quad (42)$$

and

$$A = \frac{\overline{xy} \cdot \bar{x} - \bar{y} \cdot \overline{x^2}}{\overline{xy} - \bar{x} \cdot \bar{y}} \quad (43)$$

Remember that we have assumed the errors are in y . Now, let us simply rename x as y and y as x . Then the equation of the line is

$$y = A + B x \quad (44)$$

and the least-squares fit, *with the errors now in x*, are just

$$B = \frac{\overline{y^2} - \bar{y}^2}{\overline{xy} - \bar{x} \cdot \bar{y}} \quad (45)$$

and

$$A = \frac{\overline{xy} \cdot \bar{y} - \bar{x} \cdot \overline{y^2}}{\overline{xy} - \bar{x} \cdot \bar{y}} \quad (46)$$

We see that $A \neq a$, $B \neq b$. Furthermore, we might have a case where the errors in the determination of x and of y were comparable. Then perhaps we would want to minimize the squares of the distance from the point to the line d_i :

$$(d_i)^2 = \frac{\delta x_i^2 \delta y_i^2}{\delta x_i^2 + \delta y_i^2} \quad (47)$$

This would give us yet another line, which would lie between the first two.

Finally, let us mention the *correlation coefficient*, r . This quantity is 0 if there is no correlation at all between the points, and reaches ± 1 for a collection of points which lie exactly on a straight line:

$$r = \frac{N \sum_{i=1}^N x_i y_i - \sum_{i=1}^N x_i \sum_{i=1}^N y_i}{\sqrt{\left[N \sum_{i=1}^N x_i^2 - \left(\sum_{i=1}^N x_i \right)^2 \right] \left[N \sum_{i=1}^N y_i^2 - \left(\sum_{i=1}^N y_i \right)^2 \right]}} \quad (48)$$