

*Since I don't know of a good reference that lays out model-fitting etc. in astronomy the way I like to think about it, I've started writing up notes. They are currently a work in progress and will periodically be updated/extended.*

The problem of “what model fits my data” is as old as experimental data <sup>1</sup>. An important category is when measurement errors have Gaussian distributions. Not only are Gaussians some of the easiest distributions to work with, but due to the central limit theorem, if we have enough data we'll end up at a Gaussian distribution anyways. The aim of this note is to show, starting from the probability distribution function (PDF) of a Gaussian, that we can compare the relative probabilities that different models would have given rise to our observed data.

## 1. From Gaussians to $\chi^2$

If we have a Gaussian random variable with mean  $\mu$  and standard deviation  $\sigma$ , then we want to know how likely we are that a realization of this variable gives us a value  $x$ . While the probability of a single value is zero for any continuous definition, we instead use the *probability density function* or PDF, where the probability of observing an event between  $x$  and  $x + \delta_x$  is  $\text{PDF}(x)\delta_x$ . Properly normalized PDFs should always be non-negative and integrate to unity. For a Gaussian, the PDF is:

$$\exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)/\sqrt{2\pi\sigma^2}$$

where the factor of  $\sqrt{2\pi\sigma^2}$  is needed so the PDF integrates to one.

In traditional model fitting, we are usually interested in figuring out that the data “should have been,” which mean we want to figure out  $\mu$ . We also usually take the errors  $\sigma$  as given. While we can certainly try to estimate the errors at the same time as a model, but it's quite a bit more complicated and so we don't try to do that here. We also usually have many data points, which may have their own errors and expected means. One simple case would be fitting a line to data taken at different times  $t$ , where the expected value of each data point would be  $at + b$ . In this case, if we measure the  $i^{\text{th}}$  data point at time  $t_i$ , then we have  $\mu_i = at_i + b$  for a slope  $a^2$  and intercept  $b$ . Each data point would have its own value, taken at its own time, but if a line is a good description of the data, they should all agree on the values for the slope and intercept.

If we have many data points, with *uncorrelated* errors, then the probability density of observing several data values is just the product of the probability density of the individual data points. In otherwords, if  $d_i$  are our  $n$  data points, then the PDF is

$$\prod_{i=1}^n \exp\left(-\frac{(d_i - \mu_i)^2}{2\sigma_i^2}\right)/\sqrt{2\pi\sigma_i^2}.$$

Of course, the product is rather awkward to work with, so let us instead work with the log of the

---

<sup>1</sup>citation needed

<sup>2</sup>We are using  $a$  here and not the usual  $m$  to avoid confusion down the line.

PDF:

$$\log \text{PDF} = \sum \frac{-(d_i - \mu_i)^2}{2\sigma_i^2} - \sum \frac{1}{2} \log(2\pi\sigma_i^2)$$

We can also pull a factor of  $-\frac{1}{2}$  to the left hand side to further simplify matters. The PDF is also often referred to as the likelihood  $\mathcal{L}$ , so we make that change as well, giving:

$$-2\log(\mathcal{L}) = \sum \frac{(d_i - \mu_i)^2}{\sigma_i^2} + \sum \log(2\pi\sigma_i^2)$$

With this expression for the likelihood in hand, we can now say how likely it is that given a model  $\mu_i$  and errors  $\sigma_i$  would have given rise to the observed data  $d_i$ . Of course, what we really want to ask is “given the data, what is the true model.” Unfortunately, that question is rather ill-posed (see *e.g.* the long literature on Bayesian statistics). Instead, the question we know how to answer is “given a model, how likely is it to have given rise to the data we got.” We can ask this question about many different models, and can compare the relative probabilities of two different models giving rise to the observed data. That is a well-defined question and the answer is just the ratio of the likelihoods/PDFs. In log space, the ratio becomes a difference, and we have for models  $\mu_i$  and  $\mu_i^\dagger$

$$-2\delta \log(\mathcal{L}) = \sum \frac{(x_i - \mu_i)^2}{\sigma_i^2} - \sum \frac{(x_i - \mu_i^\dagger)^2}{\sigma_i^2}$$

Happily, the second term in the likelihood has gone away, because it doesn’t depend on the model and so gets cancelled in the difference. Note that if we ever try to fit for the  $\sigma$ ’s, though, we would need to keep that term around.

Looking at the likelihood difference expression, we can see that the full information we need to know for a given model is just the sum, which is called  $\chi^2$

$$\chi^2 \equiv \sum \frac{(x_i - \mu_i)^2}{\sigma_i^2}$$

Once I have the value for  $\chi^2$  for a set of models in-hand, I can quickly compare the relative likelihoods by noting that  $-2\delta \log(\mathcal{L}) = \delta \chi^2$ , which gives a likelihood ratio for the two models of  $\exp(-\frac{\delta \chi^2}{2})$ . If the difference between two models is large, the model with the larger  $\chi^2$  is exponentially less likely to have given arise to the observed data. Similarly, if  $\delta \chi^2$  is small (where small means much less than 1), then the relative probabilities are very similar.

In short, we use  $\chi^2$  because for Gaussian random variables, everything we need to know about the relative probabilities of two models producing the observed data is encapsulated in  $\chi^2$ .

## 2. Searching for the “Best” Model

While it is useful to be able to compare two different models, usually what we want to do is find the model that does the “best” job fitting the data. In more general cases, even defining what we mean by this question can get quite complicated, and reasonable people can disagree. We can, however, get to an answer by defining the “best” model as the one that had the highest likelihood

of producing the observed data, which is usually described as the frequentist approach. Since we know that low values of  $\chi^2$  mean those models have relatively higher probabilities of producing our observed data, we can say that the “best” model is the one that minimizes  $\chi^2$ . Very often, we have models that depend on one or more continuous parameters. One example of this that we have already seen is fitting a line to data, where we can set the slope and intercept to be whatever we want. The problem then becomes, of all the infinite possible values, how do we find the ones that minimize  $\chi^2$ ?

In searching for the parameter values that minimize  $\chi^2$ , which from now on we will refer to as the “best-fit” values, we could directly differentiate  $\chi^2$  as we have written it, and possibly find a solution. However, it is usually simpler to understand what is going on, and how best to proceed, if we re-write  $\chi^2$  in the language of linear algebra. We can start by putting the observed data points and their expected values into column vectors  $d$  and  $\mu$ . Furthermore, we can define a diagonal matrix  $N$  where  $N_{i,i} = \sigma_i^2$ . In this case,  $N^{-1}$  is just trivially  $N_{i,i}^{-1} = 1/\sigma_i^2$ . If we then take  $N^{-1}(d - \mu)$  we can see that this makes a column vector, with  $i^{th}$  element just  $(d_i - \mu_i)/\sigma_i^2$ . Recalling that a row vector times a column vector is just a dot product, we can take  $(d - \mu)^T (N^{-1}(d - \mu))$ . This becomes  $\sum \frac{(d_i - \mu_i)^2}{\sigma_i^2}$ , which was just our original expression for  $\chi^2$ . So, in linear algebra notation, we have

$$\chi^2 = (d - \mu)^T N^{-1} (d - \mu)$$

. Of course, we usually have some form for what the data values should be, and usually that depends on some set of model parameters. In this case,  $\mu_i = \langle d_i \rangle = A_i(m)$  where  $A_i$  is some function that depends on both the model parameters and which data point we’re looking at. Again, in the case of fitting a line to data, we need to know both the slope/intercept *and* the  $x$ -value of a point in question before we can predict the  $y$ -value. Of course, we can as usual turn these things into vectors, giving us

$$\chi^2 = (d - A(m))^T N^{-1} (d - A(m))$$

where if we knew the true model parameters, we’d have  $\langle d \rangle = A(m_{true})$ . If we were arbitrarily good at math, we could differentiate  $\chi^2$  with respect to  $m$ , and find the global minimum, which would give us our best fit parameters. Unfortunately, this is usually very hard to do.

### 3. Linear Least Squares

Fortunately, there is a very broad class of models where we can analytically solve for our best-fit parameters. If we make the important simplification that we’ll restrict ourselves to models that depend *linearly* on the model parameters, it turns out we can indeed solve for the global minimum of  $\chi^2$ . What do we mean by a linear dependence? We mean

$$\langle d \rangle = Am \tag{1}$$

for a potentially arbitrary matrix  $A$  and fit parameters  $m$ . Generally, we decide on  $A$ , which corresponds to selecting which class of models we’re going to look at, and then finding the parameters that give us the best model in that whole class. Going back to the line-fitting example where  $\langle d_i \rangle = at_i + b$ , we could write  $A$  as an  $n$ -data by 2 matrix where the first column is  $t_i$  and the

second column is just a vector of ones. Our fit parameters  $m$  (for model parameters) is then just a column vector consisting of the slope and intercept  $m = [ab]$ . If you carry out the multiplication  $Am$  you can see that the individual elements are  $at_i + b$ , so line-fitting fits nicely into this formalism. It's important to note that once we said we were going to fit a line to our data, that specified  $A$ . We could have instead said we'd fit a quadratic or cubic polynomial, or sum of exponentials, or sum of sines and cosines, and that each of the choices would have specified a different version of  $A$ . For now, we will not worry about how to select between different choices of  $A$ , but just worry about how to pick the best parameters once we have selected  $A$ .

We can rewrite  $\chi^2$  to reflect the fact we have restricted ourselves to the linear case:

$$\chi^2 = (d - Am)^T N^{-1} (d - Am) \quad (2)$$

. We now have to find the values of  $m$  that minimize  $\chi^2$ . As usual, we do this by taking the derivative of  $\chi^2$  with respect to  $m$ . As shown in the appendix, this is just

$$\nabla \chi^2 = -2A^T N^{-1} (d - Am)$$

. At the minimum of  $\chi^2$ , we know this must be equal to zero and so we have

$$-2A^T N^{-1} (d - Am) = 0$$

or

$$A^T N^{-1} Am = A^T N^{-1} d \quad (3)$$

This is the fundamental equation of linear least-squares and people who analyze data would do well to ruminate on it periodically. When tackling large problems, I find it more useful to think about this equation in terms of operators.  $A$  takes a model and maps it to data.  $N^{-1}$  inverse variance-weights data.  $A^T$  takes data and projects it back into a something that looks like a model, but is not normalized. As an example, if we are using least-squares to solve for a map of the sky given data from a telescope scanning repeatedly over a patch of sky,  $A^T$  times a vector of ones would give us the “hit counts” – the map of how many times each map pixel was observed by a detector. We can then interpret Equation 3 as follows: The right-hand side is the inverse variance-weighted data, projected onto our fit parameters. The left hand side maps parameters into data, then inverse variance weights them and projects onto parameters. Our maximum likelihood solution is where these two are equal. If we somehow *knew* (or guessed) the solution, we could check it very quickly, just by operating on the data/model.

It is tempting to re-write Equation 3 as

$$m = (A^T N^{-1} A)^{-1} A^T N^{-1} d$$

Indeed, often that is exactly what we *should* do. However, in a wide range of common cases,  $A^T N^{-1} A$  will not lend itself to inversion. Either it is too large to invert (let alone construct) or is singular. Often it will be both. In such cases, iterative solutions to Equation 3 are preferred, where we guess a solution, run it through the left-hand side of Equation 3, and compare it to the right-hand side. Based on the difference, we update our guess and repeat until we're close enough.

### 3.1. Bias

: We have shown that Equation 3 gives the maximum-likelihood fit to our data provided that we know the noise in the data, and that our model can describe the data (Equation 1). Let's calculate the expected difference between our recovered model  $m$  and the true one  $m_t$ . We will also set  $d = d_t + n$  where  $d_t$  are the noise-free data, and  $n$  is the (unknown) realization of the noise we happened to get. Since  $d_t = Am_t$ , it is obviously true that:

$$A^T N^{-1} A m_t = A^T N^{-1} d_t$$

Note that this depends on Equation 1 being true, but critically *not* on  $N$  describing the variance in the data. We can then subtract this result from Equation 3 giving:

$$A^T N^{-1} A (m - m_t) = A^T N^{-1} (d - d_t) \quad (4)$$

Of course  $d - d_t$  is just  $n$ , so (assuming invertability):

$$m - m_t = (A^T N^{-1} A)^{-1} A^T N^{-1} n$$

*Usually* this will average to zero, since the noise is positive and negative with equal probability, meaning our maximum-likelihood parameters are unbiased. This is relied on Equation 1 being true, but not on  $N$  being an accurate description of the variance. The closer  $N$  is to the true variance, the closer  $m$  will be to the maximum-likelihood solution, but it should remain unbiased. One note of caution, however: extreme caution must be used when estimating the noise directly from the data, since we can inadvertently couple *errors* in the noise model to signal in the data, which can make the right-hand side of Equation 4 non-zero. A simple example would be using  $d^2$  as our variance when estimating the mean of the data. If the mean is positive, data that happened to fluctuate high will be assigned higher variance, and hence lower weight, than data that happened to fluctuate low. The result will be an estimate of the mean biased low. The usual solution to this problem is to subtract the model from the data before estimating the noise, then fitting the model. After multiple iterations this process should converge, but it can sometimes be slow when the SNR per-point is high. Other solutions are possible, including making a (possibly extremely) non-optimal but unbiased estimate of the model before estimating noise for the first time. While solutions vary, just being aware of the issue is most of the battle.

### 3.2. Model Variances

While the maximum-likelihood model is unbiased (subject to the warnings from the previous section), the parameters will still have errors. How do we describe these? If  $\delta_m \equiv m - m_t$ , then from Equation 4 we have

$$\begin{aligned} \delta_m \delta_m^T &= (A^T N^{-1} A)^{-1} A^T N^{-1} n \left[ (A^T N^{-1} A)^{-1} A^T N^{-1} n \right]^T \\ &= (A^T N^{-1} A)^{-1} A^T N^{-1} n n^T N^{-1} A (A^T N^{-1} A)^{-1} \end{aligned}$$

where we have used the fact that  $N$  is symmetric. If we take the expectation, then  $\langle nn^T \rangle = N$  which we can use to cancel out one factor of the noise matrix in the middle. This leaves us with:

$$\langle \delta_m \delta_m^T \rangle = (A^T N^{-1} A)^{-1} A^T N^{-1} A (A^T N^{-1} A)^{-1}$$

We can carry out one further cancellation, leaving us with

$$\langle \delta_m \delta_m^T \rangle = (A^T N^{-1} A)^{-1} \quad (5)$$

Fortunately, if we have directly solved our least-squares problem, we already have this matrix sitting around! While this gives us the full covariance of our model errors, note that the standard deviation of the model parameters is just the square root of the diagonal entries.

One further case of interest is when we care about the error in our prediction for the data and the truth. An instance of this would be fitting data with a polynomial, where the difference between our model and truth might be more interesting than the formal errors on the polynomial coefficients themselves. We know our estimate for the data is given by  $d = Am$ , so

$$d - d_t \equiv \delta_d = A(m - m_t) = A\delta_m$$

The expectation of this should be zero if our model is unbiased, but we can again take the covariance:

$$\langle \delta_d \delta_d^T \rangle = A\delta_m (A\delta_m)^T = A (A^T N^{-1} A)^{-1} A^T \quad (6)$$

## 4. Correlated Noise

So far, we have formulated everything assuming the errors in the data are uncorrelated  $\langle n_i n_j \rangle \propto \delta_{i,j}$ . Life is often not so kind to us, and frequently errors are correlated. A random walk is one case, where one knows that neighboring samples are highly correlated - if sample  $i$  is high, then sample  $i + 1$  must also be high. How can we deal with this? We start by adding in an extra invertible matrix  $S$  into  $\chi^2$  along with its inverse, so  $\chi^2$  remains unchanged:

$$\chi^2 = (d - Am)^T S^T S^{T,-1} N^{-1} S^{-1} S (d - Am)$$

Since  $S$  always shows up with its inverse, this *has* to have the same solution, model errors etc. as Equation 3. We can now group terms as follows:

$$\chi^2 = (Sd - SAM)^T (SNS^T)^{-1} (Sd - SAM)$$

We can now introduce new variables:

$$\tilde{d} = Sd$$

$$\tilde{A} = SA$$

$$\tilde{N} = SNS^T$$

We can now write  $\chi^2$  with the new variables:

$$\chi^2 = (\tilde{d} - \tilde{A}m)^T \tilde{N}^{-1} (\tilde{d} - \tilde{A}m)$$

By multiplying Equation 1 on the left by  $S$ , it is obviously true that

$$\langle \tilde{d} \rangle = \tilde{A}m$$

Less obvious is that

$$\tilde{N}_{ij} = \langle \tilde{n}_i \tilde{n}_j \rangle$$

We know that  $\tilde{n}_i = \sum_k S_{ik} n_k$ . So

$$\langle \tilde{n}_i \tilde{n}_j \rangle = \left\langle \sum_k S_{ik} n_k \sum_l S_{jl} n_l \right\rangle$$

However, the cross terms are all zero unless  $k = l$ , so this transforms into a single sum:

$$\langle \tilde{n}_i \tilde{n}_j \rangle = \left\langle \sum_k S_{ik} S_{jk} n_k^2 \right\rangle = \sum_k S_{ik} S_{jk} \sigma_k^2 \quad (7)$$

Now, let's see what  $\tilde{N}_{ij}$  is equal to:

$$(\text{SN})_{ik} = S_{ik} \sigma_k^2$$

Now, since  $S^T$  is the final term in  $\tilde{N}$ , we can get the rotated noise elements:

$$\tilde{N}_{ij} = \sum_k (\text{SN})_{ik} S_{kj}^T$$

But  $S_{kj}^T$  is just  $S_{jk}$ , so we have

$$\tilde{N}_{ij} = \sum_k S_{ik} S_{jk} \sigma_k^2$$

This result is just Equation 7, so as long as we can calculate the  $\langle \tilde{n}_i \tilde{n}_j \rangle$  in our new coordinate system, we can calculate  $\chi^2$  without ever referring to the original, uncorrelated data! So, our expression for  $\chi^2$  (Equation 2) is true even if our noise is correlated as long as  $\tilde{N}_{ij} = \langle \tilde{n}_i \tilde{n}_j \rangle$ . Since that is true, the entire framework we have developed for solving least-squares problems carries through directly. While writing the equations down using linear algebra started out as a notational convenience, it has become a powerful tool letting us analyze data with correlated errors as easily as uncorrelated errors.

Calculus with Linear Algebra': It is likely that many people reading this note have taken linear algebra at some point but perhaps not have done calculus with it. To first order, doing calculus with linear algebra is just like doing regular calculus, as long as one is careful not to

switch the order of any matrices. For instance, the derivative of the linear least-squares equations  $\chi^2 = (d - Am)^T N^{-1} (d - Am)$  is just

$$\nabla \chi^2 = -A^T N^{-1} (d - Am) + (d - Am)^T N^{-1} A$$

. We've played a bit loose here since we've ignored the difference between taking the derivative of  $m$  and  $m^T$ , but remember that underneath we can think about this as taking the derivative of  $\chi^2$  with respect to the individual  $m_i$ 's, each of which is a scalar. Since the two terms in the sum are just conjugates of each other, the individual elements are the same and so we can just add their elements and pick one of the transposes to use. The usual choice is to then say

$$\nabla \chi^2 = -2A^T N^{-1} (d - Am)$$

.