

Maximum Likelihood

Nuno Carvalho

July 11, 2025

A lot of this text is adapted from the textbook [Lynch & Walsh, 1998], particularly Appendix 4, as well as some back-and-forth with ChatGPT to help consolidate ideas. The point of this is to re-explain (restricted) maximum likelihood in a way that makes more sense to me, and hopefully to you. Here, I focus more on the theory and abstraction of ML/REML, and in other entries, I plan to apply it to a particular class of problem (e.g. variance component estimation).

1 Maximum Likelihood (ML)

1.1 Likelihood

When we construct a model of some real world phenomenon, we define some process that describes how data, \mathbf{y} , is generated according to some set of parameters, Θ . Typically, our model is not deterministic since we do not understand every single cause of the real world phenomenon with perfect accuracy. Instead, our data is generated probabilistically, which means given some set of parameter values, different possible data values can be generated with differing probabilities. We can also think of the reverse: Given our observed data, there are many different combinations of parameter values that could have generated that data, each with differing likelihoods.

We can define a **likelihood function** that tells us the probability density of some data \mathbf{y} being generated by some set of parameter values Θ :

$$\ell(\Theta|\mathbf{y})$$

Note that because we actually do know what the data is (it is observed), we are asking the question: "How well does an input set of parameter values, Θ , explain the observed data, \mathbf{y} . This does not necessarily tell us the **probability** of the parameter values given our data, since we are not computing the likelihood across all possible parameters (weighted by some prior), as we would in Bayesian statistics.

Because our data is (usually) composed of multiple data points, and our model often describes the probability of a single data point being generated, our likelihood for the

entire data will often be the product of multiple individual likelihoods:

$$\ell(\boldsymbol{\Theta}|\mathbf{y}) = \prod_i^N \ell(\boldsymbol{\Theta}|y_i)$$

Note that this makes the crucial assumption that each data point is independent of each other, such that the likelihood of one data point is not affected by the values of the other datapoints. Because sums are typically easier to work with than products, it's useful to instead compute the natural logarithm of the likelihood, which per the properties of logarithms, turns those products into sums. This defines the **log-likelihood**:

$$L(\boldsymbol{\Theta}|\mathbf{y}) = \sum_i^N \ln(\ell(\boldsymbol{\Theta}|y_i))$$

Because logarithms are a monotonic function, it means that whatever parameter values maximize the likelihood also maximize the log-likelihood, and same for the minimum.

1.1.1 Example of the likelihood of a function

(Adapted from Example 1 of Appendix 4 of [Lynch & Walsh, 1998]). Let's say we have N data points: $\mathbf{y} = [y_1, y_2, \dots, y_N]$, which we believe were generated from some normal distribution with a variance of $\sigma^2 = 1$ but an unknown mean. We know the probability density function (PDF) for a normal distribution, which is a function of the mean μ and variance σ .

$$\Pr(y_i, \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{y_i - \mu}{\sigma}\right)^2\right]$$

In our case, we know $\sigma = 1$ and we also observe the data \mathbf{y} . What we don't know is the mean, μ , which is the only non-fixed parameter. We can plug in our known parameter σ , compute the product of the PDF for all our data points, and take the logarithm. This results in a log-likelihood function that describes how well an inputted μ value explains our observed data:

$$L(\mu|\mathbf{y}) = -\frac{n}{2} \ln(2\pi) - \frac{1}{2} \sum_i^N (y_i - \mu)^2$$

1.2 Maximum Likelihood Estimation

Once we have described the likelihood of some set of parameter values given our observed data, we can then ask the question: What set of parameter values maximizes the likelihood of our data? The rationale is that if we want to estimate these unknown parameters, our best bet for a single point estimate is the combination of parameter values that are most likely to produce the observed data. The set of parameter values $\hat{\boldsymbol{\Theta}}$ that maximizes the likelihood of our observed data is called the **maximum likelihood estimate (MLE)**.

From calculus, we know that the maximum of a multivariable function can be identified by taking the partial derivatives of the function with respect to each parameter and then

solving for the parameter values that make them equal to zero. We define this first partial derivative of a log-likelihood function to be the **Score** (\mathbf{S}) of the likelihood function:

$$\mathbf{S}(\boldsymbol{\Theta}) = \frac{\partial \mathbf{L}(\boldsymbol{\Theta})}{\partial \boldsymbol{\Theta}}$$

Note that if $\boldsymbol{\Theta}$ is a vector of M parameters, then $\mathbf{S}(\boldsymbol{\Theta})$ is also a vector of length M . The MLE, $\hat{\boldsymbol{\Theta}}$, is therefore the solution to the following:

$$\mathbf{S}(\hat{\boldsymbol{\Theta}}) = 0$$

1.2.1 Example of the MLE of a function

Continuing the previous example, we can compute the score of the likelihood function, which in this case is just a single derivative:

$$\begin{aligned} \mathbf{S}(\mu) &= \frac{\partial \mathbf{L}(\mu)}{\partial \mu} \\ &= \frac{\partial}{\partial \mu} \left[-\frac{n}{2} \ln(2\pi) - \frac{1}{2} \sum_i^N (y_i - \mu)^2 \right] \\ &= \frac{\partial}{\partial \mu} \left[-\frac{n}{2} \ln(2\pi) - \frac{1}{2} \sum_i^N (y_i^2 - 2y_i\mu + \mu^2) \right] \\ &= -\frac{1}{2} \sum_i^N (-2y_i + 2\mu) \\ &= N(\bar{y} - \mu) \end{aligned}$$

where $\bar{y} = \frac{1}{N} \sum_i^N (y_i)$ (i.e. the sample mean). We now set $\mathbf{S}(\hat{\mu}) = 0$ and solve:

$$\begin{aligned} 0 &= \mathbf{S}(\hat{\mu}) \\ &= N(\bar{y} - \hat{\mu}) \end{aligned}$$

The solution is: $\hat{\mu} = \bar{y} = \frac{1}{N} \sum_i^N (y_i)$. Matching our intuitions, the mean parameter that maximizes the likelihood of our data is the arithmetic mean of our data sample.

1.2.2 Asymptotic efficiency of the MLE

What's nice about ML estimation is that as the sample size increases, the MLE converges to the true parameter value θ_0 (if the likelihood function matches that of the true data-generating function). This is called **asymptotic efficiency** and can be written formally as:

$$\hat{\theta}_N \rightarrow^p \theta_0$$

This should make intuitive sense, as the true data-generating function has some true distribution of observed data that it produces, which in most cases, is unique to the specific parameter value of that data-generating function. Consequently, the parameter that maximizes the likelihood of producing the true parameter's distribution of data is the true parameter itself. The more samples we collect, the more our observed data's distribution should match the true distribution of the possible data. Therefore, the MLE should get closer to the true parameter value.

1.3 The Curvature of the Likelihood

While the first derivative describes the rate of change of the likelihood function, the second derivative describes its curvature. At the MLE, a strong curvature (high magnitude) means the likelihood peaks very sharply at the MLE. Therefore, a small change in the parameter value(s) from the MLE leads to a sharp drop in the likelihood. In contrast, if there is weak curvature (low magnitude) at the MLE, then there is a wide region of parameter estimates that yield similar (albeit slightly lower) likelihoods of producing the data.

1.3.1 The Hessian

For M parameters, the second derivative of the log-likelihood function is the **Hessian** (\mathbf{H}), an $M \times M$ matrix where:

$$H_{ij} = \frac{\partial^2 L(\Theta|\mathbf{y})}{\partial \Theta_i \partial \Theta_j}$$

This is also called the gradient and can be evaluated at any set of parameter values, describing the multi-dimensional curvature of the log-likelihood at that set of parameters.

1.3.2 The Fisher information Matrix

Of mathematical interest to us is the negative expected value of the Hessian, which is called the **Fisher information matrix**:

$$\mathcal{I}(\Theta) = -E[\mathbf{H}(\Theta)]$$

Note that the expected value here is taken by averaging across all the data that is potentially generated by the inputted set of parameters, weighted by the probability of generating that data. This means that while the Hessian is a function that depends on both the parameters and the data that is observed, the Fisher information matrix only depends on the parameters. That is why $-\mathbf{H}(\Theta)$ is often called the **observed information matrix**, while $\mathcal{I}(\Theta)$ is the **expected information matrix**.

As sample size increases, the Fisher information matrix converges to the negative of the Hessian evaluated at the MLE:

$$\mathcal{I}(\hat{\Theta}) \approx \mathbf{H}(\hat{\Theta})$$

1.3.3 Asymptotic normality of the MLE

If our likelihood function follows [some regularity conditions](#), as sample size increases, the sampling distribution of the MLE (that is, the distribution of the MLE for different samples of observed data from the same underlying data-generating process) converges to normality. This is called **asymptotic normality** and can be written formally as:

$$\sqrt{N}(\hat{\theta}_N - \theta_0) \rightarrow^d \mathcal{N}(0, \sigma^2)$$

It can be [proven](#) that the variance (σ^2 in the above equation) of the MLE converges to the inverse of the Fisher information at the true parameter value:

$$\hat{\theta}_N \rightarrow^d \mathcal{N}(\theta_0, \mathcal{I}(\theta_0)^{-1})$$

This also applies when there are multiple parameters, such that the inverse of the Fisher information matrix evaluated at the MLE is the covariance matrix of the MLE. The i^{th} diagonal element is the variance of the MLE for the i^{th} parameter. The element (i,j) is the covariance of the MLE for the i^{th} and j^{th} parameters.

2 Restricted Maximum Likelihood (REML)

% To do

3 Methods of solving for the MLE

In practice, solving for the set of parameters that makes the score of the likelihood function equal to zero is very difficult. Instead, there are algorithms that start with a guess for the parameter estimates, and then iteratively move closer and closer to the parameter values that maximize the likelihood of the observed data. Eventually, the algorithm gets close enough to the MLE that the parameter estimate at each iteration no longer changes by much, and thus we decide on a cutoff point where we declare that the algorithm has converged on the MLE.

Broadly, the strategy of these methods is to randomly choose (or have the user specify) a location in the parameter landscape, where the height of the landscape at that set of parameter values is the log-likelihood of those values generating the observed data. The algorithm can't look very far from its position with confidence, but it can look at how the landscape height (the log-likelihood function) changes at its location. This can be simply the rate of change (first derivative) or also the curvature (second derivative). It then uses this information to predict the direction and distance of where the the landscape height should reach its maximum (the MLE) and goes to it. In the first few iterations, it's likely that this new predicted MLE location turns out to not explain the data very well, since the landscape is not flat at that location. The algorithm repeats this jumping procedure until it is no longer moving very much between its prediction for how the landscape changes around the true peak has become pretty accurate.

Note, however, if there are local maxima (peaks in the likelihood landscape that are not the tallest in the entire landscape), the algorithm may fail to converge on the true parameter values, instead converging on said local maxima. By re-running the algorithm under multiple starting spots, one can usually avoid this.

3.1 Quadratic iterative methods

From calculus, we know that a function can be approximated around some point through a Taylor series expansion:

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2 + \dots + \frac{1}{n!}f^{(n)}(x_0)(x - x_0)^n + \dots$$

Through some algebra, if we set $f(x) = 0$ and solve for x using the first-order approximation, we get:

$$\begin{aligned} f(x) = 0 &\approx f(x_0) + f'(x_0)(x - x_0) \\ f'(x_0)(x - x_0) &\approx -f(x_0) \\ (x - x_0) &\approx -\frac{f(x_0)}{f'(x_0)} \\ x &\approx x_0 - \frac{f(x_0)}{f'(x_0)} \end{aligned}$$

We can replace the above function with the score, which is the first derivative of the log-likelihood: $f(x) = \mathbf{S}(\boldsymbol{\Theta}) = \frac{\partial L(\boldsymbol{\Theta})}{\partial \boldsymbol{\Theta}}$. And we can replace its derivative with the second derivative of the log-likelihood: $f'(x) = \frac{\partial \mathbf{S}(\boldsymbol{\Theta})}{\partial \boldsymbol{\Theta}} = \frac{\partial^2 L(\boldsymbol{\Theta})}{\partial^2 \boldsymbol{\Theta}}$. Furthermore, the x_0 in this equation becomes an arbitrary set of parameter values we are starting from, $\hat{\boldsymbol{\Theta}}^{(0)}$, while the x represents our approximation of the MLE from that starting point, since $f(x) = 0$ and the MLE is what makes $\mathbf{S}(\boldsymbol{\Theta}) = 0$. We can then use that resulting MLE, $\hat{\boldsymbol{\Theta}}^{(1)}$, as our starting point to again approximate where the MLE actually is, $\hat{\boldsymbol{\Theta}}^{(2)}$, repeating the process that can be generalized as:

$$\hat{\boldsymbol{\Theta}}^{(k+1)} = \hat{\boldsymbol{\Theta}}^{(k)} - \left(\frac{\partial^2 L(\boldsymbol{\Theta})}{\partial^2 \boldsymbol{\Theta}} \right) \Big|_{\hat{\boldsymbol{\Theta}}^{(k)}}^{-1} \frac{\partial L(\boldsymbol{\Theta})}{\partial \boldsymbol{\Theta}} \Big|_{\hat{\boldsymbol{\Theta}}^{(k)}}$$

We can rewrite this using the derivative functions of the likelihood we've already defined: the Score $\mathbf{S}(\boldsymbol{\Theta})$ and the Hessian $\mathbf{H}(\boldsymbol{\Theta})$:

$$\hat{\boldsymbol{\Theta}}^{(k+1)} = \hat{\boldsymbol{\Theta}}^{(k)} - \mathbf{H}^{-1}(\hat{\boldsymbol{\Theta}}^{(k)}) \mathbf{S}(\hat{\boldsymbol{\Theta}}^{(k)})$$

Connecting back to the earlier analogy, using the first and second derivatives of the likelihood landscape at the algorithm's current spot, it predicts where the landscape peak should be. This becomes the algorithm's MLE for its next iteration. However, once at this new estimate, if there is still curvature in the landscape, then we must not yet be at the true MLE of the observed. The algorithm iteratively does this, and in theory, after infinite iterations, the algorithm lands on the true MLE, where $\mathbf{S}(\boldsymbol{\Theta})^{(K)} = 0$ and therefore our iterative formula becomes:

$$\hat{\boldsymbol{\Theta}}^{(K+1)} = \hat{\boldsymbol{\Theta}}^{(K)}$$

The algorithm has converged. Of course, in practice, we don't have time to do an infinite number of iterations, so we decide on some small distance between subsequent estimates for which we declare that the algorithm has converged.

The following three quadratic methods discussed here differ in what they compute to use as the second derivative of the log-likelihood.

3.1.1 Newton-Raphson

The Newton-Raphson method uses the observed information matrix evaluated the current iteration's MLE as the second derivative of the log-likelihood. This is just the negative of the Hessian evaluated at the current MLE, so it looks identical to how we defined the iterative process earlier:

$$\hat{\boldsymbol{\theta}}^{(k+1)} = \hat{\boldsymbol{\theta}}^{(k)} - \mathbf{H}^{-1}(\hat{\boldsymbol{\theta}}^{(k)})\mathbf{S}(\hat{\boldsymbol{\theta}}^{(k)})$$

3.1.2 Fisher-Scoring

The Fisher-Scoring method uses the expected information matrix evaluated the current iteration's MLE as the second derivative of the log-likelihood. So, instead of the negative Hessian, it uses the Fisher information matrix. Noting the change in sign, the iterative formula looks like:

$$\begin{aligned}\hat{\boldsymbol{\theta}}^{(k+1)} &= \hat{\boldsymbol{\theta}}^{(k)} + \mathcal{I}^{-1}(\hat{\boldsymbol{\theta}}^{(k)})\mathbf{S}(\hat{\boldsymbol{\theta}}^{(k)}) \\ &= \hat{\boldsymbol{\theta}}^{(k)} - \mathbf{E}[\mathbf{H}^{-1}(\hat{\boldsymbol{\theta}}^{(k)})]\mathbf{S}(\hat{\boldsymbol{\theta}}^{(k)})\end{aligned}$$

The Fisher-Scoring method offers some advantages compared to the Newton-Raphson method. First, it is often easier to compute the Fisher information matrix than the Hessian. The Hessian will usually be a function of both the unknown parameters we are trying to estimate (which are set to the current iteration's MLE when evaluating it) as well as unbiased estimates of said parameters that are computed directly from the observed data. Since these estimates from the data are unbiased, on expectation across different data samples (generated from the current MLE), those estimates will exactly equal their respective parameters. This often leads to terms cancelling out, simplifying the equation significantly. Second, the Fisher-Scoring method also appears to be better at converging correctly from different starting points compared to Newton-Raphson (according to Jenrich and Sampson 1986).

3.1.3 Average Information

The Average Information method uses the average of the expected and observed information matrices evaluated the current iteration's MLE as the second derivative of the log-likelihood. Defining this matrix as \mathbf{AI} , it is:

$$\mathbf{AI}(\boldsymbol{\theta}) = \frac{1}{2}(\mathcal{I}(\boldsymbol{\theta}) - \mathbf{H}(\boldsymbol{\theta}))$$

Such that the iterative equation is:

$$\begin{aligned}\hat{\boldsymbol{\theta}}^{(k+1)} &= \hat{\boldsymbol{\theta}}^{(k)} + \mathbf{AI}^{-1}(\hat{\boldsymbol{\theta}}^{(k)})\mathbf{S}(\hat{\boldsymbol{\theta}}^{(k)}) \\ &= \hat{\boldsymbol{\theta}}^{(k)} + \left(\frac{\mathcal{I}(\hat{\boldsymbol{\theta}}^{(k)}) - \mathbf{H}(\hat{\boldsymbol{\theta}}^{(k)})}{2}\right)^{-1}\mathbf{S}(\hat{\boldsymbol{\theta}}^{(k)})\end{aligned}$$

In some scenarios, taking the average of the expected and observed information matrices leads to terms cancelling out, resulting in a simpler matrix to calculate. So, using the Average Information method can be faster and converge better than the Newton-Raphson and Fisher-Scoring methods.