Biophysics Lecture Tuesday, January 15th, 2019

Presenter: Borries Demeler

Topic: Survey of Numerical Modeling

Copy of Lecture at:

https://demeler.uleth.ca/biophysics/archive/Demeler/

Survey of Numerical Modeling

We will touch on the following subjects:

- Model Building
 - Exact models
 - smoothing
- Theory of Fitting
 - What is a good fit and how can it be measured?
- Parameter Estimation
- Optimization how do I fit the model to the data?
 - Linear vs. non-linear least squares
 - Linearization of non-linear systems
 - Brute force methods Grid searches
 - Stochastic Methods Genetic algorithms
- Effect of noise on analysis
- Statistical analysis and Monte Carlo approaches
- Parallel Implementation

Data Fitting

Modeling involves the description of some observable data (experimental measurements) using a mathematical equation that describes the underlying physical properties of the experiment.

First, we need to identify a general, mathematical model that can represent the observed data. The *parameters* of the model describe the specifics of the data.

Second, we need to determine the *values* of the parameters in the model that best fit our data. This is accomplished by a fitting algorithm that minimizes the difference between the data and model. Generally, an initial estimate is required that is then improved.

Finally, we need to estimate the error in the parameters we determined in the fitting process and obtain the *confidence intervals*.

Model Building

To build a model, one needs to understand the physical properties of the observed process. Many processes can be described by differential equations. When solved, these equations describe a linear or nonlinear model:

Example – radioactive decay:

Hypothesis: The rate of decay is proportional to the number of nuclei present.

Mathematical model: $\frac{\partial N}{\partial t} = -a N$

Solve: $N = N_0 e^{-a(t-t_0)} + b$

You start with some experimental data...



Absorption data from multiple concentrations fitted to a sum of Gaussian functions

You start with some experimental data...



Absorption data from multiple concentrations fitted to a sum of Gaussian functions

Method of Least Squares

Fitting Data to a Model by the Method of Least Squares:

Any observable process that influences the measurement needs to be accounted for in order for the model to yield meaningful results. The object is then to minimize the residuals between the model and the data:

$$MIN \sum_{i=1}^{n} \left(\frac{Data_i - (Model_i)}{uncertainty_i} \right)^2 = MIN \chi^2$$

Extracting parameters from a simulated solution by fitting the model to experimental data is called an *"Inverse Problem"*

A non-parametric fit is used to smooth data for display, where the intrinsic model is of little interest, and hence the parameters are not needed.

You start with some experimental data...



Absorption data from multiple concentrations fitted to a sum of Gaussian functions

Method of Least Squares

$$MIN \sum_{i=1}^{n} \left(\frac{Data_i - (Model_i)}{uncertainty_i} \right)^2 = MIN \chi^2$$

Assumptions made in the Method of Least Squares:

The model is a truthful representation of reality

- All error is associated with the dependent variable. We can scale the reliability of each observation with an uncertainty factor σ_i .
- All experimental noise is considered to be of Gaussian distribution

Experimental Uncertainties

The uncertainty of a measurement can be determined by repeating the experiment several times. Each time, a slightly different value is obtained for the experimental observation. Assuming a Gaussian distribution of errors in the measurement, one can determine the standard deviation σ of the distribution of measurement values, and use σ to set error bars on a measurement and to scale the contribution of a datapoint to the sum of the residuals.

The standard deviation can be calculated by using this formula:

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (x_i - \overline{x})^2}$$

where \overline{x} is the average of all measurements.



Fitting Basics

For a straight line we have: y = a + bx

The least squares equation Is given by:

\sum^{n}	$D_i - M_i$)2
$\sum_{i=1}$	σ_i)

The distances are measured perpendicular to the data.

The object is to find the equation of the straight line that minimizes the distance between the straight line and the data points.

$$\sum_{i=1}^{n} \left(\frac{D_i - a - b x_i}{\sigma_i} \right)^2 = \chi^2(a, b)$$



Fitting Basics

What is an error surface?

Each parameter combination a, b results in a unique error when fitted to the experimental data. The optimal solution occurs where the error is the smallest. Ideally, the error surface is continuously differentiable.



Error surface for some function y = F(a, b)

The minimum in the differences occur where the derivative of our objective function with respect to the parameters is zero, so we need to differentiate it with respect to the parameters of interest, a and b:

$$\sum_{i=1}^{n} \left(\frac{D_i - (a + b x_i)}{\sigma_i} \right)^2 = \chi^2(a, b)$$
$$\frac{\partial \chi^2}{\partial a} = 0 = -2 \sum_{i=1}^{n} \left(\frac{D_i - a - b x_i}{\sigma_i^2} \right)$$
$$\frac{\partial \chi^2}{\partial b} = 0 = -2 \sum_{i=1}^{n} \left(\frac{x_i (D_i - a - b x_i)}{\sigma_i^2} \right)$$

This leads to a system of linear equations:

$$\sum_{i=1}^{n} \frac{a}{\sigma_{i}^{2}} + \sum_{i=1}^{n} \frac{b x_{i}}{\sigma_{i}^{2}} = \sum_{i=1}^{n} \frac{D_{i}}{\sigma_{i}^{2}}$$
$$\sum_{i=1}^{n} \frac{a x_{i}}{\sigma_{i}^{2}} + \sum_{i=1}^{n} \frac{b x_{i}^{2}}{\sigma_{i}^{2}} = \sum_{i=1}^{n} \frac{D_{i} x_{i}}{\sigma_{i}^{2}}$$

Or, in matrix form:

$$\begin{bmatrix} \sum_{i=1}^{n} \frac{1}{\sigma_i^2} & \sum_{i=1}^{n} \frac{x_i}{\sigma_i^2} \\ \sum_{i=1}^{n} \frac{x_i}{\sigma_i^2} & \sum_{i=1}^{n} \frac{x_i^2}{\sigma_i^2} \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^{n} \frac{D_i}{\sigma_i^2} \\ \sum_{i=1}^{n} \frac{D_i x_i}{\sigma_i^2} \\ \sum_{i=1}^{n} \frac{D_i x_i}{\sigma_i^2} \end{bmatrix}$$

Let

$$A = \begin{bmatrix} \sum_{i=1}^{n} \frac{1}{\sigma_{i}^{2}} & \sum_{i=1}^{n} \frac{x_{i}}{\sigma_{i}^{2}} \\ \sum_{i=1}^{n} \frac{x_{i}}{\sigma_{i}^{2}} & \sum_{i=1}^{n} \frac{x_{i}^{2}}{\sigma_{i}^{2}} \end{bmatrix}, X = \begin{bmatrix} x_{1} \\ x_{2} \end{bmatrix}, \text{ and } B = \begin{bmatrix} \sum_{i=1}^{n} \frac{D_{i}}{\sigma_{i}^{2}} \\ \sum_{i=1}^{n} \frac{D_{i} x_{i}}{\sigma_{i}^{2}} \end{bmatrix}$$

(where $x_1 = a$ (i.e., intercept) and $x_2 = b$ (i.e., slope))

In matrix notation: AX = B, with solution $A^{-1}AX = A^{-1}B = X$

The equations can be solved either by inverting A or by using Cramer's Rule:

$$x_1 = \frac{b_1 a_{22} - b_2 a_{12}}{a_{11} a_{22} - a_{12} a_{21}}, \ x_2 = \frac{b_2 a_{11} - b_1 a_{21}}{a_{11} a_{22} - a_{12} a_{21}}$$

Goodness of fit

The quality of the fit is determined by the randomness of the residuals and the root mean square deviation (RMSD).

The randomness of the residuals can be measured by determining the runs. Runs (R) are the number of consecutive positive (p) or negative (n) residuals from the mean.



$$R_{T} = \frac{R - \overline{R}}{\sigma_{R}}$$
$$\overline{R} = \frac{2np}{n+p} + 1$$
$$\sigma_{R}^{2} = \frac{2np(2np - n - p)}{(n+p)^{2}(n+p-1)}$$

The R_{τ} value is a measure of randomness and can be compared to a normal table to find out the probability of the test being random

Linear Models

The equation of a straight line is considered a *linear equation*:

$$y = a + bx$$

This equation is *linear* in the coefficients that are fitted. The equation doesn't have to be that of a straight line to be considered linear.

 $y = a + bx + cx^{2} + dx^{3} + ex^{4} + ...$

As long as the coefficients are linear, the equation is considered a linear fitting equation, no matter how wildly nonlinear the terms of the independent variable are:

$$y = a + b(x - sin(x^3)) + c e^{-(4-3x)} + d ln(3x^4) + ...$$

In general, we can write for any linear equation:

$$y = a_1 X_1 + a_2 X_2 + a_3 X_3 + a_4 X_4 + ...$$

where X_i can be any nonlinear term.

Linear Models

Linearization of a nonlinear equation:

Turn $y = ae^{bx}$ into a linear function of the form:

y = a + bx

take log on both sides:

Fitting the log of y reduces the nonlinear equation to a linear equation, $y^* = a^* + bx$, where $y^* = \ln y$ and $a^* = \ln a$.

$$\ln y = \ln a + bx$$

Parameter Constraints

Sometimes, we may want to constrain the value of a parameter – for example, we don't want the amplitude of an exponential to turn negative during fitting:

$$y = a e^{bx} + c = e^{\ln(a) + bx} + c$$

By making the transformation to fitting the log of a number we can assure that the number itself will never be negative (negative amplitudes don't make sense in many physical models).

Nonlinear Regression

Why is it such a big deal if an equation is linear or nonlinear? It turns out that nonlinear functions need to be fitted using iterative approaches, while linear functions can be fitted in a single iteration, so it helps to have the objective function in a linear form. For nonlinear systems, a Jacobian is defined.

The idea: Iteratively improve the parameter estimates by following along the gradient of the error function in the direction of maximum "improvement". This requires knowledge of the partial derivatives for each parameter at each point in the experiment. We build the Jacobian matrix:

$$= \begin{pmatrix} \frac{\partial X}{\partial a_1} \end{pmatrix}_{x_1} \begin{pmatrix} \frac{\partial X}{\partial a_2} \end{pmatrix}_{x_1} \cdots \begin{pmatrix} \frac{\partial X}{\partial a_n} \end{pmatrix}_{x_1} \\ \begin{pmatrix} \frac{\partial X}{\partial a_1} \end{pmatrix}_{x_2} \begin{pmatrix} \frac{\partial X}{\partial a_2} \end{pmatrix}_{x_2} \cdots \begin{pmatrix} \frac{\partial X}{\partial a_n} \end{pmatrix}_{x_2} \\ \cdots \cdots \cdots \\ \begin{pmatrix} \frac{\partial X}{\partial a_1} \end{pmatrix}_{x_m} \begin{pmatrix} \frac{\partial X}{\partial a_2} \end{pmatrix}_{x_m} \cdots \begin{pmatrix} \frac{\partial X}{\partial a_n} \end{pmatrix}_{x_m} \end{pmatrix}$$

Nonlinear Regression

Our equation is: $J * R = \Delta y$, with R = a - g

where J is the Jacobian matrix, g is the current parameter estimate, a is the adjustment made to the parameter estimate in the current iteration, this is the value we need to find. Δy is the difference between the experimental data and the model

Solve for *a*: $J^{T}J R = J^{T} \Delta y$, substitute $J^{T} \Delta y = B$, Option 1: use inverse: $(J^{T}J)^{-1}J^{T}J R = R = (J^{T}J)^{-1}J^{T} \Delta y$ Option 2: $J^{T}J$ is positive definite, so use Cholesky decomposition: $J^{T}J = LL^{T}$ $L(L^{T}R) = B$, substitute $L^{T}R = Z$, to get LZ = B and solve for Z using forward substitution, then solve for R using backward substitution: $L^{T}R = Z$, then solve for a to get the adjustment for the parameter.

Iterate until converged.

Optimization Methods

Linear Optimization:

Straight line fits Generalized linear least squares - single iteration fitting of objective functions of the type:

$$y = a_0 + \sum_{i=1}^n a_i X_i, -\infty < a_i < +\infty$$

NNLS (non-negative constrained least squares):

$$y = a_0 + \sum_{i=1}^n a_i X_i, \ a_i \ge 0$$

Multidimensional spectrum analyses (HPC recommended) Non-parametric fits (B-splines, polynomial smoothing, etc)

Optimization Methods

Nonlinear Optimization using Gradient Descent Methods for functions of the type:

 $y = F(a_i, x_i)$

Levenberg-Marquardt (stable, robust, works well even if initial guesses are rather far away from optimum) Gauss-Newton methods Quasi-Newton (works well near optimum) Conjugate gradients Tangent approximation methods (derivatives are not required) Neural networks



Problem with nonlinear least squares optimization:

For multi-component systems, the nonlinear least squares fitting algorithm gets easily stuck in local minima and the solution depends on the starting points. Problem gets worse with more parameters (i.e., multiple components).

Optimization Methods

Stochastic Methods

Monte Carlo Simulated Annealing Random walk Genetic Algorithms

Optimization Methods

Comparison Stochastic vs Deterministic Fitting Methods:

Stochastic:

- Large search space possible
- Generally slow converging
- Excellent convergence properties if given enough time
- Compute-intensive
- Suitable for many parameters
- Good for ill-conditioned error surfaces
- Derivatives not needed

Deterministic:

- Small search space
- Suitable for a few parameters only
- Well-conditioned error surface
- Very fast converging
- Requires derivatives

Noise & Data Modeling Considerations

Fitting of noisy data prevents unique solutions – multiple solutions are possible.

We need to *minimize noise* when modeling data.

There are three ways to reduce or eliminate noise:

- 1. fit the noise
- 2. maintain an exceptionally well tuned instrument
- 3. design your experiment to optimize the quality of the data

There are two noise types:

1. Systematic noise: Signal comes from a systematic source that is not part of the parametric model (finger print on the lens of a camera) and is highly correlated with some feature of the experiment.

2. Stochastic (random noise): Noise is (hopefully) Gaussian in distribution and uncorrelated to any feature of the experiment

(1) can often be fitted and accounted for, which (2) must be minimized.

Different Types of Noise - Systematic

Time Invariant noise: Noise is different for each radial position, but the same offset for each scan, and hence time independent.



Different Types of Noise - Systematic

Radially Invariant noise: Noise is different for each scan, but each radial position is offset by the same amount throughout the scan



Different Types of Noise - Stochastic

Stochastic (random noise): Noise is different for each radial and time point and it is (hopefully) Gaussian in distribution:



Experimental and Simulated Data



Tale of 2 noisy vectors



Intensity vs. Absorbance

Optical system considerations

Intensity measurements record the intensity of light passing through one channel. Absorbance measurements record the intensity of light passing through one channel, then record the intensity of the light passing through the reference channel, and subtract it from the first channel. Each channel recording contains the (nearly) same amount of time invariant noise, but different amount of stochastic noise. Subtraction of time invariant noise will eliminate it:

$$scan_{1} = signal_{1} + N_{s1} + N_{ti}$$

$$scan_{2} = signal_{2} + N_{s2} + N_{ti}$$

$$scan_{1} - scan_{2} = signal_{1} - signal_{2} + N_{s1} + N_{s2}$$

$$N_{sl} + N_{s2} \approx N_{sl} \sqrt{2}$$





Factors that affect Accuracy – Time-invariant Noise



How can we deal with uncertainty? A Recipe for Optimal Resolution:

Our solution is affected by random noise, time-invariant noise, the available signal, and its ratio to the random noise.

Improve signal to noise by:

Obtain lots of high quality data

Exploit the entire dynamic range of the acquisition system

Reduce noise and only use a well-functioning instrument

Remove systematic noise from experimental data

- Replace fitting parameters with experimentally determined values from separate experiments
- Explore the parameter space with a grid method, then parsimoniously regularize solution with GA, and use Monte Carlo to explore confidence regions

Perform global fits for multiple experimental conditions to improve signal



Remember:

you cannot get reliable answers if you start with low quality input data!



The Monte Carlo method is a stochastic approach that can be used to identify the effect noise has on the reliability of determined parameters. With the Monte Carlo approach the statistical confidence limits of each measured parameter can be determined.

Recipe for Monte Carlo:

Obtain a best-fit solution from model function fit and confirm that the residuals are random and without systematic deviation

Generate new synthetic Gaussian noise with the same quality as was observed in the original experiment and add it to the best-fit solution

Re-fit the solution

Repeat (2-4) at least 100 times and collect all parameter values

Calculate statistics from Monte Carlo distribution for each parameter

Generation of synthetic noise:

```
Method 1 – use Bootstrapping:
```

Permute a percentage of residuals. Take any residual, positive or negative, and place it elsewhere in the data. Assumption: The likelihood of a residual's magnitude is the same anywhere in the dataset.

Method 2 – generate Gaussian noise:

Run a 5-7 point Gaussian kernel over the residual vector from the best fit and use the Box-Müller algorithm to generate random new noise with the same quality at that position based on the variance obtained from the 5point kernel.

Generation of synthetic Gaussian noise with Method 2:



Take the absolute value of the residual values from 5-7 points, and smooth them with a Gaussian kernel (most weight on the center point). The average residual value is fed into the Box – Müller function to generate a new random residual that has a Gaussian probability distribution with a mean of the average residual value.

If there is a lot of noise in some area of the scan, the new data will have noise locally equivalent to the original data.

Generation of Gaussian noise is preferred because the quality of the noise varies at different points in the cell and it also varies with absorbance (more absorbance = more noise).

Box – Müller function:

```
float box muller(float m, float s)
                                        /* normal random variate generator */
                                         /* mean m, standard deviation s
                                                                             */
{
    float x1, x2, w = 2.0;
    while (w >= 1.0)
    {
        x1 = 2.0 * ranf() - 1.0;
        x2 = 2.0 * ranf() - 1.0;
        w = x1 * x1 + x2 * x2;
    }
    w = sqrt((-2.0 * log(w)) / w);
    return( m + x1 * w * s);
}
float ranf()
{
    int N = 1;
    float temp = 0.0;
    temp = (((float) rand() / ((float) RAND MAX + 1) * N));
    return(temp);
}
```

Obtaining Confidence Intervals

Once a Monte Carlo analysis is completed, a frequency distribution of parameter values for each parameter is obtained, and statistics for the distribution can be calculated:



Regularization:

Occam's Razor: The simplest solution describing the data is the preferred solution. Instead of regularizing the solution and introducing infinitely many solutes with different probability, we want to *REDUCE* the solution space to find the solution that has the smallest number of possible solutes. For probabilities of solutes inspect the GA evolution profile which provides MC statistics.



Homework 1: For the dataset shown on the right calculate the equations for a and b for a straight line fit using Cramer's rule ($y = c_1 + c_2 x$). Assume a standard deviation of 1 for each measurement. Show your work. Compare your answer by fitting with a plotting program. Show results.

Х	У	σ
2	10.6	1
4	12.1	1
6	14.5	1
8	20.8	1
10	17.3	1
12	24.7	1
14	29.1	1

$$A = \begin{bmatrix} \sum_{i=1}^{n} \frac{1}{\sigma_i^2} & \sum_{i=1}^{n} \frac{x_i}{\sigma_i^2} \\ \sum_{i=1}^{n} \frac{x_i}{\sigma_i^2} & \sum_{i=1}^{n} \frac{x_i^2}{\sigma_i^2} \end{bmatrix}, X = \begin{bmatrix} c_1 \\ c_2 \end{bmatrix}, \text{ and } B = \begin{bmatrix} \sum_{i=1}^{n} \frac{D_i}{\sigma_i^2} \\ \sum_{i=1}^{n} \frac{D_i x_i}{\sigma_i^2} \end{bmatrix}$$

$$c_1 = \frac{b_1 a_{22} - b_2 a_{12}}{a_{11} a_{22} - a_{12} a_{21}}, \quad c_2 = \frac{b_2 a_{11} - b_1 a_{21}}{a_{11} a_{22} - a_{12} a_{21}}$$

Homework assignment: For each of the following situations, rank the appropriateness of each fitting approach from best to worst or not applicable: non-parametric, gradient descent, grid method, or stochastic, specify whether each model is linear or nonlinear in its fitting parameters, propose a function if non-parametric, and justify your answer for your ranking.

- 1. Smoothing a set of noisy data with a recognizable wave pattern
- 2. Fitting parameters a and b from function:

$$F(x) = a_0 + a_1 \log(ax) + a_2 \sin(b - x^2) + c$$

- 3. Fitting parameters a_i from function: $F(x) = \sum_{i=0}^{100} a_i e^{b_i (x^2 - c^2)}, \text{ where } b_i = \frac{i}{100}$
- 4. Fitting parameters *a_i* from function:

$$F(x) = a_0 + a_1 \log(bx) + a_2 \sin(a_3 - x^2) + c$$

5. Fitting parameters *a*, *b* and *c* from function:

$$F(x) = \log(ax) + \sin(bx^2) + c$$

6. Fitting parameters $a_1 - a_8$:

$$F(x) = a_1 \log(a_{2x}) + \frac{a_3 \sin(\sqrt{2} a_4 x^2)}{a_7 x^3} + a_5 e^{\left(a_6 \frac{x^2}{2}\right)} + a_8$$