Likelihood Fits

Craig Blocker Brandeis August 23, 2004

Outline

I. What is the question?
II. Likelihood Basics
III. Mathematical Properties
IV. Uncertainties on Parameters
V. Miscellaneous
VI. Goodness of Fit
VII. Comparison With Other Fit Methods

What is the Question?

Often in HEP, we make a series of measurements and wish to deduce the value of a fundamental parameter.

For example, we may measure the mass of many B \circledast J/y K_S decays and then wish to get the best estimate of the B mass.

Or we might measure the efficiency for detecting such events as a function of momentum and then wish to derive a functional form.

The question is: What is the best way to do this? (And what does "best" mean?)

Likelihood Method

P(X|**a**) • Probability of measuring X on a given event. **a** is a parameter or set of parameters on which P depends.

Suppose we make a series of measurements, yielding a set of X_i's. The likelihood function is defined as

$$L = \bigcup_{i=1}^{N} P(\mathbf{X}_i | \mathbf{a})$$

The value of **a** that maximizes *L* is known as the Maximum Likelihood Estimator (MLE) of **a**, which we will denote as \mathbf{a}^* .

Note that we often work with ln(*L*).

Example

Suppose we have N measurements of a variable x, which we believe is Gaussian, and wish to get the best estimate of the mean and width. [Note that this and other examples will be doable analytically - usually we must use numerical methods].

$$P(\mathbf{x} \mid \mathbf{m}, \mathbf{s}) = \frac{1}{\sqrt{2\mathbf{p}\mathbf{s}}} e^{-\frac{(\mathbf{x} - \mathbf{m})^2}{2\mathbf{s}^2}}$$
$$\ln L = \overset{N}{\overset{N}{\underset{i=1}{a}}} \ln P(\mathbf{x}_i) = -N \ln(\sqrt{2\mathbf{p}\mathbf{s}}) - \overset{N}{\overset{N}{\underset{i=1}{a}}} \frac{(\mathbf{x}_i - \mathbf{m})^2}{\mathbf{s}^2}$$

Maximizing with respect to mand s gives

$$\mathbf{m}^{*} = \frac{\mathbf{\dot{a}} \mathbf{x}_{i}}{N} = \overline{\mathbf{x}}$$

$$\mathbf{s}^{*^{2}} = \frac{\mathbf{\dot{a}} \left(\mathbf{x}_{i} - \mathbf{m}^{*} \right)^{2}}{N} = \overline{\mathbf{x}^{2}} - \overline{\mathbf{x}}^{2}$$

$$\overset{\mathbf{N}}{\longrightarrow} \text{Discussed later.}$$

Warning

Neither $L(\mathbf{a})$ nor $\ln(L)$ is a probability distribution for **a**.

A frequentist would say that such a statement is just nonsense, since parameters of nature have definite values.

A Bayesian would say that you can convert *L*(**a**) to a "probability" distribution in **a** by applying Bayes thereom, which includes a prior probability distribution for **a**.

Bayesian versus Frequentist statistics is a can of worms that I won't open further in this talk.

Bias, Consistency, and Efficiency

What does "best" mean?

We want estimator to be close to the true value.

Unbiased **P** $\overline{\mathbf{a}^*} = \mathbf{a}_0$

Consistent Þ Unbiased for large N

Efficient $\mathbf{P} \left(\mathbf{a}^* - \mathbf{a}_0 \right)^2$ is minimal for large N

Maximum Likelihood Estimators are NOT necessarily unbiased but are consistent and efficient for large N. This makes MLE's powerful and popular (although we must be aware that we may not be in the large N limit). 6

Bias Example

Consider again, the mean and width of a Gaussian.

$$\overline{\boldsymbol{m}}^{*} = \overline{\frac{1}{N} \mathbf{a} \mathbf{x}_{i}} = \frac{1}{N} \mathbf{a} \overline{\mathbf{x}_{i}}^{-} = \frac{1}{N} \mathbf{a} \mathbf{m} = \mathbf{m}$$

$$\overline{\boldsymbol{s}}^{*2} = \overline{\frac{1}{N} \mathbf{a} (\mathbf{x}_{i} - \mathbf{m})^{2}} = \frac{N}{N-1} \boldsymbol{s}_{0}^{2}$$

Note that the MLE of the mean is unbiased, but for the width is not (although it is consistent). However,

$$\tilde{\mathbf{s}}^2 = \frac{\dot{\mathbf{a}}\left(\mathbf{x}_i - \mathbf{m}^*\right)^2}{N-1}$$

is unbiased.

In this case, we could find the bias analytically - in most cases we must look for it numerically.

Bias Example 2

Bias can depend upon choice of parameter. Consider an exponential lifetime distribution. We can use either the averge lifetime **t** or the decay width **G** as the parameter.

$$\mathbf{P}(\mathbf{t}) = \frac{1}{\mathbf{t}} e^{-\frac{\mathbf{t}}{\mathbf{t}}} = \mathbf{G} e^{-\mathbf{G} \mathbf{t}}$$



Uncertainty on Parameters

Just as important as getting an estimate of a parameter is knowing the uncertainty of that estimate.

The maximum likelihood method also provides an estimate of the uncertainty.

For one parameter, L becomes Gaussian for large N. Thus,

$$\ln L @ \ln L^* + \frac{1}{2} \frac{\P^2 \ln L}{\P a^2} \bigg|_{a=a^*} (a - a^*)^2$$

$$P Da^2 \circ \overline{(a^* - a_0)^2} = -\frac{1}{\frac{\P^2 \ln L}{\|a^2\|_{a=a^*}}}$$
USUALLY Write this as $\alpha = \alpha^* + \Delta \alpha$

We usually write this as $\alpha = \alpha^* \pm \Delta \alpha$

If
$$a = a^* \pm Da$$
, $\ln L = \ln L^* - \frac{1}{2}$

Note that this is a statement about the probability of the measurements, not the probability of the true value.

Uncertainty Example



a

Asymmetric Uncertainties

Sometimes ln*L* may not be parabolic and there may be asymmetric uncertainties.



Note: the DlnL = 1/2 interval does NOT always give a 68% confidence interval (see counterexample in handout).

Correlations

If there are multiple parameters, things are more complicated due to possible correlations.



For example, a fit to the linear function y = a x + b will have correlations between parameters a and b

Numerically, \mathbf{a}^{\pm} are given by where $\mathbf{D}\mathbf{ln}(\mathbf{L}) = 1/2$ and $\mathbf{ln}(\mathbf{L})$ is maximized wrt other parameters

12

Covariance matrix V is given by $V_{ij} = (\mathbf{a}_i - \overline{\mathbf{a}_i})(\mathbf{a}_j - \overline{\mathbf{a}_j})$ V is equal to U⁻¹, where $U_{ij} = -\frac{\P^2 \ln L}{\P \mathbf{a}_i \P \mathbf{a}_j}$

Normalization

Sometimes, people will say they don't need to normalize their probability distributions. This is sometimes true.

For the Gaussian example, if we omitted the normalization factor of $1/\sqrt{2ps}$ we get the mean correct but not the width. In general, if the normalization depends on any of the parameters of interest, it must be included.

My advice is always normalize (and always check the normalization).

Extended Likelihood

Suppose we have a Gaussian mass distribution with a flat background and wish to determine the number of events in the Gaussian. $(M_{-}M_{0})^{2}$

$$\mathbf{P} = \mathbf{f}_{\mathrm{S}} \frac{1}{\sqrt{2\mathbf{p}}\mathbf{s}} e^{-\frac{(\mathbf{M} - \mathbf{M}_{0})}{2\mathbf{s}^{2}}} + (1 - \mathbf{f}_{\mathrm{S}}) \frac{1}{\mathbf{D}\mathbf{M}}$$

where f_S is the fraction of signal events and **D**M is the mass range of the fit.

We can fit for f_S and get Df_S . Nf_S is a good estimate of the number of events in the Gaussian, but NDf_S is not a good estimate of the variation on the number of signal events.

We can fix this by adding a Poisson term in the total number of events. This is called an Extended Likelihood fit. We could also use $\mathbf{DN}_{S}^{2} = \mathbf{N}^{2}\mathbf{Df}_{S}^{2} + \mathbf{f}_{S}^{2}\mathbf{DN}^{2} = \mathbf{N}^{2}\mathbf{Df}_{S}^{2} + \mathbf{Nf}_{S}^{2}$ 14

Extended Likelihood 2

Instead of f_s , we use **m** and **m**_G, the expected number of signal and background events. N is the observed total number of events.

$$L = \frac{e^{-(m_{3}+m_{BG})}(m_{g}+m_{BG})^{N}}{N!} \bigotimes_{i=1}^{N} \bigotimes_{\substack{e=1\\e=m_{g}}}^{e} \frac{m_{g}}{m_{g}} G(M_{i}|M_{0},s) + \frac{m_{BG}}{m_{g}} H_{BG} + \frac{m_{BG}}{m_{g}} \frac{1}{DM} \overset{u}{u}$$
$$= \frac{e^{-(m_{g}+m_{BG})}}{N!} \bigotimes_{i=1}^{N} \bigotimes_{e=1}^{N} G(M_{i}|M_{0},s) + m_{BG} \frac{1}{DM} \overset{u}{u}$$
$$\ln(L) = -m_{g} - m_{BG} - \ln(M!) + \underset{i}{a} \ln \underset{e=1}{e} M_{G} + \frac{m_{BG}}{DM} \overset{u}{u}$$

If you are not interested in the uncertainty on N_S (for example, your are measuring a lifetime and not a cross section), I recommend not doing an extended likelihood fit.

Constrained Fits

Suppose there is a parameter in the likelihood that is somewhat known from elsewhere. This information can be incorporated in the fit.

For example, we are fitting for the mass of a particle decay with resolution **s**. Suppose the Particle Data Book lists the mass as $M_0 \pm s_M$. We can incorporate this into the likelihood function as

$$L = \frac{1}{\sqrt{2\mathbf{p}}\mathbf{s}_{\mathrm{M}}} e^{-\frac{(\mathbf{M} - \mathbf{M}_{0})^{2}}{2\mathbf{s}_{\mathrm{M}}^{2}}} \bigvee_{i} \frac{\mathbf{\hat{e}}}{\mathbf{\hat{e}}} \frac{1}{2\mathbf{p}\mathbf{s}} e^{-\frac{(\mathbf{m}_{i} - \mathbf{M})^{2}}{2\mathbf{s}^{2}}} \underbrace{\mathbf{\hat{u}}}{\mathbf{\hat{u}}}$$

This is known as a constrained fit.

Constrained Fits 2

Let **D**M be the uncertainty on M that could be determined by the fit alone.

If **D**M >> \mathbf{s}_{M} , constraint will dominate, and you might as well just fix M to M₀. For example, you never see a constrained fit to \hbar in an HEP experiment.

If $\mathbf{D}\mathbf{M} \ll \mathbf{s}_{\mathbf{M}}$, constraint does very little. You have a better measurement than the PDG. You should do an unconstrained fit and PUBLISH.

Constrained fit is most useful if \mathbf{s}_{M} and **DM** are comparable.

Simple Monte Carlo Tests

It is possible to write simple, short, fast Monte Carlo programs that generate data for fitting. Can then look at fit values, uncertainties, and pulls. These are often called "toy" Monte Carlos to differentiate them from complicated event and detector simulation programs.

- ★ Tests likelihood function.
- \star Tests for bias.
- **★** Tests that uncertainty from fit is correct.

This does NOT test the correctness of the model of the data. For example, if you think that some data is Gaussian distributed, but it is really Lorentzian, then the simple Monte Carlo test will not reveal this.

Simple Monte Carlo Tests 2

Generate exponential ($\mathbf{t} = 0.5$ and N = 1000).

Fit.

Repeat many times (1000 times here). Histrogram \mathbf{t} , \mathbf{s}_t , and pulls.



Simple Monte Carlo Tests 3



Goodness of Fit

Unfortunately, the likelihood method does not, in general, provide a measure of the goodness of fit (as a c^2 fit does).

For example, consider fitting lifetime data to an exponential.

$$L = \bigotimes_{i}^{N} \frac{1}{t} e^{-\frac{t_{i}}{t}}$$
$$t^{*} = \frac{\mathbf{\dot{a}} t_{i}}{N}$$
$$L(t^{*}) = -N \bigotimes_{k}^{\infty} 1 + \ln \frac{\mathbf{\dot{a}} t_{i}}{N} \overset{\mathbf{\ddot{o}}}{\vdots}$$

Thus the value of *L* at the maximum depends only on the number of events and average value of the data.

Goodness of Fit 2

Fit to exponential

Plot log(L*) for

(1) exponential Monte Carlo and

(2) Gaussian data



Goodness of Fit 3



Other Types of Fits

Chi-square:

If data is binned and uncertainties are Gaussian, then maximum likelihood is equivalent to a c^2 fit.

Binned Likelihood:

If data is binned and not Gaussian, can still do a binned likelihood fit. Common case is when data are Poisson distributed.

$$P_{i} = \frac{e^{-m} (m)^{n_{i}}}{n_{i}!}$$
$$\ln L = a \ln P_{i}$$
bins

Chi-square:

Goodness of fit.

Can plot function with binned data.

Data should be Gaussian, in particular, **c**² doesn't work well with bins with a small number of events.

Binned likelihood:

Goodness of fit

Can plot function with binned data.

Still need to be careful of bins with small number of events

(don't add in too many zero bins).

Unbinned likelihood:

Usually most powerful.

Don't need to bin data.

Works well for multi-dimensional data.

No goodness of fit estimate.

Can't plot fit with data (unless you bin data).

Generate 100 values for Gaussian with m=0, s=1. Fit unbinned likelihood and c^2 to SAME data. Repeat 10,000 times.







Numerical Methods

Even slight complications to the probability make analytic methods intractable.

Also, likelihood fits often have many parameters (perhaps scores) and can't be done analytically.

However, numerical methods are still very effective.

MINUIT is a powerful program from CERN for doing maximum likelihood fits (see references in handout).

Systematic Uncertainties

When fitting for one parameter, there often are other parameters that are imperfectly known.

It is tempting to estimate the systematic uncertainty due to these parameters by varying them and redoing the fit.

Because of statistical variations, this overestimates the systematic uncertainty (often called double counting).

Best way to estimate such systematics is probably with a high statistics Monte Carlo program.

Potentially Interesting Web Sites

CDF Statistics Committee page: www-cdf.fnal.gov/physics/statistics/statistics_home.html

Lectures by Louis Lyons: www-ppd.fnal.gov/EPPOffice-w/Academic_Lectures

Summary

Maximum Likelihood methods are a powerful tool for extracting measured parameters from data.

However, it is important to understand their proper use and avoid potential problems.