

Lecture 9

MAXIMUM-LIKELIHOOD ESTIMATION AND FISHER INFORMATION

Harrison H. Barrett
University of Arizona

OUTLINE

- Estimation as mapping
- Performance metrics: Bias and variance, MSE and EMSE
- Fisher information and Cramér-Rao bounds
- Maximum-likelihood estimation: Why and how?
- Textbook examples (univariate and bivariate)
- A multivariate toolkit
- The EM algorithm

Next lecture: Applications in astronomy

References

Barrett and Myers, Foundations of Image Science

For this lecture:

Chap. 13, Statistical Decision Theory

For background:

Appendix C, Probability

Chap. 8, Stochastic Properties of Objects and Images

Notation and terminology

The $M \times 1$ vector \mathbf{g} describes the random data.

The PDF on \mathbf{g} is characterized by the $K \times 1$ parameter vector $\boldsymbol{\theta}$ and denoted $\text{pr}(\mathbf{g}|\boldsymbol{\theta})$. This PDF describes the *sampling distribution* of \mathbf{g} for a given $\boldsymbol{\theta}$.

Once a data vector is measured, $\text{pr}(\mathbf{g}|\boldsymbol{\theta})$ can be regarded as a function of $\boldsymbol{\theta}$, sometimes called the *likelihood*:

$$L(\boldsymbol{\theta}|\mathbf{g}) = \text{pr}(\mathbf{g}|\boldsymbol{\theta})$$

Note that $L(\boldsymbol{\theta}|\mathbf{g})$ is *not* a PDF on $\boldsymbol{\theta}$.

An *estimate* of the parameter is denoted $\hat{\boldsymbol{\theta}}$; in most cases the estimate is a deterministic function of the data, so we can also write it as $\hat{\boldsymbol{\theta}}(\mathbf{g})$. Since \mathbf{g} is random (even for a given $\boldsymbol{\theta}$), so is $\hat{\boldsymbol{\theta}}(\mathbf{g})$.

Estimation as mapping

$$\theta \longrightarrow g \longrightarrow \hat{\theta}$$

$$\text{pr}(\theta) \longrightarrow \text{pr}(g|\theta) \longrightarrow \text{pr}(\hat{\theta}|\theta)$$

Performance metrics

For a given true value of θ , there are two basic kinds of error:

- Random error

Also called precision, estimation variance, noise, jitter, ...

- Systematic error

Also called accuracy, estimation bias, calibration error, ...

Both kinds can be quantified by the conditional distribution of the estimates:

$$\text{pr}(\hat{\theta}|\theta)$$

Bias of the estimate of a scalar parameter

Consider a scalar parameter θ and an estimate $\hat{\theta}(\mathbf{g})$. The conditional mean of the estimate can be computed from the likelihood:

$$\overline{\hat{\theta}} = \left\langle \hat{\theta}(\mathbf{g}) \right\rangle_{\mathbf{g}|\theta} = \int d^M g \, \text{pr}(\mathbf{g}|\theta) \hat{\theta}(\mathbf{g}), \quad (13.274)$$

where the subscript $\mathbf{g}|\theta$ indicates that the average is over randomness in the data when the parameter has fixed value equal to θ .

If we know the conditional PDF of $\hat{\theta}$ itself, we can write

$$\overline{\hat{\theta}} = \int d\hat{\theta} \, \text{pr}(\hat{\theta}|\theta) \hat{\theta}, \quad (13.275)$$

We denote the conditional bias by $b(\theta)$, where

$$b(\theta) = \overline{\hat{\theta}} - \theta. \quad (13.276)$$

The overline indicates an average over all realizations of the data, given the true value θ . An *unbiased* estimate is one for which $b(\theta) = 0$ for all θ .

Estimability

A parameter is said to be *estimable* or *identifiable* with respect to some data set if there exists an unbiased estimator of it for all true values of the parameter.

Examples:

- Stopped clock
- Integral of grey levels over ROI
- Pixel value

One way of dealing with the issue of estimability is to average the bias over all possible true values of the parameter, defining an average bias by

$$\bar{\bar{\hat{\theta}}} = \int d\theta \, \text{pr}(\theta) \int d^M g \, \text{pr}(g|\theta) \hat{\theta}(g), \quad (13.277)$$

Variance and MSE of the estimate of a scalar parameter

Bias specifies accuracy or systematic error. Variance specifies precision or random error. It is defined by:

$$\text{Var}(\hat{\theta}) = \sigma_{\hat{\theta}}^2 = \left\langle |\hat{\theta}(\mathbf{g}) - \bar{\hat{\theta}}|^2 \right\rangle_{\mathbf{g}|\theta} = \int d^M g \, \text{pr}(\mathbf{g}|\theta) |\hat{\theta} - \bar{\hat{\theta}}|^2. \quad (13.279)$$

The variance describes the fluctuations of the estimate about the mean of the estimate, *not the true value*.

Fluctuations around the true value are described by the *mean-square error*:

$$\text{MSE}(\theta) = \left\langle |\hat{\theta} - \theta|^2 \right\rangle_{\mathbf{g}|\theta}. \quad (13.280)$$

The *ensemble mean-square error* or *EMSE* is found by taking an additional average over the parameter:

$$\text{EMSE} = \left\langle \left\langle |\hat{\theta} - \theta|^2 \right\rangle_{\mathbf{g}|\theta} \right\rangle_{\theta}. \quad (13.281)$$

Vector parameters – bias

A P -dimensional parameter vector θ has an estimate $\hat{\theta}$ with mean $\langle \hat{\theta} \rangle$ given by

$$\bar{\hat{\theta}}(\mathbf{g}) = \int d^M g \, \text{pr}(\mathbf{g}|\theta) \hat{\theta}(\mathbf{g}) = \int d^P \hat{\theta} \, \text{pr}(\hat{\theta}|\theta) \hat{\theta}. \quad (13.282)$$

The bias $\mathbf{b}(\theta)$ is now a vector quantity:

$$\mathbf{b}(\theta) \equiv \bar{\hat{\theta}} - \theta \equiv \int_{\infty} d^M g \, [\hat{\theta}(\mathbf{g}) - \theta] \text{pr}(\mathbf{g}|\theta) = \int_{\infty} d^P \hat{\theta} \, [\hat{\theta} - \theta] \text{pr}(\hat{\theta}|\theta). \quad (13.283)$$

The average bias is now written $\bar{\mathbf{b}} = \langle \mathbf{b}(\theta) \rangle_{\theta}$.

Vector parameters – variance and covariance

If we denote the mean of the p^{th} element of the random vector $\hat{\boldsymbol{\theta}}$ by $\langle \hat{\theta} \rangle_p = \bar{\hat{\theta}}_p$, the variance of the p^{th} element is given by

$$\text{Var}(\hat{\theta}_p) \equiv \left\langle \left[\hat{\theta}_p - \langle \hat{\theta}_p \rangle \right] \left[\hat{\theta}_p - \langle \hat{\theta}_p \rangle \right]^* \right\rangle_{\mathbf{g}|\boldsymbol{\theta}}$$

$$= \int_{\infty} d^M g \left| \hat{\theta}_p(\mathbf{g}) - \langle \hat{\theta}_p(\mathbf{g}) \rangle \right|^2 \text{pr}(\mathbf{g}|\boldsymbol{\theta}) = \int_{\infty} d^P \theta \left| \hat{\theta}_p - \langle \hat{\theta}_p \rangle \right|^2 \text{pr}(\hat{\boldsymbol{\theta}}|\boldsymbol{\theta}), \quad (13.284)$$

and the full covariance matrix is given by

$$\mathbf{K}_{\hat{\boldsymbol{\theta}}} = \left\langle \left(\hat{\boldsymbol{\theta}} - \bar{\hat{\boldsymbol{\theta}}} \right) \left(\hat{\boldsymbol{\theta}} - \bar{\hat{\boldsymbol{\theta}}} \right)^{\dagger} \right\rangle = \langle \Delta \hat{\boldsymbol{\theta}} \Delta \hat{\boldsymbol{\theta}}^{\dagger} \rangle. \quad (13.285)$$

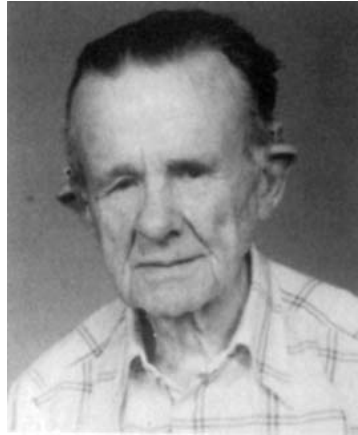
The MSE in the vector case is:

$$\text{MSE} = \left\langle ||\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}||^2 \right\rangle_{\mathbf{g}|\boldsymbol{\theta}} = \int_{\infty} d^M g \left| |\hat{\boldsymbol{\theta}}(\mathbf{g}) - \boldsymbol{\theta}| \right|^2 \text{pr}(\mathbf{g}|\boldsymbol{\theta}) = \text{tr} \left[\mathbf{K}_{\hat{\boldsymbol{\theta}}} \right] + \text{tr} \left[\mathbf{b} \mathbf{b}^{\dagger} \right]. \quad (13.287)$$

Key point:

There are fundamental limits to the variance and covariance of any estimator.

Carl Harald Cramér



Born: 25 Sept 1893 in Stockholm, Sweden

Died: 5 Oct 1985 in Stockholm, Sweden

Seminal book: Mathematical Methods of Statistics (Uppsala, 1945)

Cramér-Rao bound for a scalar random variable

The variance of any unbiased estimate of a scalar must satisfy

$$\text{Var}\{\hat{\theta}\} \geq \frac{1}{\left\langle \left[\frac{\partial}{\partial \theta} \ln \text{pr}(\mathbf{g}|\theta) \right]^2 \right\rangle_{\mathbf{g}|\theta}}, \quad (13.372)$$

where the denominator is called the *Fisher information*.

Similarly, for a biased estimator,

$$\text{Var}\{\hat{\theta}\} \geq \frac{\left(\frac{db(\theta)}{d\theta} + 1 \right)^2}{\left\langle \left[\frac{\partial}{\partial \theta} \ln \text{pr}(\mathbf{g}|\theta) \right]^2 \right\rangle_{\mathbf{g}|\theta}}. \quad (13.377)$$

This result tells us that a stopped clock can have zero variance (infinite precision)!

Fisher information in the vector case

For estimation of a scalar parameter, the Fisher information F is defined by:

$$F = \left\langle \left[\frac{\partial}{\partial \theta} \ln \text{pr}(\mathbf{g}|\theta) \right]^2 \right\rangle_{\mathbf{g}|\theta} .$$

For a vector parameter with P components, the Fisher information \mathbf{F} is a $P \times P$ Hermitian matrix with components:

$$\begin{aligned} F_{jk} &= \left\langle \left[\frac{\partial}{\partial \theta_j} \ln \text{pr}(\mathbf{g}|\boldsymbol{\theta}) \right] \left[\frac{\partial}{\partial \theta_k} \ln \text{pr}(\mathbf{g}|\boldsymbol{\theta}) \right] \right\rangle_{\mathbf{g}|\boldsymbol{\theta}} \\ &= \int_{\infty} d^M g \text{pr}(\mathbf{g}|\boldsymbol{\theta}) \left[\frac{1}{\text{pr}(\mathbf{g}|\boldsymbol{\theta})} \frac{\partial}{\partial \theta_j} \text{pr}(\mathbf{g}|\boldsymbol{\theta}) \right] \left[\frac{1}{\text{pr}(\mathbf{g}|\boldsymbol{\theta})} \frac{\partial}{\partial \theta_k} \text{pr}(\mathbf{g}|\boldsymbol{\theta}) \right] . \end{aligned} \quad (13.361)$$

CR bound in the vector case

The nn component of the covariance matrix of any random vector is the variance of the n^{th} component. For any unbiased estimate,

$$[\mathbf{K}_{\hat{\theta}}]_{nn} = \text{Var}\{\hat{\theta}_n\} \geq [\mathbf{F}^{-1}]_{nn}. \quad (13.371)$$

Note that inversion of the Fisher information is required to find the lower bound on the variance of a component of the estimate.

Other forms of the CR bound set limits on the covariance matrix of the estimate.

Terminology: An unbiased estimator that meets the CR bound is *efficient*.

Maximum-likelihood estimation

So far we have not talked about ways of actually finding an estimate. One general method is *maximum-likelihood estimation*:

$$\hat{\theta}_{\text{ML}} \equiv \underset{\theta}{\operatorname{argmax}} \operatorname{pr}(\mathbf{g}|\theta). \quad (13.348)$$

This procedure can be written equivalently as

$$\hat{\theta}_{\text{ML}} = \underset{\theta}{\operatorname{argmax}} \ln[\operatorname{pr}(\mathbf{g}|\theta)]. \quad (13.349)$$

Note that we are *not* maximizing the probability of θ ; we are choosing the value of θ that maximizes the probability of occurrence of the \mathbf{g} that we actually observed.

What's the score?

If there are no constraints on the estimate, $\hat{\theta}_{\text{ML}}$ occurs at a point where the gradient (with respect to θ) of the log-likelihood is zero:

$$\nabla_{\theta} \ln[\text{pr}(g|\theta)] = 0 \text{ at } \theta = \hat{\theta}_{\text{ML}}$$

The gradient of the log-likelihood is a random vector called the *score*:

$$s(g) = \nabla_{\theta} \ln[\text{pr}(g|\theta)] . \quad (13.358)$$

Basically, the score tells us how sensitive the log-likelihood is to changes in the parameters.

The Fisher information matrix is the covariance matrix of the score.

Why ML?

An ML estimate is:

- Efficient if an efficient estimate exists
- Asymptotically efficient (as you get more or better data)
- Asymptotically unbiased
- Asymptotically consistent
- Usually easy to compute
- A way of rigorously enforcing agreement with the data
- A way of doing estimation with no prior information

Why not ML?

ML estimation is:

- A way of rigorously enforcing agreement with the data
- A way of doing estimation with no prior information

Data are noisy. Rigorous agreement with noisy data will give noisy estimates – *even though that is the best you can do without bias!*

You always have some prior information – and you should use it, even though it might introduce bias.

One way to use a prior $\text{pr}(\theta)$ is with the *weighted likelihood*:

$$\hat{\theta}_{\text{WL}} \equiv \underset{\theta}{\operatorname{argmax}} \text{pr}(\mathbf{g}|\theta) \text{pr}(\theta) = \underset{\theta}{\operatorname{argmax}} \text{pr}(\theta|\mathbf{g}).$$

This estimate is also called the *maximum a posteriori* or MAP estimate, but that is a subject for another lecture.

TEXTBOOK EXAMPLES OF ML ESTIMATION

- Probability of heads in coin flipping
- Rate of a Poisson random process
- Mean and variance of a normal PDF

Example 1: Flipping a coin

How do you know a coin is “fair” (probability of heads = $\frac{1}{2}$)? Flip it many times and see how many heads you get.

Suppose you flip a coin 10 times and get 7 heads. What is the ML estimate of the probability of heads?

Need the Bernoulli or binomial probability law. If the probability of heads in a single flip is p , the probability of n heads in N flips is

$$\Pr(n|p) = \binom{N}{n} p^n (1 - p)^{N-n}, \quad (\text{C.161})$$

where the binomial coefficient is given by

$$\binom{N}{n} \equiv \frac{N!}{n! (N - n)!}. \quad (\text{C.162})$$

The mean and variance of a binomial random variable are given by

$$\langle n \rangle = Np \quad \text{and} \quad \sigma^2 = Np(1 - p). \quad (\text{C.163})$$

Example 1 – Computing the ML estimate

Given an observed number of heads n , maximize $\Pr(n|p)$ (or its logarithm) with respect to p to find the ML estimate of p .

$$\Pr(n|p) = \binom{N}{n} p^n (1 - p)^{N-n}, \quad (\text{C.161})$$

$$\ln \Pr(n|p) = \text{constant} + n \ln p + (N - n) \ln(1 - p),$$

$$\frac{d}{dp} \ln \Pr(n|p) = \frac{n}{p} - \frac{N - n}{1 - p} = 0 \text{ at } p = \hat{p}_{ML},$$

Thus

$$\hat{p}_{ML} = \frac{n}{N} = \frac{\text{Number of heads observed}}{\text{Number of flips}}$$

Example 1 – Bias, variance and efficiency

$$\hat{p}_{ML} = \frac{n}{N} = \frac{\text{Number of heads observed}}{\text{Number of flips}}$$

Recall:

$$\langle n \rangle = Np \quad \text{and} \quad \sigma^2 = Np(1 - p). \quad (\text{C.163})$$

Thus the bias and variance of the ML estimate are:

$$\langle \hat{p}_{ML} \rangle = \frac{\langle n \rangle}{N} = p, \quad \text{Var}(\hat{p}_{ML}) = \frac{\text{Var } n}{N^2} = \frac{p(1 - p)}{N}$$

Is the estimate efficient?

$$F = \left\langle \left[\frac{n}{p} - \frac{N - n}{1 - p} \right]^2 \right\rangle = \frac{1}{p^2(1 - p)^2} \langle [n - Np]^2 \rangle = \frac{\text{Var } n}{p^2(1 - p)^2} = \frac{N}{p(1 - p)}$$

Thus

$$\text{Var}(\hat{p}_{ML}) = F^{-1}$$

and the estimate is efficient.

Example 1 – Repeated trials

Now suppose you repeat the experiment J times, observing n_j heads out of N flips in the j^{th} repetition. Denote the set of observations by the vector \mathbf{n} with component $\{n_j, j = 1, \dots, J\}$. Since the repetitions are independent, the likelihood function is

$$\Pr(\mathbf{n}|p) = \prod_{j=1}^J \Pr(n_j|p) = \prod_{j=1}^J \binom{N}{n_j} p^{n_j} (1-p)^{N-n_j},$$

and the log-likelihood is

$$\ln \Pr(\mathbf{n}|p) = \sum_{j=1}^J [n_j \ln p + (N - n_j) \ln(1 - p)] + \text{constant},$$

There is still just one parameter (p) to estimate.

$$\frac{d}{dp} \sum_{j=1}^J [n_j \ln p + (N - n_j) \ln(1 - p)] = 0 \text{ at } p = \hat{p}_{ML},$$

from which

$$\hat{p}_{ML} = \frac{\sum_j n_j}{NJ} = \frac{\text{Total number of heads observed}}{\text{Total number of flips}}$$

Thus one long run of NJ flips gives the same ML estimate (and the same Fisher information) as J short runs of N flips each. The total number of heads, $\sum_j n_j$, is a *sufficient statistic* for estimation of p .

Example 2: Estimating the rate of a Poisson random process

Consider photons that arrive at an ideal counting detector at an average rate of a photons/sec. If the photons are independent, the random number of counts observed in time T has a probability law given by

$$\Pr(n|a) = \exp(-aT) \frac{(aT)^n}{n!}.$$

Suppose we observe a particular value of n in one trial of duration T and want to use this single observation to estimate the rate a . The log-likelihood is

$$\ln \Pr(n|a) = -aT + n \ln aT - \ln n!.$$

Thus

$$\frac{d}{da} \ln \Pr(n|a) = -T + \frac{n}{a} = 0 \text{ at } a = \hat{a}_{ML} = \frac{n}{T}.$$

Example 2 – bias, variance and efficiency

We know:

$$\hat{a}_{ML} = \frac{n}{T}, \quad \bar{n} = \text{Var}(n) = aT$$

Thus

$$\langle \hat{a}_{ML} \rangle = \frac{\bar{n}}{T} = a,$$

so the ML estimate is unbiased. Is it efficient? The variance is

$$\text{Var} [\hat{a}_{ML}] = \frac{\text{Var}(n)}{T^2} = \frac{a}{T},$$

and the Fisher information is

$$F = \frac{1}{a^2} \langle [n - aT]^2 \rangle = \frac{\text{Var}(n)}{a^2} = \frac{T}{a}.$$

Thus $\text{Var} [\hat{a}_{ML}] = F^{-1}$ and the estimate is efficient.

Example 3 – Independent samples of a normal RV

Let x be a univariate normal random variable, so that

$$\text{pr}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{1}{2\sigma^2} (x - \bar{x})^2 \right]$$

and let $\{x_j, j = 1, \dots, J\}$ be samples drawn independently from this distribution. The samples are said to be i.i.d. (independent and identically distributed).

The PDF for the whole set of samples, represented by the $J \times 1$ vector \mathbf{x} , is

$$\text{pr}(\mathbf{x}) = (2\pi\sigma^2)^{-J/2} \prod_{j=1}^J \exp \left[-\frac{1}{2\sigma^2} (x_j - \bar{x})^2 \right] .$$

The log of this PDF is

$$\ln \text{pr}(\mathbf{x}) = -\frac{J}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{j=1}^J (x_j - \bar{x})^2 .$$

Example 3 – estimation of the mean

Suppose we want to estimate the mean \bar{x} , assuming for the moment that the variance is known. Then the PDF on the last slide gets reinterpreted as the likelihood $\text{pr}(\mathbf{x}|\bar{x})$ for estimation of the scalar \bar{x} . Differentiating the log-likelihood yields:

$$\frac{d}{d\bar{x}} \ln \text{pr}(\mathbf{x}|\bar{x}) = \frac{1}{\sigma^2} \sum_{j=1}^J (x_j - \bar{x}) = 0,$$

or

$$\hat{\bar{x}}_{ML} = \frac{1}{J} \sum_{j=1}^J x_j = \text{arithmetic average of the data},$$

Efficient, unbiased, all that good stuff.

Example 3 – estimation of the mean and variance

Suppose now we want to estimate both \bar{x} and σ^2 . The same PDF as before is now the likelihood $\text{pr}(\mathbf{x}|\bar{x}, \sigma^2)$ for the 2D estimation problem. Two derivatives must vanish simultaneously:

$$\frac{\partial}{\partial \bar{x}} \ln \text{pr}(\mathbf{x}|\bar{x}, \sigma^2) = 0, \quad \frac{\partial}{\partial \sigma^2} \ln \text{pr}(\mathbf{x}|\bar{x}, \sigma^2) = 0.$$

Two equations in two unknowns. Algebra is messy, but the result is

$$\hat{x}_{ML} = \frac{1}{J} \sum_{j=1}^J x_j, \quad \hat{\sigma}_{ML}^2 = \frac{1}{J} \sum_{j=1}^J [x_j - \hat{x}_{ML}]^2.$$

Now the ML estimator is *biased*. It is shown in elementary statistics texts that an unbiased estimate of the variance is given by

$$\hat{\sigma}_{unbiased}^2 = \frac{1}{J-1} \sum_{j=1}^J [x_j - \hat{x}_{ML}]^2.$$

Barrett's Theorem

The maximum-likelihood procedure in any problem is what you are most likely to do if you don't know any statistics.

A multivariate tool kit

Before proceeding to real-world examples of ML estimation, we shall derive some general expressions for log-likelihoods, scores and Fisher information matrices. Problems considered will be:

- Multivariate normal data, mean an arbitrary function of θ
- Multivariate normal data, mean a linear function of θ
- Multivariate normal data, diagonal covariance
- Poisson data, mean an arbitrary function of θ
- Poisson data, mean a linear function of θ

The general multivariate normal problem

From the last lecture, a general multivariate normal PDF has the form:

$$\text{pr}(\mathbf{g}) = \left[(2\pi)^M \det(\mathbf{K}) \right]^{-1/2} \exp \left[-\frac{1}{2} (\mathbf{g} - \bar{\mathbf{g}})^t \mathbf{K}^{-1} (\mathbf{g} - \bar{\mathbf{g}}) \right], \quad (8.185)$$

where $\bar{\mathbf{g}}$ is the mean vector and \mathbf{K} is the covariance matrix of \mathbf{g} . The most general likelihood function is obtained by letting the mean and covariance both be functions of $\boldsymbol{\theta}$:

$$\text{pr}(\mathbf{g}|\boldsymbol{\theta}) = \left[(2\pi)^M \det[\mathbf{K}(\boldsymbol{\theta})] \right]^{-1/2} \exp \left\{ -\frac{1}{2} [\mathbf{g} - \bar{\mathbf{g}}(\boldsymbol{\theta})]^t [\mathbf{K}(\boldsymbol{\theta})]^{-1} [\mathbf{g} - \bar{\mathbf{g}}(\boldsymbol{\theta})] \right\}.$$

In many problems, the noise is independent of $\boldsymbol{\theta}$, and the log-likelihood in that case is:

$$\ln \text{pr}(\mathbf{g}|\boldsymbol{\theta}) = \text{constant} - \frac{1}{2} [\mathbf{g} - \bar{\mathbf{g}}(\boldsymbol{\theta})]^t \mathbf{K}^{-1} [\mathbf{g} - \bar{\mathbf{g}}(\boldsymbol{\theta})],$$

which is the most general form of a weighted least-squares functional. Maximum likelihood is the same as minimum least-squares residual.

Basic rule of ML estimation: Gaussian noise \Rightarrow least-squares fitting

Multivariate normal data, mean a linear function of θ

In many imaging problems with Gaussian noise, we can assume that the covariance is independent of θ and that

$$\bar{\mathbf{g}}(\theta) = \mathbf{H} \theta$$

Then the log-likelihood is given by

$$\begin{aligned} \ln \text{pr}(\mathbf{g}|\theta) &= \text{constant} - \frac{1}{2}[\mathbf{g} - \mathbf{H} \theta]^t \mathbf{K}^{-1}[\mathbf{g} - \mathbf{H} \theta] \\ &= \text{constant} - \frac{1}{2} \sum_{m=1}^M \sum_{m'=1}^M [\mathbf{g} - \mathbf{H} \theta]_m [\mathbf{K}^{-1}]_{mm'} [\mathbf{g} - \mathbf{H} \theta]_{m'}, \end{aligned}$$

and the Fisher information matrix is

$$\mathbf{F} = \mathbf{H}^t \mathbf{K}^{-1} \mathbf{H}.$$

Multivariate normal data, i.i.d. noise

In many imaging problems (e.g., CCD arrays), there are no noise correlations between detector pixels, and the pixels can be modeled as identical. Then the covariance matrix is

$$\mathbf{K} = \sigma^2 \mathbf{I}.$$

With this noise assumption and $\bar{\mathbf{g}}(\boldsymbol{\theta}) = \mathbf{H} \boldsymbol{\theta}$, the log-likelihood becomes

$$\ln \text{pr}(\mathbf{g}|\boldsymbol{\theta}) = \text{constant} - \frac{1}{2\sigma^2} \|\mathbf{g} - \mathbf{H} \boldsymbol{\theta}\|^2,$$

and the Fisher information matrix is an old friend:

$$\mathbf{F} = \frac{1}{\sigma^2} \mathbf{H}^t \mathbf{H}.$$

With these assumptions, ML estimation is identical to ordinary linear least-squares regression. Maximizing $\ln \text{pr}(\mathbf{g}|\boldsymbol{\theta})$ is the same as *minimizing* the residual norm, $\|\mathbf{g} - \mathbf{H} \boldsymbol{\theta}\|^2$

Poisson data, mean an arbitrary function of θ

From the last lecture, a general multivariate Poisson data set has a probability (not PDF) given by

$$\Pr(\mathbf{g}) = \prod_{m=1}^M \exp(-\bar{g}_m) \frac{[\bar{g}_m]^{g_m}}{g_m!} .$$

Since this probability has only one parameter (the mean) for each m , the general Poisson likelihood is

$$\Pr(\mathbf{g}|\boldsymbol{\theta}) = \prod_{m=1}^M \exp[-\bar{g}_m(\boldsymbol{\theta})] \frac{[\bar{g}_m(\boldsymbol{\theta})]^{g_m}}{g_m!} ,$$

and the log-likelihood is

$$\ln \Pr(\mathbf{g}|\boldsymbol{\theta}) = \sum_{m=1}^M \{ -\bar{g}_m(\boldsymbol{\theta}) + g_m \ln[\bar{g}_m(\boldsymbol{\theta})] - \ln g_m! \} .$$

Poisson data, arbitrary mean – cont.

When the mean is an arbitrary function of θ , it turns out (derivation on request) that the Fisher information has components:

$$F_{jk} = \sum_{m=1}^M \frac{1}{\bar{g}_m(\theta)} \frac{\partial \bar{g}_m(\theta)}{\partial \theta_j} \frac{\partial \bar{g}_m(\theta)}{\partial \theta_k}$$

All you ever need to know to compute the likelihood or the Fisher information with Poisson data is $\bar{g}_m(\theta)$.

Poisson data, mean a linear function of θ

This is the usual case in image reconstruction from photon-limited data:

$$\bar{g}_m(\theta) = \sum_{n=1}^N H_{mn} \theta_n .$$

The log-likelihood is now

$$\ln \Pr(\mathbf{g}|\theta) = \sum_{m=1}^M \left\{ - \sum_{n=1}^N H_{mn} \theta_n + g_m \ln \left[\sum_{n=1}^N H_{mn} \theta_n \right] - \ln g_m! \right\} .$$

The k^{th} component of the score is

$$\frac{\partial}{\partial \theta_k} \ln \Pr(\mathbf{g}|\theta) = \sum_{m=1}^M \left\{ -H_{mk} + \frac{g_m H_{mk}}{\left[\sum_{n=1}^N H_{mn} \theta_n \right]} \right\} .$$

Now all we need is an algorithm to find nonnegative values of each θ_k so that all of these partial derivatives vanish.

The expectation-maximization (EM) algorithm

From last slide,

$$\frac{\partial}{\partial \theta_k} \ln \Pr(\mathbf{g}|\boldsymbol{\theta}) = \sum_{m=1}^M \left\{ -H_{mk} + \frac{g_m H_{mk}}{\left[\sum_{n=1}^N H_{mn} \theta_n \right]} \right\} = 0, \quad \forall k, \text{ at } \boldsymbol{\theta} = \hat{\boldsymbol{\theta}}_{\text{ML}}.$$

A tad of algebra shows that

$$\sum_{m=1}^M H_{mk} \frac{g_m}{[\mathbf{H} \boldsymbol{\theta}]_m} = \sum_{m=1}^M H_{mk}$$

Now multiply both sides by θ_k :

$$\theta_k \sum_{m=1}^M H_{mk} \frac{g_m}{[\mathbf{H} \boldsymbol{\theta}]_m} = \theta_k \sum_{m=1}^M H_{mk}$$

and rewrite as

$$\theta_k = \theta_k \frac{1}{S_k} \sum_{m=1}^M H_{mk} \frac{g_m}{[\mathbf{H} \boldsymbol{\theta}]_m}, \quad \text{where } S_k = \sum_{m=1}^M H_{mk}.$$

The EM algorithm – punchline

From last slide, the ML estimate is obtained if

$$\theta_k = \theta_k \frac{1}{S_k} \sum_{m=1}^M H_{mk} \frac{g_m}{[\mathbf{H} \boldsymbol{\theta}]_m}, \quad \text{where } S_k = \sum_{m=1}^M H_{mk}.$$

Now treat this as an iteration rule. Let $\hat{\theta}_k^{(n)}$ be the estimate of θ_k at the n^{th} iteration and write

$$\theta_k^{(n+1)} = \theta_k^{(n)} \frac{1}{S_k} \sum_{m=1}^M H_{mk} \frac{g_m}{[\mathbf{H} \hat{\boldsymbol{\theta}}^{(n)}]_m}.$$

This iteration rule is called the EM (expectation-maximization) algorithm or, in the optics literature, the Richardson-Lucy algorithm.

Discuss convergence

Still to come – applications in astronomy

- ML estimation with Shack-Hartman wavefront sensors
- The optimal generic wavefront sensor
- Image reconstruction in SCIDAR
- Estimation of fringe visibility in interferometry
- Gauss-Markov estimation for stellar photometry