

The Maximum Likelihood Estimation Method

(last modified August 14, 2009)

Xianguo Lu

xianguo.lu@desy.de

Motivation & Outlines

Analytically illustrate the maximum likelihood method, from the basic to the advanced, with MC examples; implement the method to practical data analysis.

- Part 1: all about the standard method
 1. Concepts
 2. Analysis
 3. MC Examples
- Part 2: "weighting" and the extended method
 1. Concepts
 2. Analysis
 3. Application
- Part 3: for the combined analysis

Part 1:

THE STANDARD

Concepts

- Observable

$$x \in \mathbf{X}. \quad (1)$$

- Data set

$$x_1, x_2, \dots, x_N. \quad (2)$$

Remark.

"Experiment", consist of N observations, obtaining x_1, x_2, \dots, x_N .

•

$$\text{Prob} (x' < x < x' + dx) = p (x') dx, \quad (3)$$

$p (x)$, **Probability density function (p.d.f.)**; $p (x|\theta_1, \theta_2, \dots, \theta_M)$;

$\theta_1, \theta_2, \dots, \theta_M$, **parameters**, $\underline{\theta} \equiv \text{col} (\theta_1, \theta_2, \dots, \theta_M)$;

$\underline{\theta}^*$, the **true parameters** corresponding to the data;

(Assuming p is twice differentiable with respect to θ_i) $\frac{\partial}{\partial \underline{\theta}} \equiv \text{col} \left(\frac{\partial}{\partial \theta_1}, \frac{\partial}{\partial \theta_2}, \dots, \frac{\partial}{\partial \theta_M} \right)$,
gradient operator in the parameter space.

$$\int_{\mathbf{X}} p (x|\underline{\theta}) dx = 1, \quad \langle f(x) \rangle \equiv \int_{\mathbf{X}} f(x) p (x|\underline{\theta}^*) dx. \quad (4)$$

Remark.

- Expectation for repeated *observations*.
- $\langle f(x) \rangle_{\underline{\theta}} \equiv \int_{\mathbf{X}} f(x) p (x|\underline{\theta}) dx$.

- Joint p.d.f./likelihood function

$$L(x_1, x_2, \dots, x_N | \underline{\theta}) \equiv \prod_{i=1}^N p(x_i | \underline{\theta}).$$

$$\int \cdots \int_{\mathbf{X}^N} L(x_1, x_2, \dots, x_N | \underline{\theta}) dx_1 dx_2 \dots dx_N = 1. \quad (5)$$

$$\begin{aligned} \langle g(x_1, x_2, \dots, x_N) \rangle &\equiv \\ &\int \cdots \int_{\mathbf{X}^N} g(x_1, x_2, \dots, x_N) L(x_1, x_2, \dots, x_N | \underline{\theta}^*) dx_1 dx_2 \dots dx_N. \end{aligned} \quad (6)$$

Remark.

- Expectation for repeated *experiments*.
- $\langle g(x_1, x_2, \dots, x_N) \rangle_{\underline{\theta}} \equiv \int \cdots \int_{\mathbf{X}^N} g(x_1, x_2, \dots, x_N) L(x_1, x_2, \dots, x_N | \underline{\theta}) dx_1 dx_2 \dots dx_N.$

- Log-likelihood function

$$l(x_1, x_2, \dots, x_N | \underline{\theta}) \equiv \ln L(x_1, x_2, \dots, x_N | \underline{\theta}) \quad (7)$$

$$= \sum_{i=1}^N \ln p(x_i | \underline{\theta}). \quad (8)$$

- Estimators

$$\hat{\underline{\theta}}(x_1, x_2, \dots, x_N). \quad (9)$$

Remark.

Also random.

- Bias

$$\underline{b}_{\hat{\underline{\theta}}} \equiv \langle \hat{\underline{\theta}} \rangle - \underline{\theta}^*. \quad (10)$$

- Covariance matrix

$$\underline{\underline{V}}_{\hat{\underline{\theta}}} \equiv \left\langle \left(\hat{\underline{\theta}} - \langle \hat{\underline{\theta}} \rangle \right) \left(\hat{\underline{\theta}} - \langle \hat{\underline{\theta}} \rangle \right)^T \right\rangle. \quad (11)$$

- Minimum variance bound, MVB

$$\underline{\underline{V}}_{\hat{\underline{\theta}}} \geq \left(\frac{\partial \langle \hat{\theta}_i \rangle_{\underline{\theta}}}{\partial \theta_j} \Big|_{\underline{\theta}^*} \right) \left(\left\langle \frac{\partial l}{\partial \theta_i} \frac{\partial l}{\partial \theta_j} \Big|_{\underline{\theta}^*} \right\rangle \right)^{-1} \left(\frac{\partial \langle \hat{\theta}_j \rangle_{\underline{\theta}}}{\partial \theta_i} \Big|_{\underline{\theta}^*} \right). \quad (12)$$

[Hint: Using Schwarz's Inequality (generalized, cf. Appendix),

$\langle \underline{a} \underline{a}^T \rangle \geq \langle \underline{a} \underline{b}^T \rangle \langle \underline{b} \underline{b}^T \rangle^{-1} \langle \underline{b} \underline{a}^T \rangle$, let $\underline{a} = \hat{\underline{\theta}} - \langle \hat{\underline{\theta}} \rangle$, $\underline{b} = \frac{\partial l}{\partial \underline{\theta}} \Big|_{\underline{\theta}^*}$. (If matrix $\underline{\underline{A}}$ is positive definite (semidefinite), i.e. $\forall \underline{v} \neq \underline{0}$, $\underline{v}^T \underline{\underline{A}} \underline{v} > (\geq) 0$, we write $\underline{\underline{A}} > (\geq) \underline{0}$).]

Remark.

– For unbiased estimators, $\langle \hat{\theta}_i \rangle_{\underline{\theta}} = \theta_i$,

$$\therefore \underline{V}_{\hat{\underline{\theta}}} \geq (\delta_{ij}) \left(\left\langle \frac{\partial l}{\partial \theta_i} \frac{\partial l}{\partial \theta_j} \middle|_{\underline{\theta}^*} \right\rangle \right)^{-1} (\delta_{ji}) \quad (13)$$

$$= \left(\left\langle \frac{\partial l}{\partial \theta_i} \frac{\partial l}{\partial \theta_j} \middle|_{\underline{\theta}^*} \right\rangle \right)^{-1} = \frac{1}{N} \left(\left\langle \frac{\partial \ln p}{\partial \theta_i} \frac{\partial \ln p}{\partial \theta_j} \middle|_{\underline{\theta}^*} \right\rangle \right)^{-1}. \quad (14)$$

Let $\underline{v} = \text{col}(0, \dots, \underbrace{1}_{i\text{th component}}, \dots, 0)$, the **variance**

$$\sigma_{\hat{\theta}_i}^2 \equiv [\underline{V}_{\hat{\underline{\theta}}}]_{ii} = \underline{v}^T \underline{V}_{\hat{\underline{\theta}}} \underline{v} \geq \underline{v}^T \frac{1}{N} \left(\left\langle \frac{\partial \ln p}{\partial \theta_i} \frac{\partial \ln p}{\partial \theta_j} \middle|_{\underline{\theta}^*} \right\rangle \right)^{-1} \underline{v} = \frac{1}{N} \left[\left(\left\langle \frac{\partial \ln p}{\partial \theta_i} \frac{\partial \ln p}{\partial \theta_j} \middle|_{\underline{\theta}^*} \right\rangle \right)^{-1} \right]_{ii}. \quad (15)$$

In the single parameter case,

$$\sigma_{\hat{\theta}}^2 \geq \frac{1}{N \left\langle \left(\frac{\partial \ln p}{\partial \theta} \right)^2 \middle|_{\theta^*} \right\rangle}. \quad (16)$$

[Hint: $\frac{\partial l}{\partial \theta_i} \frac{\partial l}{\partial \theta_j} = \sum_m \frac{\partial \ln p(x_m)}{\partial \theta_i} \frac{\partial \ln p(x_m)}{\partial \theta_j} + \sum_m \sum_{n \neq m} \frac{\partial \ln p(x_m)}{\partial \theta_i} \frac{\partial \ln p(x_n)}{\partial \theta_j}$, the second term having an expectation of 0.]

- Point estimation

Criteria,

1. **Consistency:** $\hat{\theta} \xrightarrow[N \rightarrow \infty]{P} \theta^*$, ($\xrightarrow[N \rightarrow \infty]{P}$ means converging in probability as $N \rightarrow \infty$).

Remark.

If the statistics is large enough, the estimation gives the true values;

2. **Unbiasedness:** $b_{\hat{\theta}} = 0$,

Remark.

If the experiment is repeated for many times, even with low statistics, the estimators can give the true values;

3. **Efficiency:** $\sigma_{\hat{\theta}_i}^2$ is as small as possible,

Remark.

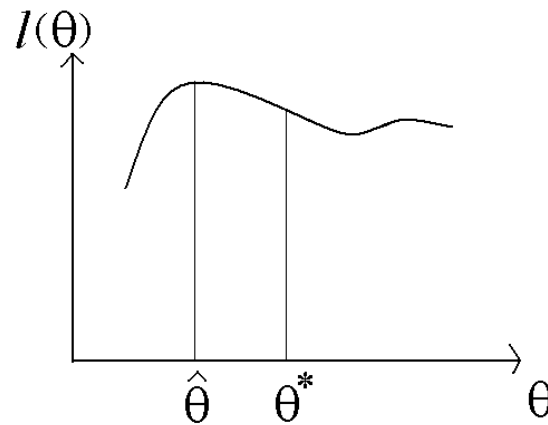
The larger efficiency, the less repeated experiments are required.

...

Analysis

Maximum Likelihood Estimators:

the ones maximizing the likelihood function L .



Why? The one happened is the most probable?

Need analytical proofs of the criteria; common sense is *NOT* always correct.

Analysis

- Strong Law of Large Numbers

$$\frac{1}{N} \sum_{i=1}^N f(x_i) \xrightarrow[N \rightarrow \infty]{P} \langle f(x) \rangle. \quad (17)$$

Remark.

- Large sample limit.
- $l = \sum_{i=1}^N \ln p(x_i | \underline{\theta}) \xrightarrow[N \rightarrow \infty]{P} N \langle \ln p(x | \underline{\theta}) \rangle$, independent of x in the large sample limit.

$$\left\langle \frac{\partial \ln p}{\partial \underline{\theta}} \Big|_{\underline{\theta}^*} \right\rangle = \underline{0}, \quad \left\langle \frac{\partial^2 \ln p}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right\rangle = - \left\langle \left(\frac{\partial \ln p}{\partial \theta_i} \frac{\partial \ln p}{\partial \theta_j} \right) \Big|_{\underline{\theta}^*} \right\rangle. \quad (18)$$

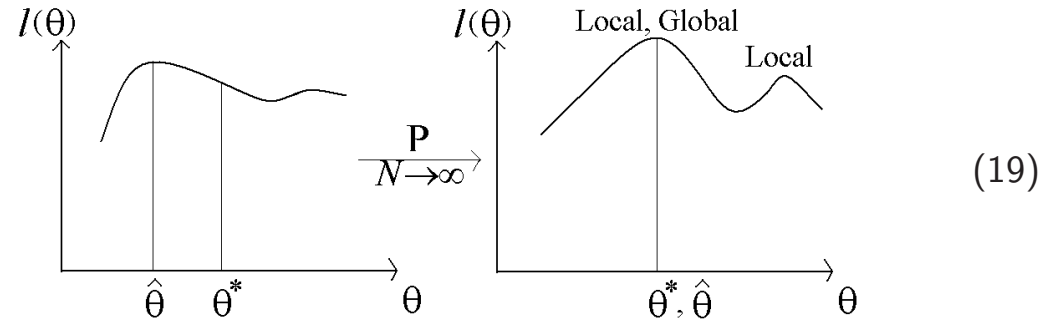
[Hint: $\left\langle \frac{\partial \ln p}{\partial \underline{\theta}} \Big|_{\underline{\theta}^*} \right\rangle = \int \left(\frac{\partial \ln p}{\partial \underline{\theta}} p \right)_{\underline{\theta}^*} dx = \int \frac{\partial p}{\partial \underline{\theta}} \Big|_{\underline{\theta}^*} dx = \left(\frac{\partial}{\partial \underline{\theta}} \int p dx \right)_{\underline{\theta}^*}$ (\mathbf{X} needs to be independent of the parameters), $\frac{\partial^2 \ln p}{\partial \theta_i \partial \theta_j} = \frac{\partial}{\partial \theta_i} \left(\frac{\partial \ln p}{\partial \theta_j} \right) = -\frac{\partial \ln p}{\partial \theta_i} \frac{\partial \ln p}{\partial \theta_j} + \frac{1}{p} \frac{\partial^2 p}{\partial \theta_i \partial \theta_j}$.]

Analysis

1. Consistency:

$$\hat{\underline{\theta}} \xrightarrow[N \rightarrow \infty]{P} \underline{\theta}^*$$

$$\Leftrightarrow \forall \underline{\delta} \neq \underline{0}, l(\underline{\theta}^* + \underline{\delta}) < l(\underline{\theta}^*).$$



Locally ($\underline{\delta} \rightarrow \underline{0}$), approximated to the second order,

$$l(\underline{\theta}^* + \underline{\delta}) = l(\underline{\theta}^*) + \left(\frac{\partial l}{\partial \underline{\theta}} \right)_{\underline{\theta}^*}^T \underline{\delta} + \frac{1}{2} \underline{\delta}^T \left(\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right) \underline{\delta}. \quad (20)$$

$$\Rightarrow \begin{cases} \frac{\partial l}{\partial \underline{\theta}} \Big|_{\underline{\theta}^*} = \underline{0}, \\ \text{Hessian } \underline{\underline{H}}(l) \equiv \left(\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right) < 0 \\ (\text{Reminder. negative definite, i.e. } \forall \underline{v} \neq \underline{0}, \underline{v}^T \underline{\underline{H}} \underline{v} < 0). \end{cases} \quad \begin{matrix} \\ \\ \text{(Second Derivative Test)} \end{matrix} \quad (21)$$

Proof. of local maximum

$$l \xrightarrow[N \rightarrow \infty]{P} N \langle \ln p \rangle \quad (22)$$

$$\therefore \begin{cases} \frac{\partial l}{\partial \underline{\theta}}|_{\theta^*} \xrightarrow[N \rightarrow \infty]{P} N \left\langle \frac{\partial \ln p}{\partial \underline{\theta}}|_{\theta^*} \right\rangle = 0 \\ H_{ij} \xrightarrow[N \rightarrow \infty]{P} N \left\langle \frac{\partial^2 \ln p}{\partial \theta_i \partial \theta_j} |_{\underline{\theta}^*} \right\rangle = -N \left\langle \left(\frac{\partial \ln p}{\partial \theta_i} \frac{\partial \ln p}{\partial \theta_j} \right)_{\underline{\theta}^*} \right\rangle. \end{cases} \quad (23)$$

Denote $\frac{\partial \ln p}{\partial \underline{\theta}}|_{\underline{\theta}^*} = \underline{d}(x)$, so

$$\underline{\underline{H}} \xrightarrow[N \rightarrow \infty]{P} -N \left\langle \underline{d} \underline{d}^T \right\rangle, \quad (24)$$

$$\forall \underline{v} \neq 0, \underline{v}^T \underline{\underline{H}} \underline{v} \xrightarrow[N \rightarrow \infty]{P} -N \left\langle \underline{v}^T \underline{d} \underline{d}^T \underline{v} \right\rangle \quad (25)$$

$$= -N \left\langle \left(\underline{v}^T \underline{d} \right)^2 \right\rangle < 0 \text{ if } \exists x' \in \mathbf{X}, \underline{v}^T \underline{d}(x') \neq 0. \quad (26)$$

□

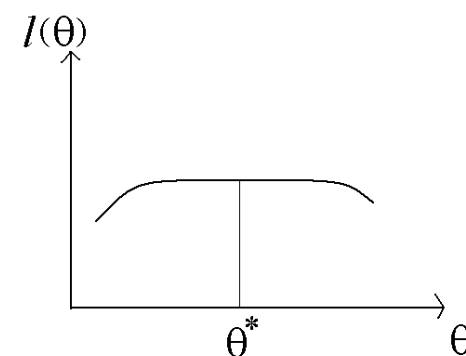
Remark.

- For a proof of global maximum, please refer to Wald's (1949) proof, e.g. in Section 17.15 of *The advanced theory of statistics*, by M. Kendall et al., 4th ed. of Vol. 2 of the 3-volume ed. (1979) ISBN: 0 85264 255 5.
- This proof is **Conditional**; it **fails if**

$$\exists \underline{v} \neq 0, \forall x \in \mathbf{X},$$

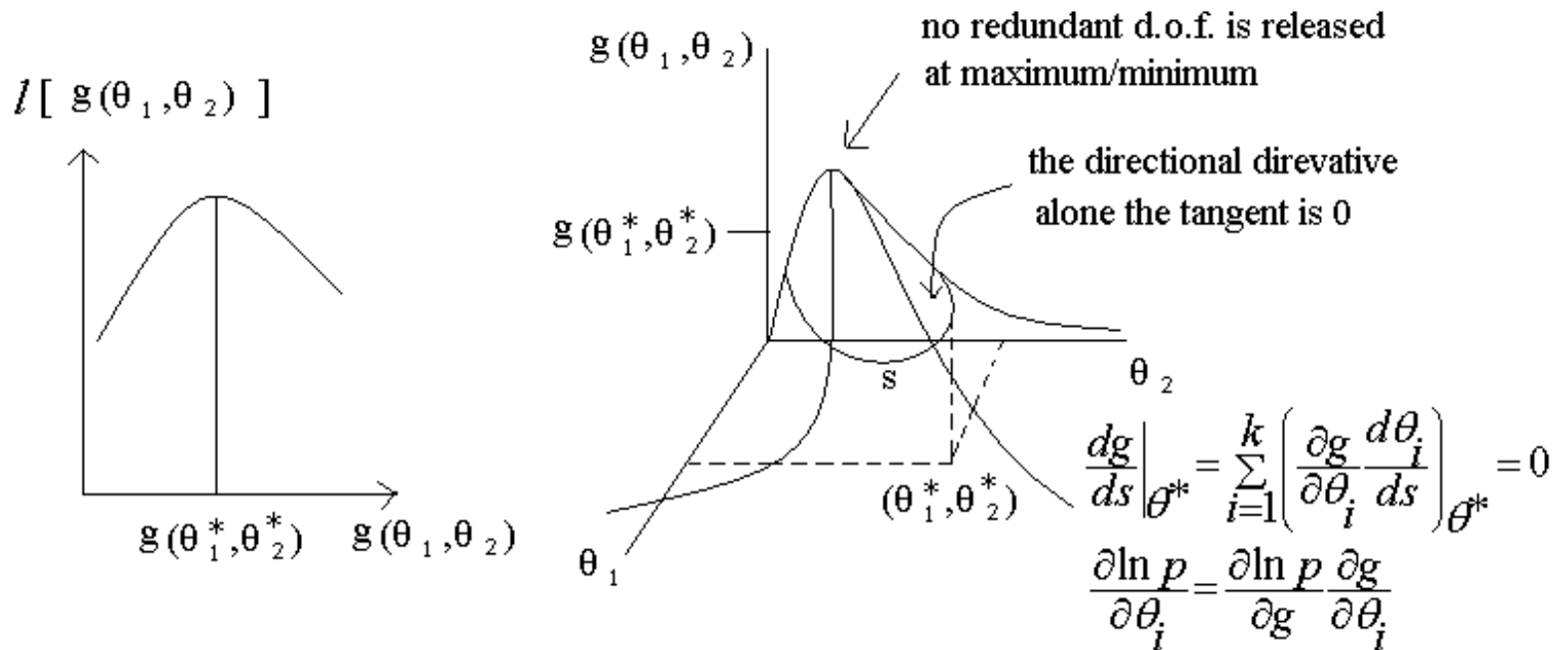
$$\underline{v}^T \underline{d} = v_1 \frac{\partial \ln p}{\partial \theta_1} \Big|_{\underline{\theta}^*} + v_2 \frac{\partial \ln p}{\partial \theta_2} \Big|_{\underline{\theta}^*} + \cdots + v_M \frac{\partial \ln p}{\partial \theta_M} \Big|_{\underline{\theta}^*} = 0,$$

i.e. $\frac{\partial \ln p}{\partial \theta_1} \Big|_{\underline{\theta}^*}, \frac{\partial \ln p}{\partial \theta_2} \Big|_{\underline{\theta}^*}, \dots, \frac{\partial \ln p}{\partial \theta_M} \Big|_{\underline{\theta}^*}$ are linearly dependent functions of x ; in the single parameter case, $\frac{\partial \ln p}{\partial \theta} \Big|_{\underline{\theta}^*} = 0$.



- The condition ensures that the maximum likelihood method is justified and imposes some **constraints on**
 - the **parameterization/construction of the p.d.f.**, e.g., the following are **inappropriate**:
 - (a) $\forall x, \exists i, \frac{\partial p}{\partial \theta_i} \Big|_{\underline{\theta}^*} = 0$;
 - (b) $\forall x, \exists i, \frac{\partial p}{\partial \theta_i} \equiv 0$, that is, p is independent of θ_i , typically when normalizing the p.d.f.;

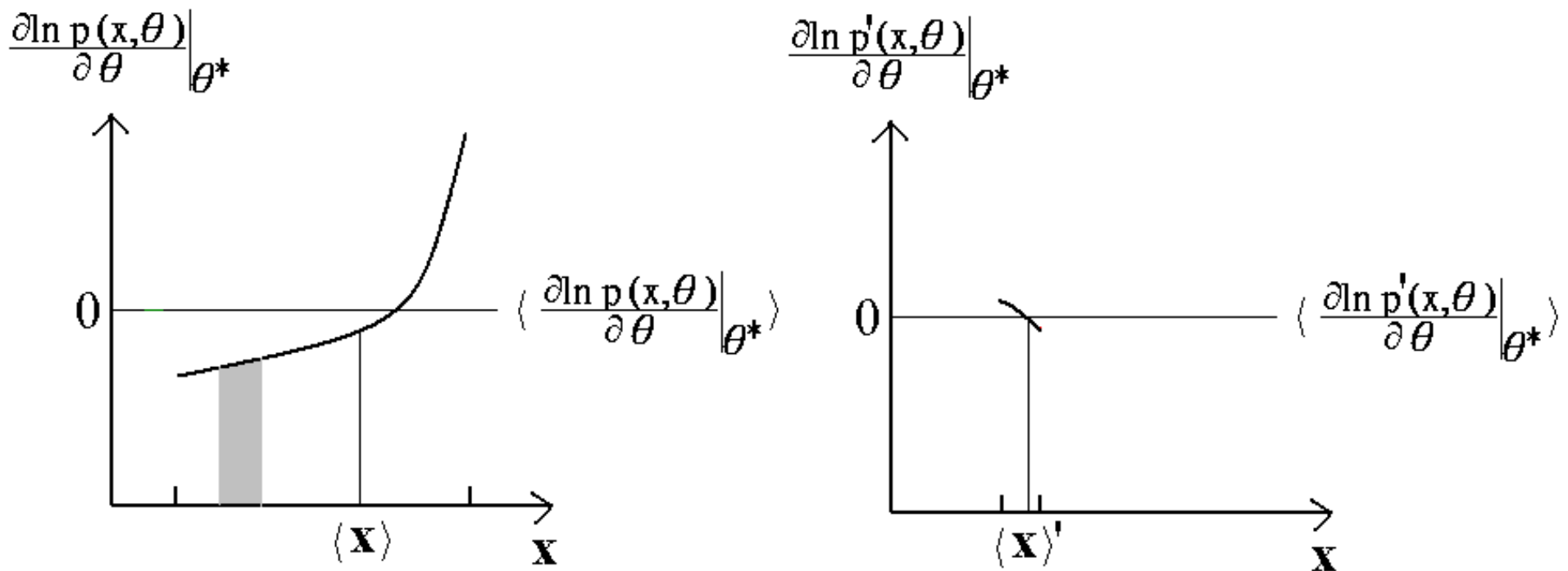
- (c) $p(x|\theta_1, \theta_2, \dots, \theta_M) = p[x|g(\theta_1, \theta_2, \dots, \theta_k), \dots]$, $2 \leq k \leq M$, that is, the dependence of $\theta_1, \theta_2, \dots, \theta_k$ degenerates into a whole via $g(\theta_1, \theta_2, \dots, \theta_k)$ and thus redundant degree(s) of freedom are(is) released. (Except that $g(\theta_1^*, \theta_2^*, \dots, \theta_k^*)$ is maximum/minimum, which belongs to Situation a.) Suppose $p(x|\theta_1, \theta_2) = [1 - (\theta_1^2 + \theta_2^2)x] / \text{nor}$, and thus $g = \theta_1^2 + \theta_2^2$. In maximizing $l = l[g(\theta_1, \theta_2)]$, the program enters the region $g(\theta_1, \theta_2) = g(\theta_1^*, \theta_2^*)$ where l has its largest value, but still can not decide which (θ_1, θ_2) is the true one:



- the **selection of data samples**, e.g.,

the "spread" of x , measured by its RMS $\equiv \sqrt{\langle (x - \langle x \rangle)^2 \rangle}$, should be large, since as

$x \rightarrow$ certain fixed value, $\forall x$, $\frac{\partial \ln p}{\partial \theta_i} \rightarrow \langle \frac{\partial \ln p}{\partial \theta_i} \rangle = 0$; and thus $\forall x$, $\frac{\partial \ln p}{\partial \theta_i} \equiv 0$ holds approximately, making the method unjustified. (This constrain holds no matter whether $\langle x \rangle$ is zero.)



– the **grouping of the data sample**:

The data set can be divided into several, say K , sub-sets, each of which has its p.d.f.: $p_1(x|\underline{\theta})$, $p_2(x|\underline{\theta})$, \dots , $p_K(x|\underline{\theta})$, and correspondingly $L_J = \prod_i p_J(x_i|\underline{\theta})$, $J = 1, 2, \dots, K$. For the whole data set, the likelihood function reads

$$L_s = L_1 L_2 \cdots L_K \quad (27)$$

$$\therefore l_s = l_1 + l_2 + \cdots + l_K \quad (28)$$

$$= \sum_i^{N_1} \ln p_1(x_i|\underline{\theta}) + \sum_i^{N_2} \ln p_2(x_i|\underline{\theta}) + \cdots + \sum_i^{N_K} \ln p_K(x_i|\underline{\theta}) \quad (29)$$

$$\neq l = \left(\sum_i^{N_1} + \sum_i^{N_2} + \cdots + \sum_i^{N_K} \right) \ln p(x_i|\underline{\theta}), \text{ since generally } p_J \neq p. \quad (30)$$

One can verify (by following the same approaches as above) that maximizing l_s also gives consistency results conditionally and since

$$[H_s]_{mn} \xrightarrow[N \rightarrow \infty]{P} - \sum_{J=1}^K N_J \left\langle \left(\frac{\partial \ln p_J}{\partial \theta_m} \frac{\partial \ln p_J}{\partial \theta_n} \right)_{\underline{\theta}^*} \right\rangle, \quad (31)$$

the grouping **may be unjustified if all the sub-sets are close to any of the above inappropriate situations.**

2. Bias and Efficiency

Using Taylor's Theorem,

$$\frac{\partial l}{\partial \theta_i} \Big|_{\hat{\theta}} = \frac{\partial l}{\partial \theta_i} \Big|_{\underline{\theta}^*} + \sum_j^M \frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^\Delta} (\hat{\theta}_j - \theta_j^*) \Rightarrow \frac{\partial l}{\partial \underline{\theta}} \Big|_{\hat{\theta}} = \frac{\partial l}{\partial \underline{\theta}} \Big|_{\underline{\theta}^*} + \left(\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^\Delta} \right) (\hat{\underline{\theta}} - \underline{\theta}^*), \quad (32)$$

where $\underline{\theta}^\Delta$ are some values between $\hat{\underline{\theta}}$ and $\underline{\theta}^*$. Since $\frac{\partial l}{\partial \underline{\theta}} \Big|_{\hat{\theta}} = 0$ and $\underline{\theta}^\Delta \xrightarrow[N \rightarrow \infty]{P} \underline{\theta}^*$,

$$\therefore \left\{ \begin{array}{l} (\hat{\underline{\theta}} - \underline{\theta}^*) = \left(-\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^\Delta} \right)^{-1} \frac{\partial l}{\partial \underline{\theta}} \Big|_{\underline{\theta}^*}, \text{ if } \left(-\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^\Delta} \right) \text{ is nonsingular} \\ \left(-\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^\Delta} \right) \xrightarrow[N \rightarrow \infty]{P} \left(-\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right), \text{ independent of } x, \text{ nonsingular} \\ \text{(assuming the condition for consistency holds)} \end{array} \right. \quad (33)$$

(34)

$$\underline{b}_{\hat{\theta}} = \left\langle \left(-\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^\Delta} \right)^{-1} \frac{\partial l}{\partial \underline{\theta}} \Big|_{\underline{\theta}^*} \right\rangle \xrightarrow[N \rightarrow \infty]{P} \left(-\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right)^{-1} \left\langle \frac{\partial l}{\partial \underline{\theta}} \Big|_{\underline{\theta}^*} \right\rangle = \underline{0}. \quad (35)$$

Similarly, expanding at $\langle \hat{\underline{\theta}} \rangle$, with $\underline{\theta}_2^\Delta$ between $\hat{\underline{\theta}}$ and $\langle \hat{\underline{\theta}} \rangle$,

$$\underline{\underline{V}}_{\hat{\underline{\theta}}} = \left\langle \left(-\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}_2^\Delta} \right)^{-1} \left(\frac{\partial l}{\partial \theta_i} \frac{\partial l}{\partial \theta_j} \Big|_{\langle \hat{\underline{\theta}} \rangle} \right) \left(-\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}_2^\Delta} \right)^{-1} \right\rangle \quad (36)$$

$$\xrightarrow[N \rightarrow \infty]{P} \left(-\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right)^{-1} \left(\left\langle \frac{\partial l}{\partial \theta_i} \frac{\partial l}{\partial \theta_j} \Big|_{\underline{\theta}^*} \right\rangle \right) \left(-\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right)^{-1}. \quad (37)$$

$$\underline{\underline{V}}_{\hat{\underline{\theta}}} \xrightarrow[N \rightarrow \infty]{P} \frac{1}{N} \left(\left\langle \frac{\partial \ln p}{\partial \theta_i} \frac{\partial \ln p}{\partial \theta_j} \Big|_{\underline{\theta}^*} \right\rangle \right)^{-1} = \left(-\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right)^{-1}. \quad (38)$$

Remark.

- $\hat{\underline{\theta}}$ is *not* unbiased. As N tends to infinity, the biases tend to zero and $\underline{\underline{V}}_{\hat{\underline{\theta}}}$ tends to the MVB.
- $\underline{\underline{V}}_{\hat{\underline{\theta}}}$ can be estimated by $\left(-\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\hat{\underline{\theta}}} \right)^{-1}$ in the large sample limit.

MC Examples

1. MINUIT

- Searching *local* minima of FCN as a function of par: $\text{FCN} = -l$, $\text{par}_i = \theta_i$;
- **Error matrix:** $2 \left(\frac{\partial^2_{\text{FCN}}}{\partial \text{par}_i \partial \text{par}_j} \right)^{-1} \times \text{UP} : \text{UP} = 0.5$;
- The MINUIT command

`gMinuit->mnexcm("MIGRAD", arglist, 2, ierflg)`

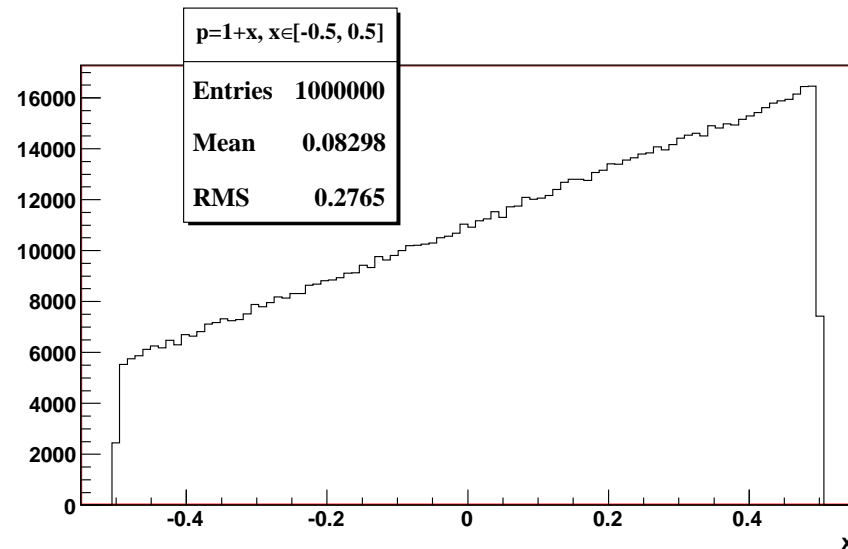
searches a *local* minimum of FCN and produces the error matrix.

2. $p = (1 + \theta x) / \text{nor}$

The prototype-p.d.f. in asymmetry analysis. E.g., in beam-spin asymmetry, $x = \text{beam pol.}$,
 $\theta = A_{LU}^{c0}$.

MC Examples

- Inappropriate parameterization $\frac{\partial p}{\partial \theta_i} \big|_{\theta^*} = 0$:



1. $p = 1 + (1 + \theta) x$, $\theta^* = 0$: $\frac{\partial p}{\partial \theta} \big|_{\theta^*} = x$, $\hat{\theta} = (-4.5 \pm 3.2) \times 10^{-3}$;
2. $p = 1 + (1 + \theta^2) x$, $\theta^* = 0$: $\frac{\partial p}{\partial \theta} \big|_{\theta^*} = 0$, $\hat{\theta} = (0.2 \pm 33.4) \times 10^{-3}$.

- The smaller RMS, the larger error:

- MC samples:

$$p = \frac{1+x}{\text{nor}}, \quad x \in [x_0 - a, x_0 + a], \quad x_0 = 0.1, \quad N = 10^6;$$

varying a (N fixed) to get samples with different RMS.

(Cf. RMS vs. a).

- Parameterization:

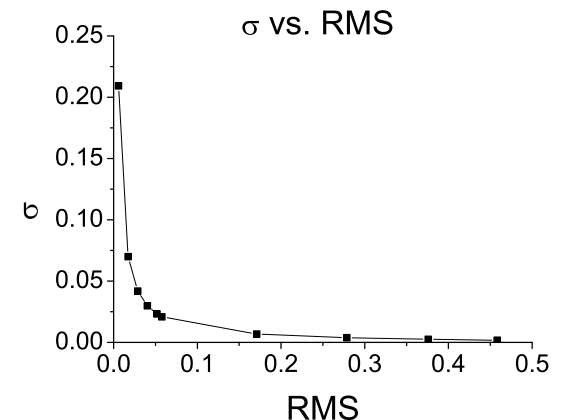
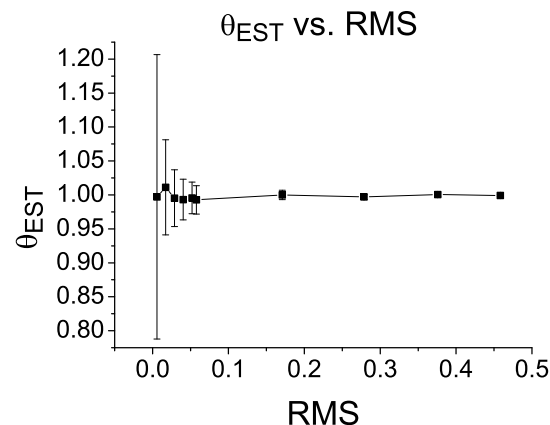
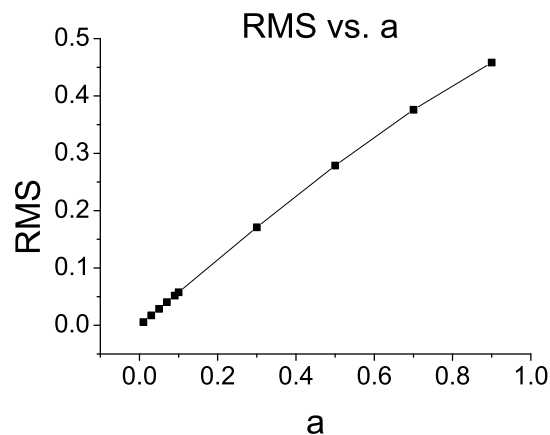
$$p = \frac{1+\theta x}{\text{nor}(\theta)}, \quad \theta^* = 1.$$

(Cf. θ_{EST} vs. RMS for fitting results.)

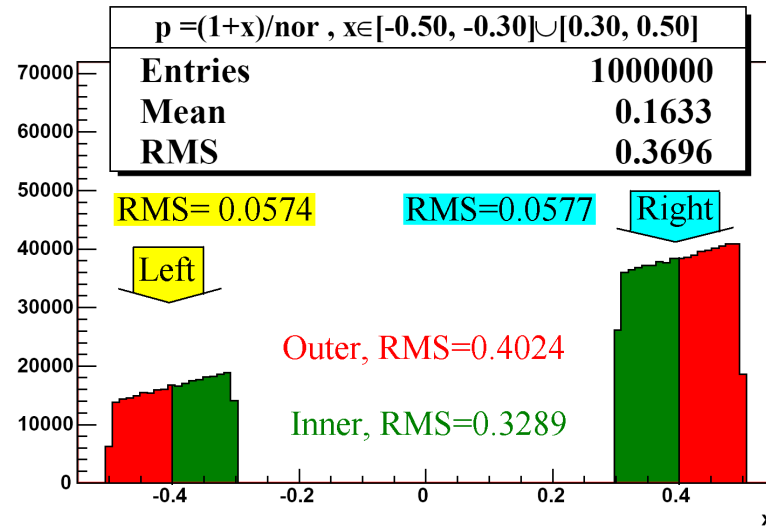
- Conclusion:

The error gets larger with smaller RMS.

(Cf. σ vs. RMS).



- Inappropriate data grouping:



– Parameterization: $p = \frac{1+\theta x}{\text{nor}(\theta)}, \theta^* = 1;$

– Results:

	no grouping	Left-Right	Inner-Outer
$\hat{\theta}$	1.0003 ± 0.0022	1.0146 ± 0.0107	1.0003 ± 0.0022

– Conclusion: inappropriate data grouping can lead to larger fitting errors.

(This explains why in BSA analysis, the constant term of the asymmetry by $l_s = l_{\text{beam pol.}>0} + l_{\text{beam pol.<0}}$ has much larger error.)

Summary for Part 1

1. Introduce *useful formulas* to deal with theoretical problems of maximum likelihood;
2. Conclude the *condition* for the justified maximum likelihood method, illustrate with MC examples and explain certain issue on practical data analysis;
3. Show the *standard procedure to obtain the covariance matrix* of the estimators, which will be used directly in the following parts.

Part 2:

WEIGHTS & THE EXTENDED

Weights – Concepts

- Weight Function $w(x)$:
the multiplicity of x ;
 - Naturally $w(x) = p(x)$ (Prob(x) if x is discrete), e.g. $N = 3$,
 $x_1 = A, x_2 = A, x_3 = B$,
 $w_0(A) = 2/3$. (may not be precise since N is very small)
 - Artificially weights – any non-negative number – can be assigned to x : "**weighting**", e.g.
 assigning weights $w(A) = 3, w(B) = 5$ to *each data point*, we have a "virtual" data sample:

$$x_1^1 = A, x_1^2 = A, x_1^3 = A, \quad (39)$$

$$x_2^1 = A, x_2^2 = A, x_2^3 = A, \quad (40)$$

$$x_3^1 = B, x_3^2 = B, x_3^3 = B, x_3^4 = B, x_3^5 = B. \quad (41)$$

$$\text{Prob}_w(A) = \frac{\text{Prob}(A)w(A)}{\text{Prob}(A)w(A) + \text{Prob}(B)w(B)} = \frac{6}{11}, N_w = w(x_1) + w(x_2) + w(x_3) = 11.$$

- Weighted p.d.f.

$$\text{Prob}_w(A) = \frac{\text{Prob}(A) w(A)}{\text{Prob}(A) w(A) + \text{Prob}(B) w(B)} \quad (42)$$

$$\Rightarrow p_w(x|\underline{\theta}) \equiv \frac{p(x|\underline{\theta}) w(x)}{\int_{\mathbf{X}} p(x|\underline{\theta}) w(x) dx}. \quad (43)$$

- Weighted likelihood function

$$N_w = w(x_1) + w(x_2) + w(x_3) \quad (44)$$

$$\Rightarrow L_w(x_1, x_2, \dots, x_N|\underline{\theta}) \equiv \prod_{i=1}^N [p_w(x_i|\underline{\theta})]^{w(x_i)} \quad (45)$$

– Weighted log-likelihood function

$$l_w(x_1, x_2, \dots, x_N | \underline{\theta}) \equiv \ln L_w(x_1, x_2, \dots, x_N | \underline{\theta}) \quad (46)$$

$$= \sum_{i=1}^N w(x_i) \ln p_w(x_i | \underline{\theta}) \quad (47)$$

$$\xrightarrow[N \rightarrow \infty]{P} N \langle w(x) \ln p_w(x | \underline{\theta}) \rangle \quad (48)$$

Remark.

Since by \sum_i^N we can only deal with real data, by the Strong Law of Large Numbers, the definitions of expectations do *NOT* change:

$$\langle f(x) \rangle \equiv \int_{\mathbf{X}} f(x) p(x | \underline{\theta}^*) dx, \quad (49)$$

$$\begin{aligned} \langle g(x_1, x_2, \dots, x_N) \rangle &\equiv \\ &\int \cdots \int_{\mathbf{X}^N} g(x_1, x_2, \dots, x_N) L(x_1, x_2, \dots, x_N | \underline{\theta}^*) dx_1 dx_2 \cdots dx_N. \end{aligned} \quad (50)$$

Weights – Analysis

1.

$$\left(\text{Reminder. } \left\langle \frac{\partial \ln p}{\partial \underline{\theta}} \middle|_{\underline{\theta}^*} \right\rangle = \underline{0}, \quad \left\langle \frac{\partial^2 \ln p}{\partial \theta_i \partial \theta_j} \middle|_{\underline{\theta}^*} \right\rangle = - \left\langle \left(\frac{\partial \ln p}{\partial \theta_i} \frac{\partial \ln p}{\partial \theta_j} \right) \middle|_{\underline{\theta}^*} \right\rangle. \right) \quad (51)$$

$$\left\langle w \frac{\partial \ln p_w}{\partial \underline{\theta}} \middle|_{\underline{\theta}^*} \right\rangle = 0, \quad \left\langle w \frac{\partial^2 \ln p_w}{\partial \theta_i \partial \theta_j} \middle|_{\underline{\theta}^*} \right\rangle = - \left\langle w \left(\frac{\partial \ln p_w}{\partial \theta_i} \frac{\partial \ln p_w}{\partial \theta_j} \right) \middle|_{\underline{\theta}^*} \right\rangle. \quad (52)$$

$$[\text{Hint: } \left\langle w \frac{\partial \ln p_w}{\partial \underline{\theta}} \middle|_{\underline{\theta}^*} \right\rangle = \int w \left(\frac{\partial \ln p_w}{\partial \underline{\theta}} p \right)_{\underline{\theta}^*} dx = \left(\int p w dx \cdot \frac{\partial}{\partial \underline{\theta}} \int p_w dx \right)_{\underline{\theta}^*}.]$$

2. Consistency *fails* if

$$\left(\text{Reminder. } \frac{\partial \ln p}{\partial \theta_1} \middle|_{\underline{\theta}^*}, \frac{\partial \ln p}{\partial \theta_2} \middle|_{\underline{\theta}^*}, \dots, \frac{\partial \ln p}{\partial \theta_M} \middle|_{\underline{\theta}^*}, \text{ are linearly dependent functions of } x. \right) \quad (53)$$

$$\frac{\partial \ln p_w}{\partial \theta_1} \middle|_{\underline{\theta}^*}, \frac{\partial \ln p_w}{\partial \theta_2} \middle|_{\underline{\theta}^*}, \dots, \frac{\partial \ln p_w}{\partial \theta_M} \middle|_{\underline{\theta}^*} \text{ are linearly dependent functions of } x. \quad (54)$$

3. (Reminder. $\underline{V}_{\hat{\underline{\theta}}} \xrightarrow[N \rightarrow \infty]{P} \left(-\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right)^{-1} \left(\left\langle \frac{\partial l}{\partial \theta_i} \frac{\partial l}{\partial \theta_j} \Big|_{\underline{\theta}^*} \right\rangle \right) \left(-\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right)^{-1} = \left(-\frac{\partial^2 l}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right)^{-1}.$)

$$\underline{V}_{\hat{\underline{\theta}}_w} \xrightarrow[N \rightarrow \infty]{P} \left(-\frac{\partial^2 l_w}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right)^{-1} \left(\left\langle \frac{\partial l_w}{\partial \theta_i} \frac{\partial l_w}{\partial \theta_j} \Big|_{\underline{\theta}^*} \right\rangle \right) \left(-\frac{\partial^2 l_w}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right)^{-1} \quad (55)$$

$$= \left(-\frac{\partial^2 l_w}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right)^{-1} \left(\sum_{k=1}^N w^2(x_k) \left[\frac{\partial \ln p_w(x_k|\underline{\theta})}{\partial \theta_i} \frac{\partial \ln p_w(x_k|\underline{\theta})}{\partial \theta_j} \right] \Big|_{\underline{\theta}^*} \right) \left(-\frac{\partial^2 l_w}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right)^{-1}. \quad (56)$$

Remark.

(a) Suppose $w(x) \equiv w_0 > 0$ constant,

$$p_w = p, \quad l_w = w_0 l, \quad N_w = \sum_i^N w_0 = w_0 N, \quad (57)$$

$$\Rightarrow \left(-\frac{\partial^2 l_w}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right)^{-1} = \frac{1}{w_0 N} \left(\left\langle \frac{\partial \ln p}{\partial \theta_i} \frac{\partial \ln p}{\partial \theta_j} \Big|_{\underline{\theta}^*} \right\rangle \right)^{-1} = \frac{1}{N_w} \left(\left\langle \frac{\partial \ln p}{\partial \theta_i} \frac{\partial \ln p}{\partial \theta_j} \Big|_{\underline{\theta}^*} \right\rangle \right)^{-1}. \quad (58)$$

$\left(-\frac{\partial^2 l_w}{\partial \theta_i \partial \theta_j} \Big|_{\underline{\theta}^*} \right)^{-1}$, defined as the covariance matrix by MINUIT, which takes into account the artificially arbitrary statistical inflation, is *not* the corresponding matrix in Eq. 56.

(b) If, e.g., $\frac{\partial^2 p_w}{\partial \theta_i \partial \theta_j} \big|_{\underline{\theta}^*} = 0$,

$$\left[-\frac{\partial^2}{\partial \theta_i \partial \theta_j} \sum_k^N w^2(x_k) \ln p_w(x_k | \underline{\theta}) \right]_{\hat{\underline{\theta}}} \xrightarrow[N \rightarrow \infty]{P} \sum_{k=1}^N w^2(x_k) \left[\frac{\partial \ln p_w(x_k | \underline{\theta})}{\partial \theta_i} \frac{\partial \ln p_w(x_k | \underline{\theta})}{\partial \theta_j} \right]_{\underline{\theta}^*}. \quad (59)$$

This provides a method to estimate the covariance matrix in Eq. 56.

$$[\text{Hint: } \langle w^2 \frac{\partial^2 \ln p_w}{\partial \theta_i \partial \theta_j} \big|_{\underline{\theta}^*} \rangle = -\langle w^2 (\frac{\partial \ln p_w}{\partial \theta_i} \frac{\partial \ln p_w}{\partial \theta_j})_{\underline{\theta}^*} \rangle + \langle w^2 (\frac{1}{p_w} \frac{\partial^2 p_w}{\partial \theta_i \partial \theta_j})_{\underline{\theta}^*} \rangle,$$

$$\langle w^2 (\frac{1}{p_w} \frac{\partial^2 p_w}{\partial \theta_i \partial \theta_j})_{\underline{\theta}^*} \rangle = \int w^2 (p \frac{1}{p_w} \frac{\partial^2 p_w}{\partial \theta_i \partial \theta_j})_{\underline{\theta}^*} dx = (\int p w dx) \cdot \frac{\partial^2}{\partial \theta_i \partial \theta_j} (\int w p_w dx) \neq 0 \text{ generally.}]$$

- (c) Generally the weighted covariance matrix in Eq. 56 does *not* equal to the matrix by the standard method, $\underline{\underline{V}}_{\hat{\theta}_w} \neq \underline{\underline{V}}_{\hat{\theta}}$. The difference arises from the inhomogeneity of the weights.

$$\underline{\underline{V}}_{\hat{\theta}_{w \equiv w_0}} = \underline{\underline{V}}_{\hat{\theta}}, \quad \underline{\underline{V}}_{\hat{\theta}_{k \cdot w}} = \underline{\underline{V}}_{\hat{\theta}_w}, \quad k > 0 \text{ constant.}$$

Example:

MC sample: $p(x) = \frac{1+x}{\text{nor}}$, $x \in \mathbf{X}$; $p(x|\theta) = \frac{1+\theta x}{\text{nor}(\theta)}$, $\theta^* = 1$, $w(x) = x$. For different \mathbf{X} the homogeneity of $w(x)$ changes and thus $\sigma_{\hat{\theta}_w} \simeq \sigma_{\hat{\theta}}$ does not hold in the same approximation level:

\mathbf{X}	$\hat{\theta}_{w(x)=x}$	$\hat{\theta}$
[0.00, 0.20]	0.9785 \pm 0.0245	0.9928 \pm 0.0209
[0.07, 0.13]	1.0000 \pm 0.0698	1.0113 \pm 0.0700

Since $\underline{\underline{V}}_{\hat{\theta}}$ is the minimum variance bound, the larger the difference between $\underline{\underline{V}}_{\hat{\theta}_w}$ and $\underline{\underline{V}}_{\hat{\theta}}$ is, the worse is the estimation of this weighted method. Only when

$$\frac{w_i}{w_j} \simeq 1, \quad \forall i \neq j, \quad (60)$$

the estimation is acceptable.

Weights – Application

1. Eliminating the parameter-dependence of the normalization factor:

- If $p(x|\underline{\theta}) = f(\underline{\theta}) P(x|\underline{\theta})$,

$$p(x|\underline{\theta}) = \frac{P(x|\underline{\theta})}{\int_{\mathbf{X}} P(x|\underline{\theta}) dx}. \quad (61)$$

$P(x|\underline{\theta})$, the **Extended p.d.f.**;

$\text{nor}(\underline{\theta}) \equiv \int_{\mathbf{X}} P(x|\underline{\theta}) dx$, the **normalization factor**.

Usually $P(x|\underline{\theta})$ can be easily obtained while $\text{nor}(\underline{\theta})$ needs more complex calculation. A parameter-*independent* nor can be neglected in the minimization procedure

- $(\frac{\partial \ln p(x|\underline{\theta})}{\partial \underline{\theta}} = \frac{\partial \ln P(x|\underline{\theta})}{\partial \underline{\theta}} - \frac{\partial \ln \text{nor}}{\partial \underline{\theta}}, \frac{\partial \ln \text{nor}}{\partial \underline{\theta}} = 0).$

$$p_w(x|\underline{\theta}) = \frac{w(x) P(x|\underline{\theta})}{\int_{\mathbf{X}} w(x) P(x|\underline{\theta}) dx}. \quad (62)$$

$w(x)$ can be suitably chosen so that $\text{nor}_w(\underline{\theta}) \equiv \int_{\mathbf{X}} w(x) P(x|\underline{\theta}) dx$ is independent of $\underline{\theta}$ and thus equivalently $l_w = \sum_i w(x_i) \ln P(x_i|\underline{\theta})$.

Remark.

Since we would not trade "efficiency" (the criterion!) for convenience, we **require that the weights should be as homogenous as possible**.

- Example

MC sample: $\text{Prob}(x) = \frac{1+x}{\text{nor}}$, $x \in \{-0.3, 0.2\}$, $N = 10^6$, which indicates

$$N_{-0.3} = \frac{1-0.3}{1-0.3+1+0.2}N = 368,421, \quad N_{0.2} = \frac{1+0.2}{1-0.3+1+0.2}N = 631,579.$$

Method 1: Simple SML fitting, i.e. without weight:

$$\text{Prob}(x|\theta) = \frac{1+\theta x}{\text{nor}(\theta)}, \quad \theta^* = 1; \quad \hat{\theta} = 1.0000 \pm 0.0035. \quad (63)$$

Method 2: Weighting:

$$(a) \quad \text{nor}_w(\theta) = \sum_x (1+\theta x) w(x) = w(-0.3) + w(0.2) + \theta [-0.3w(-0.3) + 0.2w(0.2)],$$

we can choose $w(-0.3) = 2$, $w(0.2) = 3$ so that $\text{nor}_w = 5$, independent of θ .

$$l_w(\theta) = \sum_{i=1}^N w(x_i) \ln(1+\theta x_i), \quad \theta^* = 1; \quad \hat{\theta}_w = 1.0000 \pm 0.0023. \quad (64)$$

This error is given by $\left(-\frac{\partial^2 l_w}{\partial \theta^2}\right)^{-1}_{\hat{\theta}_w}$, which is *not* the corresponding weighted error (cf. Weights – Analysis Remark. 3a).

(b) Since

$$\text{Prob}_w(\underline{\theta}) \propto (1 + \theta x), \quad \frac{\partial^2 \text{Prob}_w}{\partial \theta^2} = 0, \quad (65)$$

we can use

$$l_{w2}(\theta) = \sum_{i=1}^N w^2(x_i) \ln(1 + \theta x_i) \quad (66)$$

to evaluate the weighted error (cf. Weights – Analysis Remark 3b):

$$\sigma_{\hat{\theta}_w}^2 = \frac{-\frac{\partial^2 l_{w2}}{\partial \theta^2} |_{\hat{\theta}_w}}{\left(-\frac{\partial^2 l_w}{\partial \theta^2} |_{\hat{\theta}_w}\right)^2}, \quad \sigma_{\hat{\theta}_w} = 0.0035. \quad (67)$$

(c) Because the weights $w(-0.3) = 2$, $w(0.2) = 3$ are approximately homogenous, this estimation is acceptable (cf. Weights – Analysis Remark 3c).

(d) To get the correct $\sigma_{\hat{\theta}_w}^2$ (Eq. 67), one needs

- i. "MIGRAD" $-l_w$ to estimate $\hat{\theta}_w$ and record the output error σ_M ;
- ii. Squared the weights;
- iii. "HESSE" $-l_{w^2}$ and record the output error σ_H ;

The MINUIT command

`gMinuit->mnexcm("HESSE", arglist, 0, ierflg)`

calculates the error matrix with the *current* parameter values *without minimization*.

iv. Finally, $\sigma_{\hat{\theta}_w} = \sigma_M / \sigma_H$.

v. For multi-parameter cases, one just needs to use the covariance matrix instead and

$$\text{follow } \underline{V}_{\hat{\theta}} \xrightarrow[N \rightarrow \infty]{P} \left(-\frac{\partial^2 l_w}{\partial \theta_i \partial \theta_j} \Big|_{\hat{\theta}} \right)^{-1} \left(-\frac{\partial^2 l_{w^2}}{\partial \theta_i \partial \theta_j} \Big|_{\hat{\theta}} \right) \left(-\frac{\partial^2 l_w}{\partial \theta_i \partial \theta_j} \Big|_{\hat{\theta}} \right)^{-1}.$$

(e) The condition Eq. 65 is *important*. One can see that by the same procedure as above using another parameterization $\text{Prob}_w(\underline{\theta}) \propto (1 + \theta^2 x) \left(\frac{\partial^2 \text{Prob}_w}{\partial \theta^2} \propto x \neq 0 \right)$, which makes the estimation in Weights – Analysis Remark 3b unjustified. Since the weights are approximately homogeneous, we expect $\sigma_{\hat{\theta}_w} = \sigma_{\hat{\theta}}$ if the weighted covariance matrix is correctly estimated, which however contradicts with the following results:

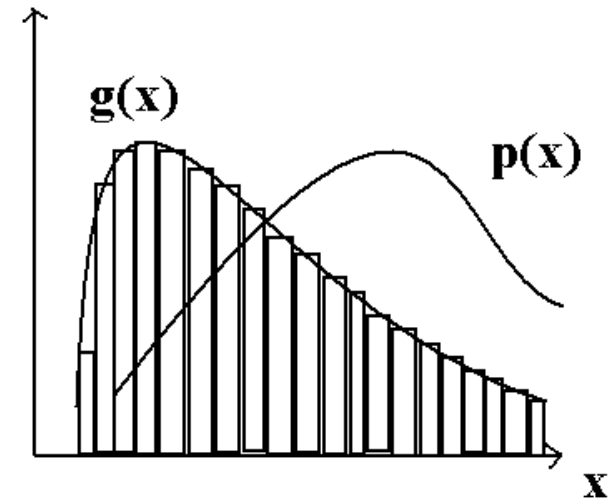
- i. No weighting: $\text{Prob}(x|\theta) = \frac{1+\theta^2 x}{\text{nor}(\theta)}$, $\theta^* = 1$; $\hat{\theta} = 1.0000 \pm 0.0017$.
- ii. Weighting: $l_w(\theta) = \sum_{i=1}^N w(x_i) \ln(1 + \theta^2 x_i)$, $\theta^* = 1$; $\hat{\theta}_w = 1.0000 \pm 0.0014$ (final).

2. Deal with observation efficiency:

- **Efficiency**, the probability that a process with x is observed, $e(x)$.

$p(x|\underline{\theta})$, p.d.f. for the *happened* processes,
 $\int_{\mathbf{X}} p(x|\underline{\theta}) dx = 1$;

$g(x|\underline{\theta})$, p.d.f. for the *observed* processes,
 $g(x|\underline{\theta}) \equiv \frac{e(x)p(x|\underline{\theta})}{\int_{\mathbf{X}} e(x)p(x|\underline{\theta}) dx}$.



(68)

- If we need to **investigate the observed** data sample but g is not known analytically, we can **weight the observed** data to generate a virtual one – the one from **the happened** processes – and then try to obtain information corresponding to the observed.

Let $w(x) = \frac{1}{e(x)}$,

$$g_w = \frac{wep}{\int w e p dx} = p, \quad (69)$$

$$l_w = \sum_{i=1}^N w(x_i) \ln g_w(x_i|\underline{\theta}) = \sum_{i=1}^N w(x_i) \ln p(x_i|\underline{\theta}), \quad (70)$$

$$\underline{V}_{\hat{\underline{\theta}}_w} \xrightarrow[N \rightarrow \infty]{P} \left(-\frac{\partial^2 l_w}{\partial \theta_i \partial \theta_j} \Big|_{\hat{\underline{\theta}}} \right)^{-1} \left(\sum_{k=1}^N w^2(x_k) \left[\frac{\partial \ln p(x_k|\underline{\theta})}{\partial \theta_i} \frac{\partial \ln p(x_k|\underline{\theta})}{\partial \theta_j} \right]_{\hat{\underline{\theta}}} \right) \left(-\frac{\partial^2 l_w}{\partial \theta_i \partial \theta_j} \Big|_{\hat{\underline{\theta}}} \right)^{-1}. \quad (71)$$

Again, this method is **inappropriate if $e(x)$ is severely inhomogeneous**.

The Extended

- If $P(x|\underline{\theta})$ is so chosen that $\mathbb{N}(\underline{\theta}) [\equiv \text{nor}(\underline{\theta}) = \int P dx]$ is the **theoretical number of the observed events**, by assuming that N follows a Poisson distribution,

$$\frac{e^{-\mathbb{N}(\underline{\theta})} [\mathbb{N}(\underline{\theta})]^N}{N!}, \quad (72)$$

we have the **extended likelihood function**,

$$L_{\text{ext}} = \frac{e^{-\mathbb{N}(\underline{\theta})} [\mathbb{N}(\underline{\theta})]^N}{N!} \prod_{i=1}^N p(x_i|\underline{\theta}) = \frac{e^{-\mathbb{N}(\underline{\theta})} [\mathbb{N}(\underline{\theta})]^N}{N!} \prod_{i=1}^N \frac{P(x_i|\underline{\theta})}{\mathbb{N}(\underline{\theta})}, \quad (73)$$

and thus the **extended log-likelihood function** (equivalently),

$$l_{\text{ext}} = \sum_{i=1}^N \ln P(x_i|\underline{\theta}) - \mathbb{N}(\underline{\theta}). \quad (74)$$

- The weighted extended likelihood function:

$$L_{\text{ext},w} = \frac{e^{-\mathbb{N}_w(\underline{\theta})} [\mathbb{N}_w(\underline{\theta})]^{N_w}}{N_w!} \prod_{i=1}^N [p_w(x_i|\underline{\theta})]^{w(x_i)} \quad (75)$$

$$= \frac{e^{-\mathbb{N}_w(\underline{\theta})} [\mathbb{N}_w(\underline{\theta})]^{N_w}}{N_w!} \prod_{i=1}^N \left[\frac{w(x_i) P(x_i|\underline{\theta})}{\mathbb{N}_w(\underline{\theta})} \right]^{w(x_i)}, \quad (76)$$

where $N_w = \sum_{i=1}^N w(x_i)$.

•

$$N \left\langle w \frac{\partial \ln P}{\partial \underline{\theta}} \middle| \underline{\theta}^* \right\rangle - \frac{\partial \mathbb{N}_w}{\partial \underline{\theta}} \bigg|_{\underline{\theta}^*} \xrightarrow[N \rightarrow \infty]{P} 0, \quad (77)$$

$$N \left\langle w \frac{\partial^2 \ln P}{\partial \theta_i \partial \theta_j} \middle| \underline{\theta}^* \right\rangle - \frac{\partial^2 \mathbb{N}_w}{\partial \theta_i \partial \theta_j} \bigg|_{\underline{\theta}^*} \xrightarrow[N \rightarrow \infty]{P} -N \left\langle w \left(\frac{\partial \ln P}{\partial \theta_i} \frac{\partial \ln P}{\partial \theta_j} \right) \middle| \underline{\theta}^* \right\rangle. \quad (78)$$

[Hint: $\frac{N}{\mathbb{N}} - 1 \xrightarrow[N \rightarrow \infty]{P} 0$.]

Remark.

For unweighted case, let $w(x) = 1$.

- Consistency *fails* if

$$\left. \frac{\partial \ln P}{\partial \theta_1} \right|_{\underline{\theta}^*}, \left. \frac{\partial \ln P}{\partial \theta_2} \right|_{\underline{\theta}^*}, \dots, \left. \frac{\partial \ln P}{\partial \theta_M} \right|_{\underline{\theta}^*} \text{ are linearly dependent.} \quad (79)$$

Remark.

Different from the one for the standard method and thus has certain unique applications (cf. Example).

•

$$V_{\hat{\underline{\theta}}_{\text{ext}}} \xrightarrow[N \rightarrow \infty]{P} \frac{1}{N} \left(\left\langle \left. \frac{\partial \ln P}{\partial \theta_i} \frac{\partial \ln P}{\partial \theta_j} \right|_{\underline{\theta}^*} \right\rangle \right)^{-1} = \left(- \left. \frac{\partial^2 l_{\text{ext}}}{\partial \theta_i \partial \theta_j} \right|_{\underline{\theta}^*} \right)^{-1}, \quad (80)$$

$$V_{\hat{\underline{\theta}}_{\text{ext},w}} \xrightarrow[N \rightarrow \infty]{P} \left(- \left. \frac{\partial^2 l_{\text{ext},w}}{\partial \theta_i \partial \theta_j} \right|_{\underline{\theta}^*} \right)^{-1} \left(\sum_{k=1}^N w^2(x_k) \left[\left. \frac{\partial \ln P(x_k)}{\partial \theta_i} \frac{\partial \ln P(x_k)}{\partial \theta_j} \right]_{\underline{\theta}^*} \right) \left(- \left. \frac{\partial^2 l_{\text{ext},w}}{\partial \theta_i \partial \theta_j} \right|_{\underline{\theta}^*} \right)^{-1} \quad (81)$$

$$= \left(- \left. \frac{\partial^2 l_{\text{ext},w}}{\partial \theta_i \partial \theta_j} \right|_{\underline{\theta}^*} \right)^{-1} \left(- \left. \frac{\partial^2 l_{\text{ext},w^2}}{\partial \theta_i \partial \theta_j} \right|_{\underline{\theta}^*} \right) \left(- \left. \frac{\partial^2 l_{\text{ext},w}}{\partial \theta_i \partial \theta_j} \right|_{\underline{\theta}^*} \right)^{-1}. \quad (82)$$

Remark.

$V_{\hat{\underline{\theta}}_{\text{ext},w}}$ is more accessible than $V_{\hat{\underline{\theta}}_w}$ in general cases.

- Example

1. MC sample:

$$p(x) = \frac{1+x}{\text{nor}}, \quad x \in [-0.6, 0.8], \quad N = 10^6. \quad (83)$$

2. Fitting with the standard method, SML,

$$p(x|\theta) = \frac{1+\theta x}{\text{nor}(\theta)}, \quad \theta^* = 1; \quad \hat{\theta} = 1.0004 \pm 0.0026. \quad (84)$$

3. Fitting with the extended method, EML,

$$P(x|\theta_0, \theta_1) = \theta_0 + \theta_1 x, \quad \theta_0^* = \theta_1^* = \frac{N}{\int 1+x dx} = 6.4935 \times 10^5; \quad (85)$$

$$\begin{aligned} \hat{\theta}_0 &= (6.4933 \pm 0.0067) \times 10^5 \\ \hat{\theta}_1 &= (6.4959 \pm 0.0165) \times 10^5, \quad \underline{\underline{V_{\hat{\theta}}}} = \begin{pmatrix} 4.443 \times 10^5 & 1.911 \times 10^5 \\ 1.911 \times 10^5 & 2.738 \times 10^6 \end{pmatrix}; \end{aligned} \quad (86)$$

$$\therefore \frac{\hat{\theta}_1}{\hat{\theta}_0} = 1.0004 \pm 0.0026. \quad (87)$$

Remark.

$$p = \frac{P}{\int P dx}, \quad \text{with } P(x|\theta_0, \theta_1) = \theta_0 + \theta_1 x, \quad (88)$$

is an inappropriate parameterization (Situation c) for the standard method, but not for the extended.

- Comparison between SML and EML:

- Procedure:

- a. Apply SML and EML fits to a MC sample of a fixed / Poisson-fluctuated size.

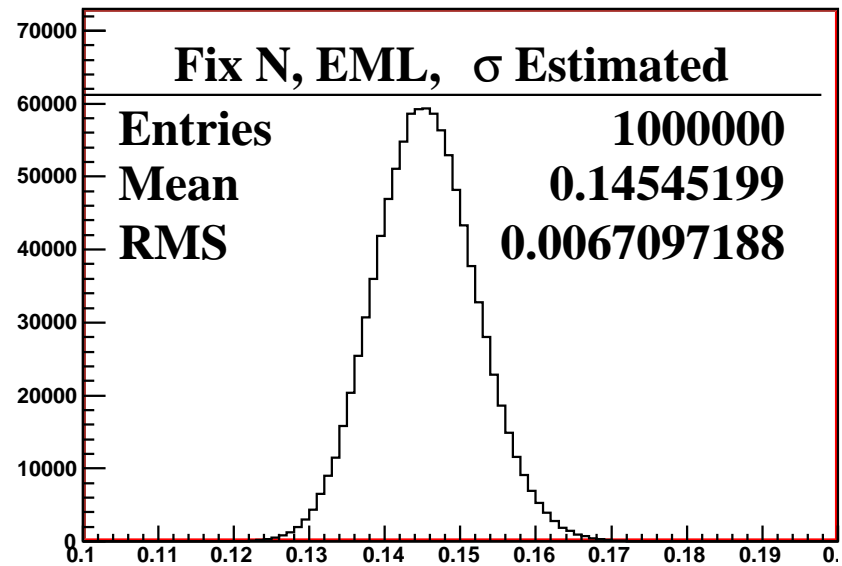
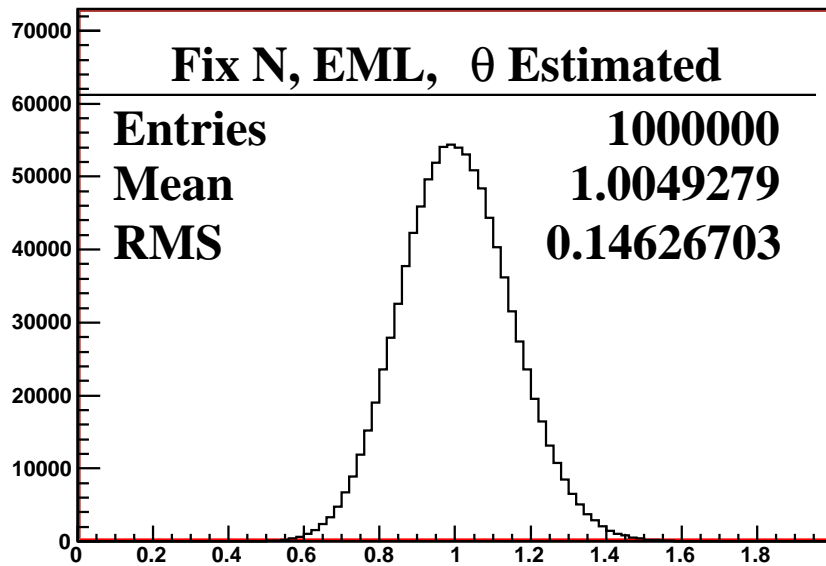
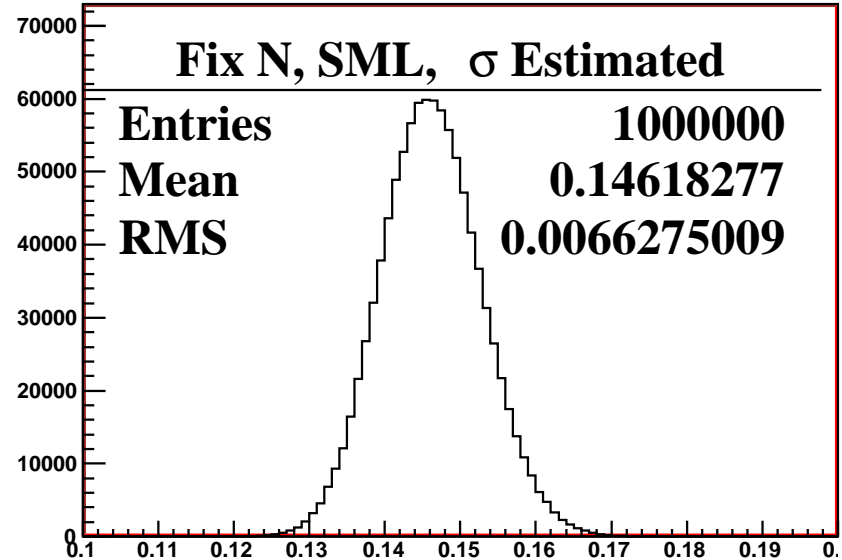
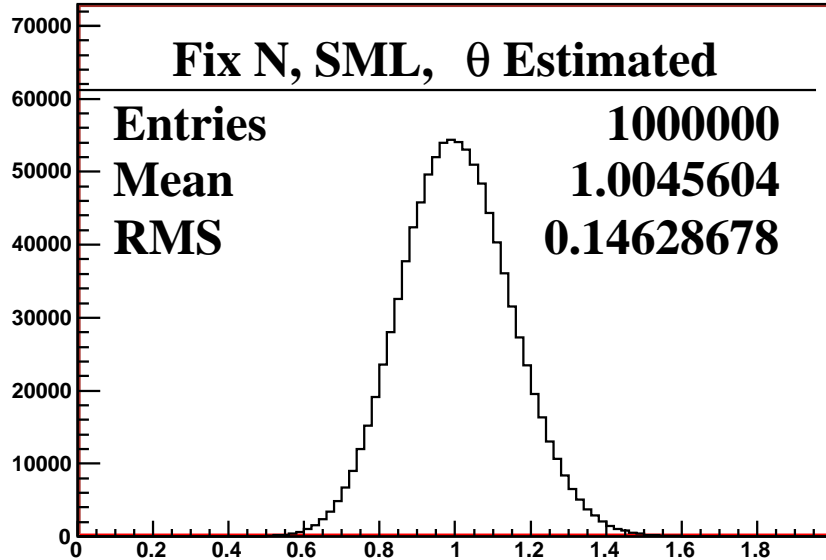
$$\text{MC: } p(x) = \frac{1+x}{\text{nor}}, \quad x \in [-0.3, 0.8], \quad N = 10^3 / \mathbb{N} = 10^3; \quad (89)$$

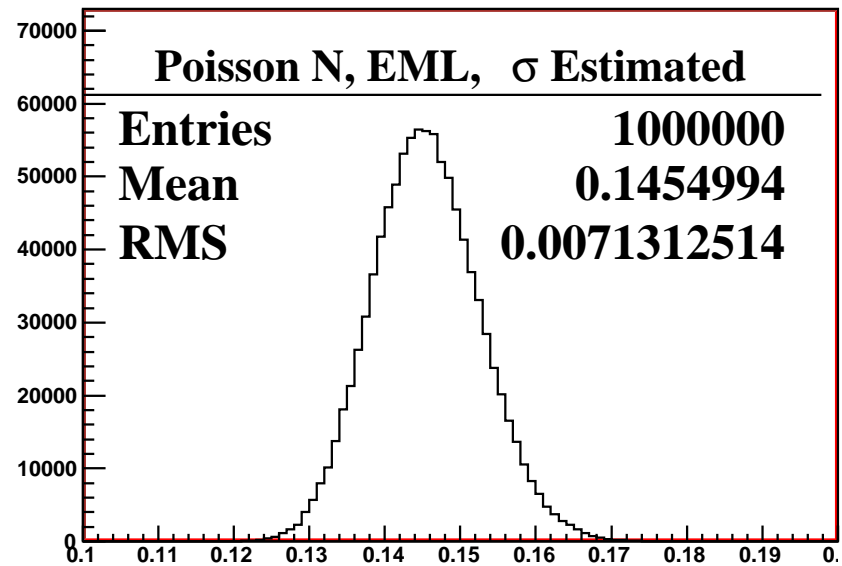
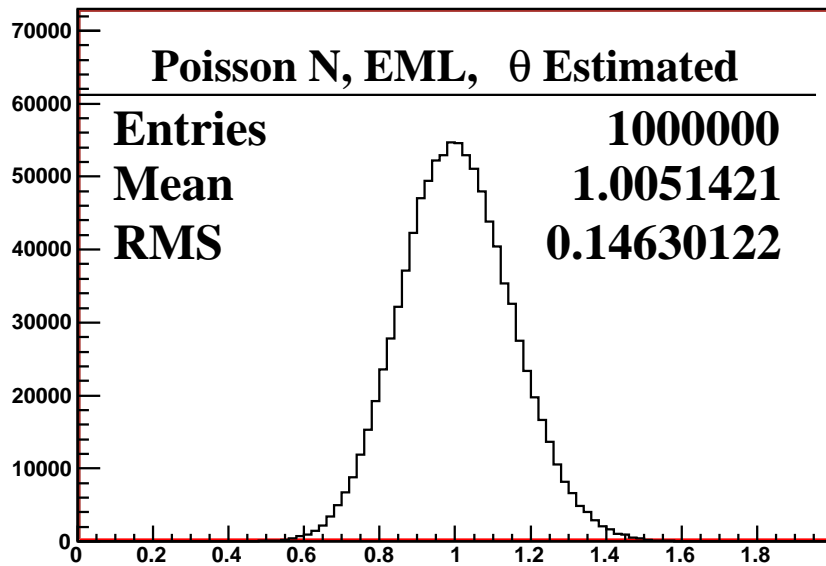
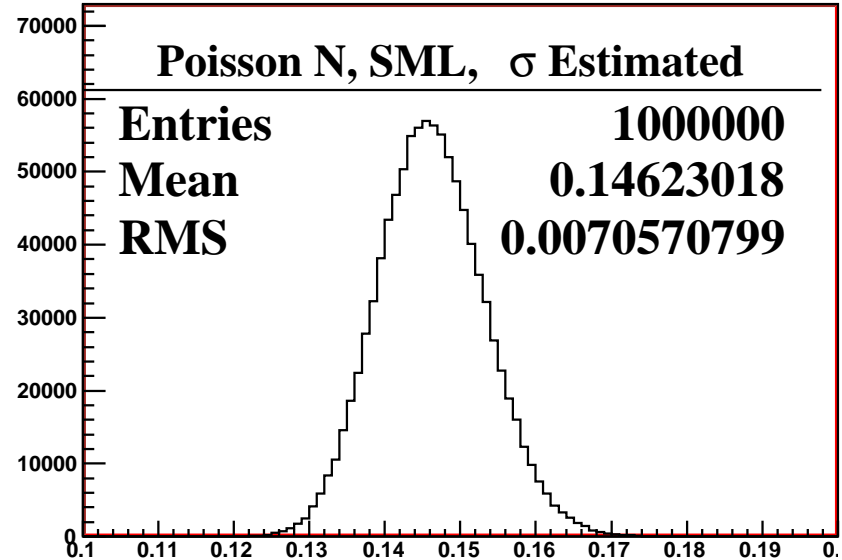
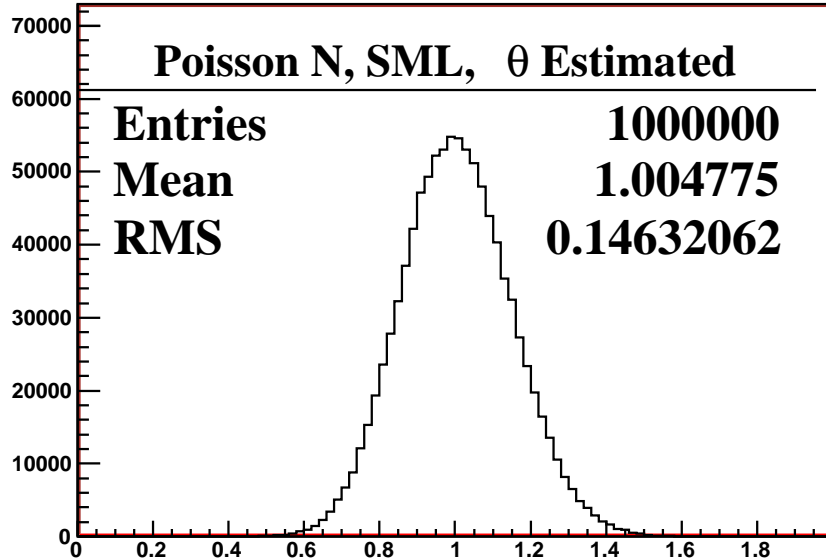
$$\text{SML: } p(x|\theta) = \frac{1+\theta x}{\text{nor}(\theta)}, \quad \theta_{\text{SML}} \equiv \theta, \quad \theta_{\text{SML}}^* = 1; \quad (90)$$

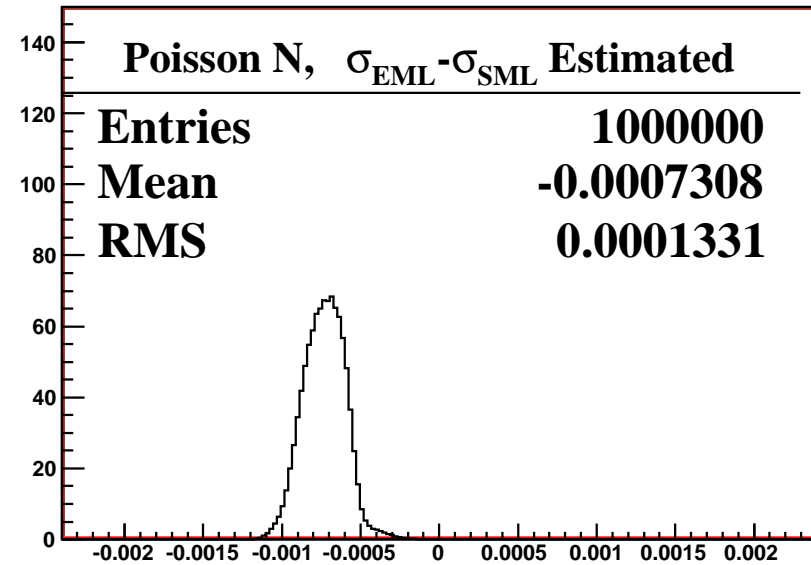
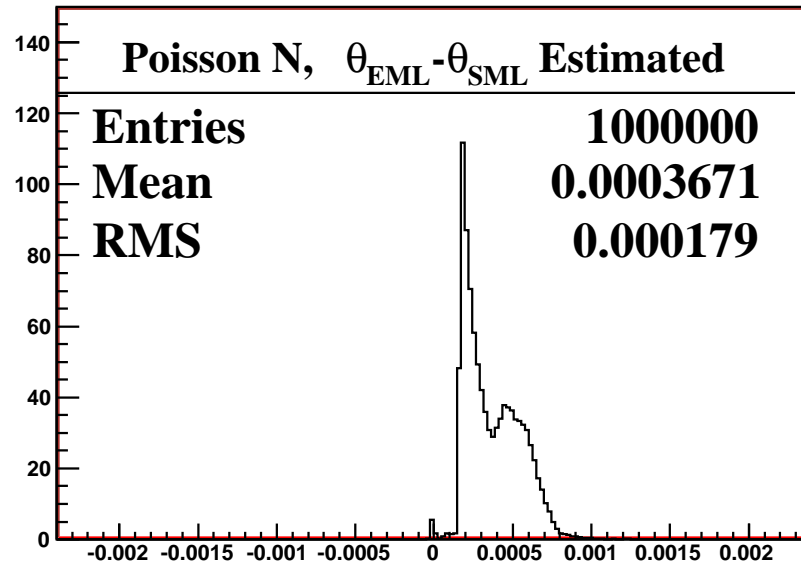
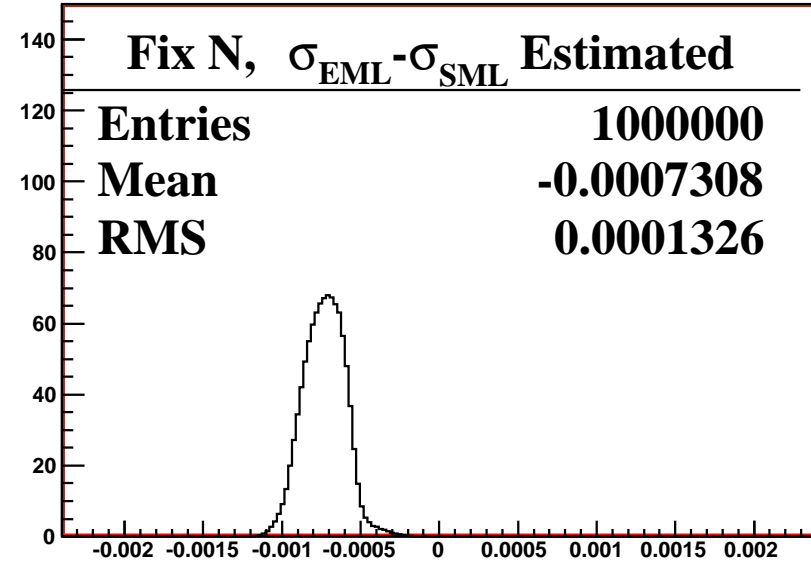
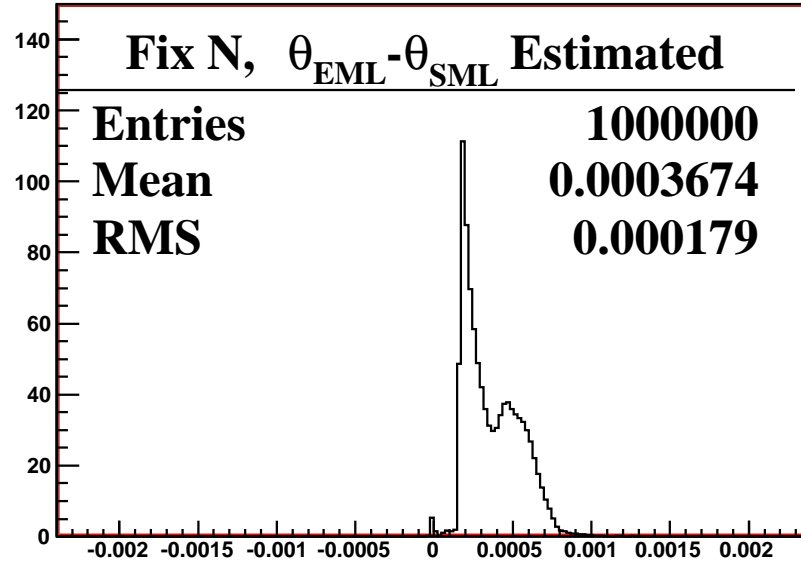
$$\text{EML: } P(x|\theta_0, \theta_1) = \theta_0 + \theta_1 x, \quad \theta_{\text{EML}} \equiv \frac{\theta_1}{\theta_0}, \quad \theta_{\text{EML}}^* = 1. \quad (91)$$

- b. Repeat a. with different random samples (10^6 times for each type of sample).

- c. Compare the distributions of the estimated parameters θ_{SML} , θ_{EML} and the estimated errors.







- Conclusion: SML and EML are consistent to a large extend.

Part 3:

COMBINED ANALYSIS

Combined Analysis – p.d.f.

- **Cross section** for exclusive process with longitudinally polarized beam and unpolarized target:

$$\sigma_{\text{LU}}(\phi, \mathbf{x}) = \sigma_{\text{UU}}^0(\phi, \mathbf{x}) \left[1 + \eta A_{\text{C}}(\phi, \mathbf{x}) + \lambda A_{\text{LU}}^{\text{DVCS}}(\phi, \mathbf{x}) + \eta \lambda A_{\text{LU}}^{\text{T}}(\phi, \mathbf{x}) \right], \quad (92)$$

η , **beam charge**;

λ , **beam pol.**;

ϕ , the angle between the production and scattering planes;

\mathbf{x} , the **set of kinematic variables** $\{-t_c, x_B, Q^2, \dots\}$ *excluding* ϕ .

- The **number density** of exclusive events in the τ - ϕ - \mathbf{x} **space** (not considering the acceptance effect and detection efficiency):

$$n(\tau, \phi, \mathbf{x}) = \mathcal{L}(\tau) \sigma_{\text{LU}}(\phi, \mathbf{x}) = \mathcal{L}(\tau) \sigma_{\text{UU}}^0(\phi, \mathbf{x}) \underline{b}^T(\tau) \underline{a}(\phi, \mathbf{x}), \quad (93)$$

τ , **time**; $\mathcal{L}(\tau)$, **luminosity**;

$$\underline{b}(\tau) \equiv \text{col} [1, \eta(\tau), \lambda(\tau), \eta(\tau)\lambda(\tau)], \quad (94)$$

$$\underline{a}(\phi, \mathbf{x}) \equiv \text{col} \left[1, A_{\text{C}}(\phi, \mathbf{x}), A_{\text{LU}}^{\text{DVCS}}(\phi, \mathbf{x}), A_{\text{LU}}^{\text{T}}(\phi, \mathbf{x}) \right]. \quad (95)$$

- Sum rules for discrete variables.
 x , y are independent variables, coupled as data point (x, y) .

$$\sum_y \sum_x f(x, y) n(x, y) = \sum_{i=1}^N f[(x, y)_i];$$

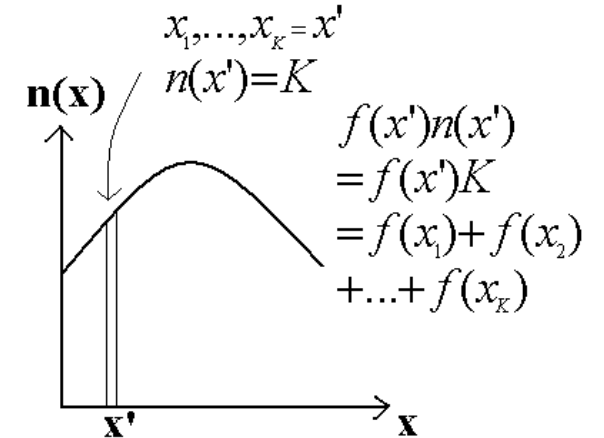
$$\stackrel{\text{OR}}{=} \sum_y \left[\sum_{i=1}^{\mathcal{N}(y)} f(x_i, y) \right] = \sum_y \mathcal{N}(y) \left[\frac{1}{\mathcal{N}(y)} \sum_{i=1}^{\mathcal{N}(y)} f(x_i, y) \right]$$

$$= \sum_y \mathcal{N}(y) \mathbb{f}(y) = \sum_{j=1}^N \mathbb{f}(y_j),$$

$$\mathbb{f}(y) \equiv \frac{1}{\mathcal{N}(y)} \sum_{i=1}^{\mathcal{N}(y)} f(x_i, y), \text{ (x, y are independent)}$$

$$= \frac{1}{\mathfrak{N}} \sum_{i=1}^{\mathfrak{N}} f(x_i, y), \text{ the "x-averaged } f\text{"}.$$

$$\sum_x f(x) n(x) = \sum_{i=1}^{\mathcal{N}} f(x_i);$$



- P.d.f.** for the exclusive process

$$p(\tau, \phi, \mathbf{x}) = \frac{n(\tau, \phi, \mathbf{x})}{\sum_{\tau} \sum_{\phi} \sum_{\mathbf{x}} n(\tau, \phi, \mathbf{x})}, \quad (96)$$

$$\mathbb{p}(\tau, \phi) = \frac{\sum_{\mathbf{x}} n(\tau, \phi, \mathbf{x})}{\sum_{\tau} \sum_{\phi} \sum_{\mathbf{x}} n(\tau, \phi, \mathbf{x})}. \quad (97)$$

For the application of the likelihood method, the parameterization of the p.d.f., which originates from that of $\underline{a}(\phi, \mathbf{x})$, requires that the number density and the total number $\mathbf{N} \equiv \sum_{\tau} \sum_{\phi} \sum_{\mathbf{x}} n(\tau, \phi, \mathbf{x})$ preserve their dependence on $\underline{a}(\phi, \mathbf{x})$. The following is to calculate the analytical form of $\mathbf{N}(\underline{a})$.

1. \mathcal{L} can be evaluated by the number density of the DIS events,

$$n_{\text{DIS}}(\tau) = \mathcal{L}(\tau)\sigma_{\text{DIS}}, \quad (98)$$

σ_{DIS} , the total cross section of the DIS process. Thus

$$n(\tau, \phi, \mathbf{x}) = n_{\text{DIS}}(\tau)r(\phi, \mathbf{x})\underline{b}^T(\tau)\underline{a}(\phi, \mathbf{x}), \quad r(\phi, \mathbf{x}) \equiv \frac{\sigma_{\text{UU}}^0(\phi, \mathbf{x})}{\sigma_{\text{DIS}}}. \quad (99)$$

2. Define **summation operators** with respect to different ranges of time τ ,

$$\underline{\sum}_{\tau} \equiv \begin{pmatrix} \sum_{\forall \tau, \eta(\tau)=+1, \lambda(\tau)>0} \\ \sum_{\forall \tau, \eta(\tau)=+1, \lambda(\tau)<0} \\ \sum_{\forall \tau, \eta(\tau)=-1, \lambda(\tau)>0} \\ \sum_{\forall \tau, \eta(\tau)=-1, \lambda(\tau)<0} \end{pmatrix}. \quad (100)$$

The number densities in ϕ and \mathbf{x} over the corresponding τ ranges,

$$\underline{\mathcal{N}}(\phi, \mathbf{x}) \equiv \underline{\sum}_{\tau} n(\tau, \phi, \mathbf{x}) = \begin{pmatrix} \vec{\mathcal{N}}^+(\phi, \mathbf{x}) \\ \overleftarrow{\mathcal{N}}^+(\phi, \mathbf{x}) \\ \vec{\mathcal{N}}^-(\phi, \mathbf{x}) \\ \overleftarrow{\mathcal{N}}^-(\phi, \mathbf{x}) \end{pmatrix}, \quad (101)$$

$\rightarrow (\leftarrow), \quad + (-)$: positive (negative) beam polarization and charge.

The **beam-state matrix** – a constant matrix encoding the beam-charge and -polarization state, ($\underline{b}(\tau) \equiv \text{col}[1, \eta(\tau), \lambda(\tau), \eta(\tau)\lambda(\tau)]$)

$$\underline{\underline{\mathbb{B}}} \equiv \sum_{\tau} n_{\text{DIS}}(\tau) \underline{b}^T(\tau) = \begin{pmatrix} \vec{N}_{\text{DIS}}^+ & \vec{N}_{\text{DIS}}^+ & \vec{N}_{\text{DIS}}^+ \langle \vec{\lambda}^+ \rangle & \vec{N}_{\text{DIS}}^+ \langle \vec{\lambda}^+ \rangle \\ \overleftarrow{N}_{\text{DIS}}^+ & \overleftarrow{N}_{\text{DIS}}^+ & \overleftarrow{N}_{\text{DIS}}^+ \langle \overleftarrow{\lambda}^+ \rangle & \overleftarrow{N}_{\text{DIS}}^+ \langle \overleftarrow{\lambda}^+ \rangle \\ \vec{N}_{\text{DIS}}^- & -\vec{N}_{\text{DIS}}^- & \vec{N}_{\text{DIS}}^- \langle \vec{\lambda}^- \rangle & -\vec{N}_{\text{DIS}}^- \langle \vec{\lambda}^- \rangle \\ \overleftarrow{N}_{\text{DIS}}^- & -\overleftarrow{N}_{\text{DIS}}^- & \overleftarrow{N}_{\text{DIS}}^- \langle \overleftarrow{\lambda}^- \rangle & -\overleftarrow{N}_{\text{DIS}}^- \langle \overleftarrow{\lambda}^- \rangle \end{pmatrix}. \quad (102)$$

[Hint: $\sum_x n(x) f(x) = \sum_{i=1}^N f(x_i) = N \langle f \rangle$.]

N_{DIS} the number of the corresponding DIS events,

$\langle \lambda \rangle$ the average beam polarization over the DIS events.

Generally $\underline{\underline{\mathbb{B}}}$ is non-singular.

$$\sum_{\tau} n(\tau, \phi, \mathbf{x}) = \sum_{\tau} n_{\text{DIS}}(\tau) r(\phi, \mathbf{x}) \underline{b}^T(\tau) \underline{a}(\phi, \mathbf{x}), \quad (103)$$

$$\underline{\mathcal{N}}(\phi, \mathbf{x}) = \underline{\underline{\mathbb{B}}} r(\phi, \mathbf{x}) \underline{a}(\phi, \mathbf{x}). \quad (104)$$

So the number density summed over the whole range of τ :

$$\underline{\mathcal{N}}(\phi, \mathbf{x}) \equiv \underline{s}^T \underline{\mathcal{N}}(\phi, \mathbf{x}) \quad (105)$$

$$= \underline{s}^T \underline{\underline{\mathbb{B}}} r(\phi, \mathbf{x}) \underline{a}(\phi, \mathbf{x}), \quad \underline{s} \equiv \text{col}(1, 1, 1, 1). \quad (106)$$

3. $r(\phi, \mathbf{x})$ can be solved:

$$\therefore \underline{a}(\phi, \mathbf{x}) \equiv \text{col} \left[1, A_C(\phi, \mathbf{x}), A_{\text{LU}}^{\text{DVCS}}(\phi, \mathbf{x}), A_{\text{LU}}^{\mathcal{I}}(\phi, \mathbf{x}) \right] \quad (107)$$

$$\therefore r(\phi, \mathbf{x}) = [r(\phi, \mathbf{x}) \underline{a}(\phi, \mathbf{x})]_1 \quad (108)$$

$$= \left[\underline{\underline{\mathbb{B}}}^{-1} \underline{\mathcal{N}}(\phi, \mathbf{x}) \right]_1 \quad (109)$$

$$= \sum_{k=1}^4 \left[\underline{\underline{\mathbb{B}}}^{-1} \right]_{1k} [\underline{\mathcal{N}}(\phi, \mathbf{x})]_k. \quad (110)$$

$$\therefore \mathcal{N}(\phi, \mathbf{x}) = \sum_{k=1}^4 \left[\underline{\underline{\mathbb{B}}}^{-1} \right]_{1k} \underline{s}^T \underline{\underline{\mathbb{B}}} \underline{a}(\phi, \mathbf{x}) [\underline{\mathcal{N}}(\phi, \mathbf{x})]_k. \quad (111)$$

4. The total number (with dependence on \underline{a}) over the whole range of τ , ϕ and \mathbf{x} :

$$\mathbf{N} = \sum_{\phi} \sum_{\mathbf{x}} \mathcal{N}(\phi, \mathbf{x}) = \sum_{\phi} \sum_{\mathbf{x}} \sum_{k=1}^4 \left[\underline{\mathbb{B}}^{-1} \right]_{1k} \underline{s}^T \underline{\mathbb{B}} \underline{a}(\phi, \mathbf{x}) [\underline{\mathcal{N}}(\phi, \mathbf{x})]_k \quad (112)$$

$$= \underline{s}^T \underline{\mathbb{B}} \sum_{i=1}^{\mathbf{N}} W_i \underline{a}(\phi_i, \mathbf{x}_i) \stackrel{\text{OR}}{=} \underline{s}^T \underline{\mathbb{B}} \sum_{i=1}^{\mathbf{N}} W_i \underline{a}(\phi_i). \quad (113)$$

$$\underline{a}(\phi) \equiv \frac{1}{\mathfrak{N}} \sum_{j=1}^{\mathfrak{N}} \underline{a}(\phi, \mathbf{x}_j) \quad (114)$$

$$= \text{col} \left[1, \mathbb{A}_{\text{C}}(\phi), \mathbb{A}_{\text{LU}}^{\text{DVCS}}(\phi), \mathbb{A}_{\text{LU}}^{\mathcal{I}}(\phi) \right], \text{ the "}\mathbf{x}\text{-averaged } \underline{a}\text{"}. \quad (115)$$

The upper bound \mathbf{N} is the total number (simply the sample size, no dependence on \underline{a});

$$W_i = \begin{cases} \left[\underline{\mathbb{B}}^{-1} \right]_{11} & (\eta_i = +1, \lambda_i > 0) \\ \left[\underline{\mathbb{B}}^{-1} \right]_{12} & (\eta_i = +1, \lambda_i < 0) \\ \left[\underline{\mathbb{B}}^{-1} \right]_{13} & (\eta_i = -1, \lambda_i > 0) \\ \left[\underline{\mathbb{B}}^{-1} \right]_{14} & (\eta_i = -1, \lambda_i < 0) \end{cases}. \quad (116)$$

□

- Parameterized p.d.f.,

$$p(\tau, \phi, \mathbf{x}; \underline{\theta}) = \frac{n(\tau, \phi, \mathbf{x}; \underline{\theta})}{N(\underline{\theta})} \quad (117)$$

$$= \frac{\mathcal{L}(\tau) \sigma_{UU}^0(\phi, \mathbf{x}) \underline{b}^T(\tau) \underline{a}(\phi, \mathbf{x}; \underline{\theta})}{\underline{s}^T \underline{\mathbb{B}} \sum_{i=1}^N W_i \underline{a}(\phi_i, \mathbf{x}_i; \underline{\theta})} \quad (118)$$

$$\mathbb{p}(\tau, \phi; \underline{\theta}) = \sum_{\mathbf{x}} p(\tau, \phi, \mathbf{x}; \underline{\theta}) \quad (119)$$

$$= \frac{\mathcal{L}(\tau) \bar{\sigma}_{UU}^0(\phi) \underline{b}^T(\tau) \underline{a}'(\phi; \underline{\theta})}{\underline{s}^T \underline{\mathbb{B}} \sum_{i=1}^N W_i \underline{a}(\phi_i; \underline{\theta})} \quad (120)$$

$$\simeq \frac{\mathcal{L}(\tau) \bar{\sigma}_{UU}^0(\phi) \underline{b}^T(\tau) \underline{a}(\phi; \underline{\theta})}{\underline{s}^T \underline{\mathbb{B}} \sum_{i=1}^N W_i \underline{a}(\phi_i; \underline{\theta})}. \quad (121)$$

$$\bar{\sigma}_{UU}^0(\phi) \equiv \sum_{\mathbf{x}} \sigma_{UU}^0(\phi, \mathbf{x}), \quad \underline{a}'(\phi; \underline{\theta}) \equiv \sum_{\mathbf{x}} \frac{\sigma_{UU}^0(\phi, \mathbf{x})}{\bar{\sigma}_{UU}^0(\phi)} \underline{a}(\phi, \mathbf{x}; \underline{\theta}) \simeq \underline{a}(\phi; \underline{\theta}).$$

- For a sample of a certain beam charge,

$$\sigma_{\text{LU}}(\tau, \phi, \mathbf{x}) = \sigma_{\text{UU}}(\phi, \mathbf{x}) [1 + \lambda A_{\text{LU}}(\phi, \mathbf{x})], \quad (122)$$

one just needs to modify the following definitions:

$$\underline{a}(\phi, \mathbf{x}) \equiv \begin{pmatrix} 1 \\ A_{\text{LU}}(\phi, \mathbf{x}) \end{pmatrix}, \quad \underline{\underline{\mathbb{B}}} \equiv \begin{pmatrix} \vec{N}_{\text{DIS}} & \vec{N}_{\text{DIS}} \langle \vec{\lambda} \rangle \\ \overleftarrow{N}_{\text{DIS}} & \overleftarrow{N}_{\text{DIS}} \langle \overleftarrow{\lambda} \rangle \end{pmatrix}, \quad (123)$$

$$\underline{s} \equiv \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad W_i \equiv \begin{cases} \left[\underline{\underline{\mathbb{B}}}^{-1} \right]_{11} & (\lambda_i > 0) \\ \left[\underline{\underline{\mathbb{B}}}^{-1} \right]_{12} & (\lambda_i < 0) \end{cases}, \quad (124)$$

$$\underline{a}(\phi) \equiv \begin{pmatrix} 1 \\ A_{\text{LU}}(\phi) \end{pmatrix}. \quad (125)$$

- Corrections for Detection Efficiencies

Consideration above does not take the detection efficiencies for the exclusive events and the DIS events ($e(\tau, \phi, \mathbf{x})$ and $e_{\text{DIS}}(\tau)$ respectively) into account. The following shows that under certain general conditions, no corrections are needed.

Considering the efficiencies, similar to Eq. (99), we have

$$n(\tau, \phi, \mathbf{x}) = \frac{e(\tau, \phi, \mathbf{x})}{e_{\text{DIS}}(\tau)} n_{\text{DIS}}(\tau) r(\phi, \mathbf{x}) \underline{b}^T(\tau) \underline{a}(\phi, \mathbf{x}), \quad (126)$$

where n and n_{DIS} are the *detected* number densities. If $\frac{e(\tau, \phi, \mathbf{x})}{e_{\text{DIS}}(\tau)}$ depends on τ, ϕ, \mathbf{x} weakly, it drops out when Eq. (126) is normalized.

Combined Analysis – Estimation Methods

- The **likelihood function** and **the extended** one for $\mathbb{P}(\tau, \phi; \underline{\theta})$

$$l = \sum_{i=1}^N \ln \left[\underline{b}_i^T \underline{a}(\phi_i; \underline{\theta}) \right] - N \ln \left[\underline{s}^T \underline{\mathbb{B}} \sum_{i=1}^N W_i \underline{a}(\phi_i; \underline{\theta}) \right], \text{ SML}, \quad (127)$$

$$l_{\text{ext}} = \sum_{i=1}^N \ln \left[\underline{b}_i^T \underline{a}(\phi_i; \underline{\theta}) \right] - \underline{s}^T \underline{\mathbb{B}} \sum_{i=1}^N W_i \underline{a}(\phi_i; \underline{\theta}), \text{ EML}. \quad (128)$$

where $\underline{b}_i \equiv (1, \eta_i, \lambda_i, \eta_i \lambda_i)^T$ and $(1, \lambda_i)^T$ respectively.

- **WUML** (W: weighting, U: unnormalized):

1. Note that $\underline{s} \equiv \text{col}(1, 1, 1, 1)$ is in fact the weighting vector,

$$\underline{s}^T \underline{\underline{\mathbb{B}}} = \underline{s}^T \sum_{\tau} n_{\text{DIS}}(\tau) \underline{b}^T(\tau) = \left(\sum_{\eta(\tau)=+1, \lambda(\tau)>0} + \cdots + \sum_{\eta(\tau)=-1, \lambda(\tau)<0} \right) n_{\text{DIS}}(\tau) \underline{b}^T(\tau). \quad (129)$$

Assigning different \underline{s} components is to weight the data according the beam-state. One can chose \underline{s} so that

$$\underline{s}^T \underline{\underline{\mathbb{B}}} = (k, 0, 0, 0), \quad k \text{ is any positive number}, \quad (130)$$

$$\text{i.e. } \underline{s}_w = \left[\underline{\underline{\mathbb{B}}}^T \right]^{-1} \text{col}(k, 0, 0, 0), \quad (131)$$

making $\underline{N}(\underline{\theta})$ independent of $\underline{\theta}$. $(\underline{N}(\underline{\theta}) = \underline{s}^T \underline{\underline{\mathbb{B}}} \sum_{i=1}^N W_i \underline{a}(\phi_i; \underline{\theta}))$

2. **Warning!!!**: When the weights are inhomogeneous because of, e.g., *extremely* unbalanced data sample of different states, this method is unjustified.

3.

$$l_w = \sum_{i=1}^N w_i \ln \left[\underline{b}_i^T \underline{a}(\phi_i; \underline{\theta}) \right], \text{ WUML}, \quad (132)$$

$$w_i = \begin{cases} s_{w1} & (\eta_i = +1, \lambda_i > 0) \\ s_{w2} & (\eta_i = +1, \lambda_i < 0) \\ s_{w3} & (\eta_i = -1, \lambda_i > 0) \\ s_{w4} & (\eta_i = -1, \lambda_i < 0) \end{cases}. \quad (133)$$

Since

$$\mathbb{p}_w(\tau, \phi | \underline{\theta}) \propto \underline{b}^T(\tau) \underline{a}(\phi; \underline{\theta}), \quad (134)$$

$$(135)$$

with parameterizations like $A = \theta_1 + \theta_2 \sin \phi + \dots$,

$$\frac{\partial^2 \mathbb{p}_w}{\partial \theta_i \partial \theta_j} = 0, \quad (136)$$

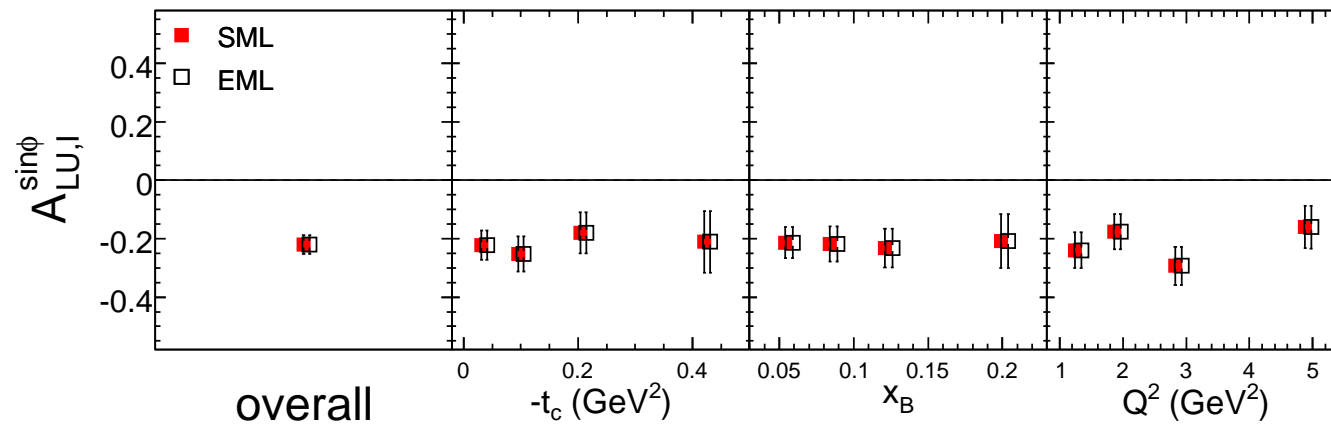
the weighted covariance matrix is accessible.

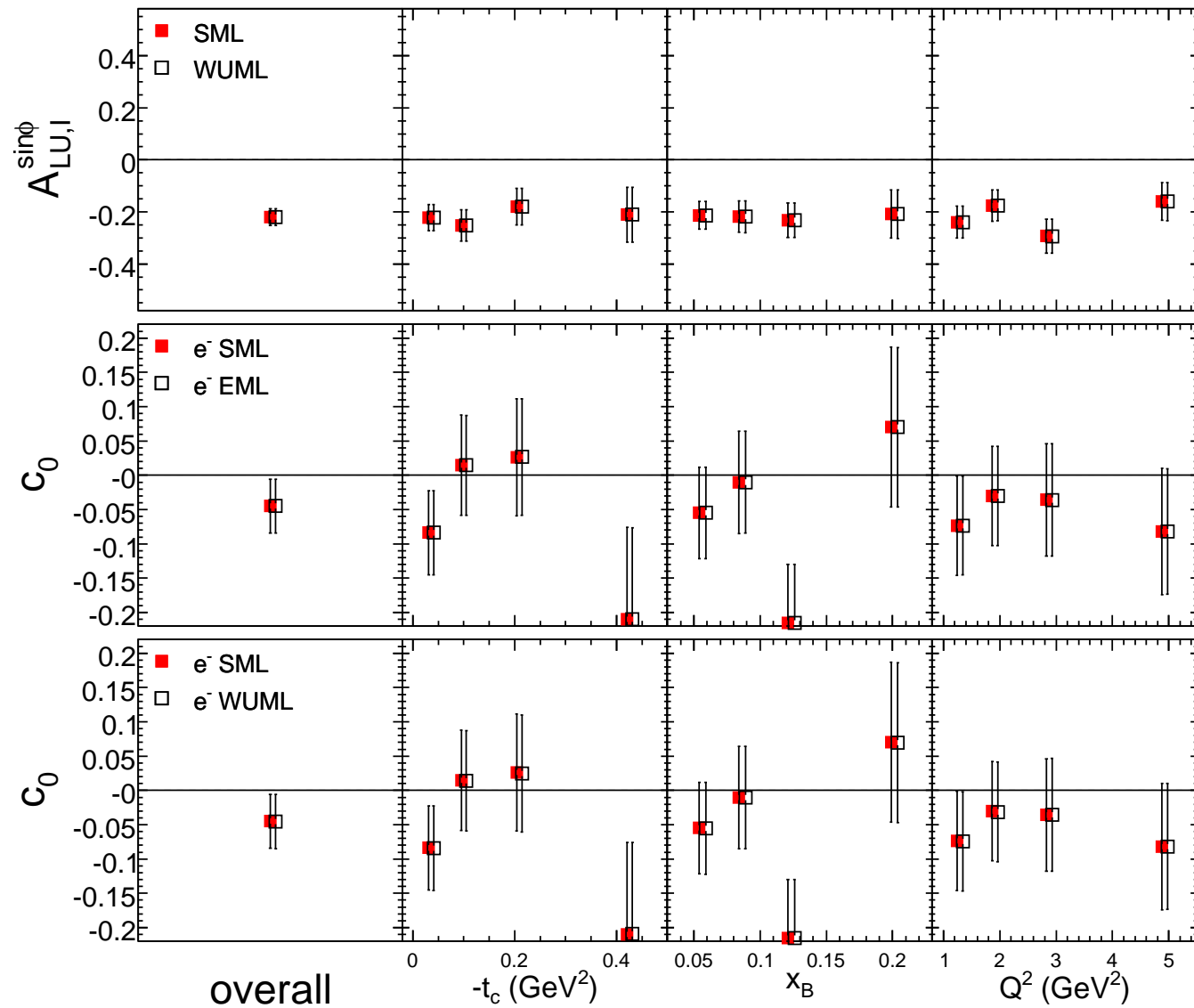
- With parameterizations

$$A_{LU}(\phi; c_0, s_1, c_1) = c_0 + s_1 \sin \phi + c_1 \cos \phi \text{ for BSA analysis;} \quad (137)$$

$$\left. \begin{aligned} A_C(\phi; c_0, s_1, c_1) &= c_0 + s_1 \sin \phi + c_1 \cos \phi; \\ A_{LU}^{DVCS}(\phi; c_0, s_1, c_1) &= c_0 + s_1 \sin \phi + c_1 \cos \phi; \\ A_{LU}^{\mathcal{I}}(\phi; c_0, s_1, c_1) &= c_0 + s_1 \sin \phi + c_1 \cos \phi; \end{aligned} \right\} \text{ for combined analysis,} \quad (138)$$

the 3 methods are applied to the 00d0+05b1 Hydrogen DVCS data, results are *very close*:





Summary for Part 2 & 3

1. Systematically illustrate the analytical properties of the weighting and the extended methods.
2. Provide the p.d.f. for combined analysis, application for SML, EML and WUML are described; results are shown.
3. *WUML may be a good choice for combined analysis since its much easier manipulation, while SML is more accurate in general cases.*

References

- [1] *The advanced theory of statistics*, by M. Kendall et al., 4th ed. of Vol. 2 of the 3-volume ed. (1979) ISBN: 0 85264 255 5.
- [2] A Survey of Maximum Likelihood Estimation, Part 1 & 2, R. H. Norden.
- [3] "Probability", "Statistics" and "Monte Carlo techniques" sections of PDG.
- [4] Statistics, Roger Barlow.
- [5] ANALYSIS OF EXPERIMENTS IN PARTICLE PHYSICS, FRANK T. SOLMITZ.
- [6] Notes on Statistics for Physicists, Revised, Jay Orear.
- [7] MINUIT Reference Manual, Version 94.1, F. James.
- [8] The Interpretation of Errors, F. James.
- [9] MINUIT User's Guide, F. James.
- [10] MINUIT Tutorial, F. James.

Appendix: Schwarz' Inequality (Generalized)

Let $\underline{a}(x_1, x_2, \dots, x_N)$, $\underline{b}(x_1, x_2, \dots, x_N)$ be real vector functions with M components, if $\langle \underline{b} \underline{b}^T \rangle$ is nonsingular,

$$\langle \underline{a} \underline{a}^T \rangle \geq \langle \underline{a} \underline{b}^T \rangle \langle \underline{b} \underline{b}^T \rangle^{-1} \langle \underline{b} \underline{a}^T \rangle. \quad (139)$$

$$(\underline{B} \geq \underline{C} : \underline{B} - \underline{C} \text{ positive semidefinite.})$$

Proof. Let \underline{A} be a constant real $M \times M$ matrix,

$$\begin{aligned} \forall \underline{v} \neq \underline{0}, \quad & \underline{v}^T \left\langle \left(\underline{a} + \underline{A} \underline{b} \right) \left(\underline{a} + \underline{A} \underline{b} \right)^T \right\rangle \underline{v} \\ &= \left\langle \underline{v}^T \left(\underline{a} + \underline{A} \underline{b} \right) \left(\underline{a} + \underline{A} \underline{b} \right)^T \underline{v} \right\rangle \\ &= \left\langle \left[\underline{v}^T \left(\underline{a} + \underline{A} \underline{b} \right) \right]^2 \right\rangle \geq 0, \end{aligned} \quad \begin{aligned} \therefore \quad & \underline{0} \leq \left\langle \left(\underline{a} + \underline{A} \underline{b} \right) \left(\underline{a} + \underline{A} \underline{b} \right)^T \right\rangle \\ &= \langle \underline{a} \underline{a}^T \rangle + \langle \underline{a} \underline{b}^T \rangle \underline{A}^T + \underline{A} \langle \underline{b} \underline{a}^T \rangle + \underline{A} \langle \underline{b} \underline{b}^T \rangle \underline{A}^T, \\ \text{let } \underline{A} &= - \langle \underline{a} \underline{b}^T \rangle \langle \underline{b} \underline{b}^T \rangle^{-1}, \\ \text{we have} \quad & \langle \underline{a} \underline{a}^T \rangle - \langle \underline{a} \underline{b}^T \rangle \langle \underline{b} \underline{b}^T \rangle^{-1} \langle \underline{b} \underline{a}^T \rangle \geq \underline{0}. \end{aligned}$$

" = " holds Iff $\underline{v}^T \left(\underline{a} + \underline{A} \underline{b} \right) \equiv 0$,
i.e. $a_i + (\underline{A} \underline{b})_i$ are linearly dependent
functions of (x_1, x_2, \dots, x_N) .

□