



Optimal Estimation Methods

(Lecture 7 – Minimum Variance Estimation & Cramér-Rao Inequality)

Dr. John L. Crassidis

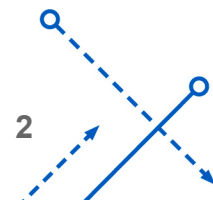
University at Buffalo – State University of New York
Department of Mechanical & Aerospace Engineering
Amherst, NY 14260-4400

johnc@buffalo.edu

<http://www.buffalo.edu/~johnc>

- Previously stated that in many cases we wish to weight different measurements differently
 - We now derive the “optimal” weighting matrix based on probability
- Two main approaches shown here
 - Minimum Variance Estimation
 - Maximum Likelihood Estimation
- Note there are others, such as Minimum Risk
- Two main types of estimators (we’ll derive both)
 - Without *a priori* estimates
 - With *a priori* estimates
- Also, we’ll derive the covariance of the estimation errors
 - Discuss the Cramér-Rao lower bound too

↖ Very Important



- Consider case without *a priori* estimates first \mathcal{V} - zero mean

- Assume a linear observation model

$$\begin{matrix} (m \times 1) \\ \tilde{\mathbf{y}} \end{matrix} = \begin{matrix} (m \times n) \\ H \end{matrix} \begin{matrix} (n \times 1) \\ \mathbf{x} \end{matrix} + \begin{matrix} (m \times 1) \\ \mathbf{v} \end{matrix}$$

Measurement Error
with $E\{\mathbf{v}\} = \mathbf{0}$ and
 $E\{\mathbf{v}\mathbf{v}^T\} = R$

- We desire an estimate as a linear combination of the measurements

$$\begin{matrix} (n \times 1) \\ \hat{\mathbf{x}} \end{matrix} = \begin{matrix} (n \times m) \\ M \end{matrix} \begin{matrix} (m \times 1) \\ \tilde{\mathbf{y}} \end{matrix} + \begin{matrix} (n \times 1) \\ \mathbf{n} \end{matrix}$$

- The minimum variance definition of “optimum” M and \mathbf{n} is that the variance of n estimates from their respective “true” values is minimized

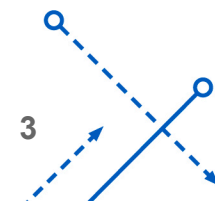
$E\{(\hat{x} - x)^2\}$ - Variance Definition

- Leads to the following loss function

$$J_i = \frac{1}{2} E \left\{ (\hat{x}_i - x_i)^2 \right\}, \quad i = 1, 2, \dots, n$$

- This clearly requires n minimizations depending upon the same M and \mathbf{n}

- Let's prove that the “uncoupled” loss function is valid



- The linear model **must** also be true when no measurement errors exist, so in this case we have $(v=0)$

$$\tilde{\mathbf{y}} \equiv \mathbf{y} = H\mathbf{x}$$

- An obvious requirement upon the desired estimator is that perfect measurements should result (if a solution is possible) when $\hat{\mathbf{x}} = \mathbf{x} = \text{true state}$
- This requirement can be written by substituting $\hat{\mathbf{x}} = \mathbf{x}$ and $\tilde{\mathbf{y}} = H\mathbf{x}$ into the linear measurement model, which gives

$$\mathbf{x} = MH\mathbf{x} + \mathbf{n}$$

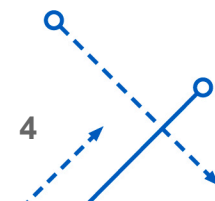
- Thus we conclude that

$$\mathbf{n} = \mathbf{0} \quad \text{and} \quad MH = I$$

- Note that $MH = I$ will also be shown for unbiased estimates
- The desired estimator then has the form

$$\hat{\mathbf{x}} = M\tilde{\mathbf{y}}$$

- Need to now find M



- The unknown M -matrix is partitioned by rows as

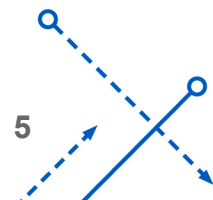
$$M = \begin{bmatrix} M_1 \\ M_2 \\ \vdots \\ M_n \end{bmatrix}, \quad M_i \equiv \{M_{i1} \ M_{i2} \ \cdots \ M_{im}\}$$

or

$$M^T = [M_1^T \ M_2^T \ \cdots \ M_n^T]$$

- The identity matrix can be partitioned by rows and columns as

$$I = \begin{bmatrix} I_1^r \\ I_2^r \\ \vdots \\ I_n^r \end{bmatrix} = [I_1^c \ I_2^c \ \cdots \ I_n^c], \quad \text{note } I_i^r = (I_i^c)^T$$



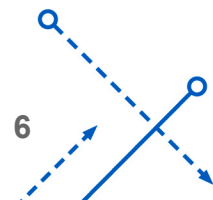
- The constraint $MH = I$ can now be written as

$$\begin{aligned} \text{or} \quad H^T M_i^T &= I_i^c, \quad i = 1, 2, \dots, n \\ M_i H &= I_i^r, \quad i = 1, 2, \dots, n \end{aligned}$$

- The i^{th} element of the estimate is given by

$$\hat{x}_i = M_i \tilde{\mathbf{y}}, \quad i = 1, 2, \dots, n \quad (1)$$

- A glance at this equation reveals that the i^{th} element of the estimate depends only upon the elements of M contained in the i^{th} row
- A similar statement holds for the constraint equations
 - The elements of the i^{th} row are independently constrained
- This “uncoupled” nature is the key feature which allows one to carry out the n separate minimizations of the loss function shown before
 - We will show another approach later that does not need this feature per se



- Substituting Eq. (1) into the loss function gives

$$\tilde{y} = Hx + v$$

$$J_i = \frac{1}{2} E \left\{ (M_i \tilde{y} - x_i)^2 \right\}, \quad i = 1, 2, \dots, n$$

- Substituting the measurement equation gives

$$J_i = \frac{1}{2} E \left\{ (M_i H x + M_i v - x_i)^2 \right\}, \quad i = 1, 2, \dots, n$$

- Incorporating the constraint gives

$$M_i H = I_i^r$$

$$J_i = \frac{1}{2} E \left\{ (I_i^r x + M_i v - x_i)^2 \right\}, \quad i = 1, 2, \dots, n$$

- Noting $I_i^r x = x_i$ gives simply

x_i 's cancel

$$J_i = \frac{1}{2} E \left\{ (M_i v)^2 \right\}, \quad i = 1, 2, \dots, n$$

$$= \frac{1}{2} E \left\{ M_i (v v^T) M_i^T \right\}, \quad i = 1, 2, \dots, n$$

$$E(v v^T) = R$$

As stated previously,
zero mean

$$E[(v-0)(v-0)^T]$$

- Note: the only random variable on the right-hand side is v

M_i can be pulled out.



- Assuming that \mathbf{v} has zero mean and $\text{cov}\{\mathbf{v}\} \equiv R = E\{\mathbf{v}\mathbf{v}^T\}$ gives

$$J_i = \frac{1}{2} M_i R M_i^T, \quad i = 1, 2, \dots, n$$

- Need to also account for constraint equations
 - Use the Lagrange multiplier approach

$$\frac{\partial}{\partial x} (x C x^T) = (C + C^T)x$$

$$C + C^T = 2C$$

\uparrow
 $I +$
 C is
 symmetric

$$J_i = \frac{1}{2} M_i R M_i^T + \underbrace{\lambda_i^T (I_i^c - H^T M_i^T)}_{\text{Constraint}}, \quad i = 1, 2, \dots, n$$

\uparrow
 symmetric

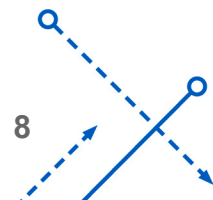
where

$$\lambda_i^T = \{\lambda_{1_i}, \lambda_{2_i}, \dots, \lambda_{n_i}\}$$

- The necessary conditions give

$$\nabla_{M_i^T} J_i = R M_i^T - H \lambda_i = \mathbf{0}, \quad i = 1, 2, \dots, n \quad (2)$$

$$\nabla_{\lambda_i} J_i = I_i^c - H^T M_i^T = \mathbf{0}, \text{ or } M_i H = I_i^r, \quad i = 1, 2, \dots, n \quad (3)$$



- From Eq. (2) we have

Assume positive definite



$$M_i = \lambda_i^T H^T R^{-1}, \quad i = 1, 2, \dots, n$$

- Substituting this equation into Eq. (3) gives

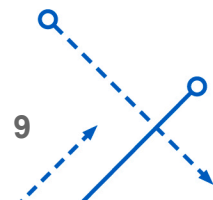
$$\lambda_i^T = I_i^r (H^T R^{-1} H)^{-1}$$

- Substituting this equation into M_i gives

$$M_i = I_i^r (H^T R^{-1} H)^{-1} H^T R^{-1}, \quad i = 1, 2, \dots, n$$

- It then follows that

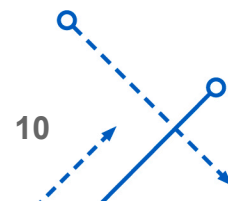
$$M = (H^T R^{-1} H)^{-1} H^T R^{-1}$$



- Substituting M into $\hat{\mathbf{x}} = M\tilde{\mathbf{y}}$ gives

$$\hat{\mathbf{x}} = (H^T R^{-1} H)^{-1} H^T R^{-1} \tilde{\mathbf{y}}$$

- This is referred to as the Gauss-Markov Theorem
- Some observations
 - The minimal variance estimator is identical to the least squares estimator provided that the weight matrix is identified as the inverse of the observation error covariance *Optimal for this loss function*
 - Also, the “sequential least squares estimation” results are seen to embody a special case “sequential minimal variance estimation”
 - It is simply necessary to employ R^{-1} as W in the sequential least squares formulation
 - But we still require R^{-1} to have the block diagonal structure assumed for W



- Another approach

- Define the error covariance matrix for an unbiased estimator

$$P = E \{ (\hat{\mathbf{x}} - \mathbf{x})(\hat{\mathbf{x}} - \mathbf{x})^T \} \quad \leftarrow \text{variance of estimation error}$$

- Minimum variance estimation is equivalent to minimizing the trace of P
- Need to also satisfy constraint $MH = I$
- Use method of Lagrange multipliers to append loss function

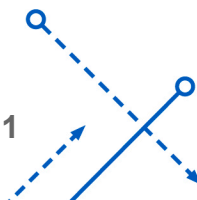
$$J = \frac{1}{2} \text{Tr} [E \{ (\hat{\mathbf{x}} - \mathbf{x})(\hat{\mathbf{x}} - \mathbf{x})^T \}] + \text{Tr} [\Lambda(I - MH)]$$

where Λ is a matrix of Lagrange multipliers

- Note: covariance can also be found using *Parallel Axis Theorem* for an unbiased estimate

x is not a random variable, put out of a interval of truth

$$\begin{aligned} E \{ (\hat{\mathbf{x}} - \mathbf{x})(\hat{\mathbf{x}} - \mathbf{x})^T \} &= E \{ \hat{\mathbf{x}} \hat{\mathbf{x}}^T \} - E \{ \mathbf{x} \} E \{ \mathbf{x} \}^T \\ &= E \{ \hat{\mathbf{x}} \hat{\mathbf{x}}^T \} - \mathbf{x} \mathbf{x}^T \end{aligned}$$



- We have

$$\begin{aligned}
 P &= E \{ \hat{\mathbf{x}} \hat{\mathbf{x}}^T \} - \mathbf{x} \mathbf{x}^T \\
 &= E \{ M \tilde{\mathbf{y}} \tilde{\mathbf{y}}^T M^T \} - \mathbf{x} \mathbf{x}^T \\
 &= E \{ (M H \mathbf{x} + M \mathbf{v})(M H \mathbf{x} + M \mathbf{v})^T \} - \mathbf{x} \mathbf{x}^T
 \end{aligned}$$

Do this

- Now use $E\{\mathbf{v}\} = \mathbf{0}$ and $E\{\mathbf{v} \mathbf{v}^T\} = R$

$$P = M R M^T + M H \mathbf{x} \mathbf{x}^T H^T M^T - \mathbf{x} \mathbf{x}^T$$

- Noting that $M H = I$ leads to

$$P = M R M^T$$

- Therefore, the loss function becomes

$$J = \frac{1}{2} \text{Tr}(M R M^T) + \text{Tr} [\Lambda(I - M H)]$$

- Again, the goal is to find M that minimizes J



- Consider the following useful identities

$$\frac{\partial}{\partial A} \text{Tr}(BAC) = B^T C^T$$

$$\frac{\partial}{\partial A} \text{Tr}(ABA^T) = A(B + B^T)$$

- Thus, we have the following necessary conditions

$$\nabla_M J = MR - \Lambda^T H^T = 0$$

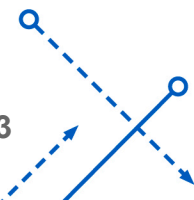
$$\nabla_\Lambda J = I - MH = 0$$

- Two equations for M and Λ^T
- Solving the first equation for M gives

$$M = \Lambda^T H^T R^{-1}$$

- Substituting this into the second equation gives

$$\Lambda^T = (H^T R^{-1} H)^{-1}$$



- Note that Λ is a symmetric matrix
 - It also has a physical interpretation
 - This is equivalent to the error-covariance matrix, which tells us about the quality of the estimate
 - The “larger” its value, the worse the estimate will be
 - This will be discussed in detail later
- Substituting Λ^T back into M gives

$$M = (H^T R^{-1} H)^{-1} H^T R^{-1}$$

- This gives exactly the same solution as before
 - Note that the “decoupling” assumptions are actually in the loss function
 - We choose to minimize the trace of the covariance, which ignores the correlations (this is the decoupling)
 - Other possible forms for the loss function can be chosen, such as minimizing the infinity norm

Unbiased Estimates (i)

\hat{x} is random variable because it depends on \tilde{y} which depends on random variables

- An estimator $\hat{x}(\tilde{y})$ is said to be an “unbiased estimator” of \mathbf{x} if $E \{ \hat{x}(\tilde{y}) \} = \mathbf{x}$ for every possible value of \mathbf{x}
 - If \hat{x} is biased, the difference $E \{ \hat{x}(\tilde{y}) \} - \mathbf{x}$ is called the “bias” of \hat{x}
- Go back to the previous estimator form

$$\hat{\mathbf{x}} = M\tilde{\mathbf{y}}$$

$$\tilde{\mathbf{y}} = H\mathbf{x} + \mathbf{v}$$

$$= MH\mathbf{x} + M\mathbf{v}$$

$$MH E(\tilde{\mathbf{y}}) = \mathbf{x} \quad M E(\mathbf{v}) = 0$$

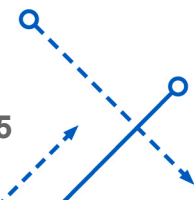
$$E(\hat{\mathbf{x}}) = \mathbf{x} = E(MH\mathbf{x}) + E(M\mathbf{v}) = MH\mathbf{x} + 0$$

- Taking the expectation of both sides and assuming zero-mean measurement error, so that $E\{\mathbf{v}\} = \mathbf{0}$, gives

$$E \{ \hat{\mathbf{x}} \} = MH\mathbf{x}$$

$$\mathbf{x} = MH\mathbf{x} \Rightarrow MH = I$$

- Thus for an unbiased estimate we must have $MH = I$
- Same result as before!



- Sample Variance Example with Data $\{\tilde{y}(t_1), \tilde{y}(t_2), \dots, \tilde{y}(t_m)\}$
 - Compute sample variance using

$$\hat{\sigma}^2 = \frac{1}{m-1} \sum_{i=1}^m [\tilde{y}(t_i) - \hat{\mu}]^2$$

- Note, many calculators give the option of dividing by m or $m-1$
- Check to see if this estimate is unbiased using $m-1$
- Defining $E\{\hat{\sigma}^2\} \equiv S^2$ with $\hat{\mu} = \frac{1}{m} \sum_{i=1}^m \tilde{y}(t_i)$ gives

$$\begin{aligned} S^2 &= \frac{1}{m-1} E \left\{ \left[\sum_{i=1}^m \tilde{y}^2(t_i) - 2\tilde{y}(t_i)\hat{\mu} + \hat{\mu}^2 \right] \right\} \\ &= \frac{1}{m-1} \left[E \left\{ \sum_{i=1}^m \tilde{y}^2(t_i) \right\} - \frac{2}{m} E \left\{ \sum_{i=1}^m \tilde{y}(t_i) \left[\sum_{i=1}^m \tilde{y}(t_i) \right] \right\} + \frac{1}{m^2} E \left\{ \sum_{i=1}^m \left[\sum_{i=1}^m \tilde{y}(t_i) \right]^2 \right\} \right] \\ &= \frac{1}{m-1} \left[\sum_{i=1}^m E \{ [\tilde{y}(t_i)]^2 \} - \frac{2}{m} E \left\{ \left[\sum_{i=1}^m \tilde{y}(t_i) \right]^2 \right\} + \frac{m}{m^2} E \left\{ \left[\sum_{i=1}^m \tilde{y}(t_i) \right]^2 \right\} \right] \\ &= \frac{1}{m-1} \left[\sum_{i=1}^m E \{ [\tilde{y}(t_i)]^2 \} - \frac{1}{m} E \left\{ \left[\sum_{i=1}^m \tilde{y}(t_i) \right]^2 \right\} \right] \end{aligned}$$



- For any random variable z the variance is computed from (using the parallel axis theorem) $\text{var}\{z\} = E\{z^2\} - E\{z\}^2$
- Then applying the variance equation gives

$$\begin{aligned}
 S^2 &= \frac{1}{m-1} \left[\sum_{i=1}^m (\sigma^2 + \mu^2) - \frac{1}{m} \left\{ \text{var} \left[\sum_{i=1}^m \tilde{y}(t_i) \right] + \left[E \left\{ \sum_{i=1}^m \tilde{y}(t_i) \right\} \right]^2 \right\} \right] \\
 &= \frac{1}{m-1} \left[m\sigma^2 + m\mu^2 - \frac{1}{m}m\sigma^2 - \frac{1}{m}m^2\mu^2 \right] \\
 &= \frac{1}{m-1} [m\sigma^2 - \sigma^2] \\
 &= \sigma^2
 \end{aligned}$$

Divide by $\frac{1}{m}$ gives
biased estimator

- Therefore, this estimator is unbiased
- However, the sample variance shown in this example does not give an estimate with the smallest mean-square-error for Gaussian (normal) distributions



- A more general definition for an unbiased estimator is

$$E \{ \hat{\mathbf{x}}_k(\tilde{\mathbf{y}}) \} = \mathbf{x} \quad \text{for all } k$$

- For the sequential estimator we wish to have the form

$$\hat{\mathbf{x}}_{k+1} = G_{k+1} \hat{\mathbf{x}}_k + K_{k+1} \tilde{\mathbf{y}}_{k+1}$$

where G_{k+1} and K_{k+1} are deterministic matrices

- Substituting the measurement equation at $k+1$ gives

$$\hat{\mathbf{x}}_{k+1} = G_{k+1} \hat{\mathbf{x}}_k + K_{k+1} H_{k+1} \mathbf{x}_{k+1} + K_{k+1} \mathbf{v}_{k+1}$$

- Taking the expectation gives

$$E \{ \hat{\mathbf{x}}_{k+1} \} = G_{k+1} E \{ \hat{\mathbf{x}}_k \} + K_{k+1} H_{k+1} \mathbf{x}_{k+1}$$

- Noting that the unbiased condition must be valid for all k leads to

$$G_{k+1} = I - K_{k+1} H_{k+1}$$

Then

$$\hat{\mathbf{x}}_{k+1} = \hat{\mathbf{x}}_k + K_{k+1} (\tilde{\mathbf{y}}_{k+1} - H_{k+1} \hat{\mathbf{x}}_k)$$

- This is exactly the sequential process
 - We have now shown that it produces unbiased estimates though

- Consider case with *a priori* estimates now

v & w are both zero mean

- Measurement model is same as before

$$\tilde{\mathbf{y}} = H\mathbf{x} + \mathbf{v}, \quad \text{with } E\{\mathbf{v}\} = \mathbf{0} \quad \text{and} \quad E\{\mathbf{v}\mathbf{v}^T\} = R$$

- Now consider an *a priori* estimate with model given by

$$\hat{\mathbf{x}}_a = \mathbf{x} + \mathbf{w}, \quad \text{with } E\{\mathbf{w}\} = \mathbf{0} \quad \text{and} \quad E\{\mathbf{w}\mathbf{w}^T\} = Q$$

- We also assume that the measurement errors and *a priori* errors are uncorrelated so that $E\{\mathbf{w}\mathbf{v}^T\} = E\{\mathbf{v}\mathbf{w}^T\} = \mathbf{0}$

- We desire to estimate \mathbf{x} as a linear combination of the measurements and *a priori* estimates as

$$\hat{\mathbf{x}} = M\tilde{\mathbf{y}} + N\hat{\mathbf{x}}_a + \mathbf{n}$$

$$E[M\mathbf{v}] = E[N\mathbf{w}] = \mathbf{0}$$

- For unbiased estimates we require

$$E\{\hat{\mathbf{x}}\} = E\{M(H\mathbf{x} + \mathbf{v})\} + E\{N(\mathbf{x} + \mathbf{w}) + \mathbf{n}\} = (MH + N)\mathbf{x} + \mathbf{n} = \mathbf{x}$$

- Then $\mathbf{n} = \mathbf{0}$ and $MH + N = I$ is required for an unbiased estimate
- So the actual form is given by

$$\hat{\mathbf{x}} = M\tilde{\mathbf{y}} + N\hat{\mathbf{x}}_a, \quad \text{subject to } MH + N = I$$



- Loss function for this case becomes

$$J = \frac{1}{2} \text{Tr} [E \{ (\hat{\mathbf{x}} - \mathbf{x})(\hat{\mathbf{x}} - \mathbf{x})^T \}] + \text{Tr} [\Lambda(I - MH - N)]$$

- Substituting the models into the estimate equation gives

$$\begin{aligned} \hat{\mathbf{x}} &= M\tilde{\mathbf{y}} + N\hat{\mathbf{x}}_a \\ &= (MH + N)\mathbf{x} + M\mathbf{v} + N\mathbf{w} \\ &= \mathbf{x} + M\mathbf{v} + N\mathbf{w} \end{aligned}$$

where the equality constraint $MH + N = I$ was used

- Then we have

$$\begin{aligned} E \{ (\hat{\mathbf{x}} - \mathbf{x})(\hat{\mathbf{x}} - \mathbf{x})^T \} &= ME \{ \mathbf{v} \mathbf{v}^T \} M^T + NE \{ \mathbf{w} \mathbf{w}^T \} N^T \\ &\quad + ME \{ \mathbf{v} \mathbf{w}^T \} N^T + NE \{ \mathbf{w} \mathbf{v}^T \} M^T \end{aligned}$$

- Use $E\{\mathbf{v} \mathbf{v}^T\} = R$, $E\{\mathbf{w} \mathbf{w}^T\} = Q$ and $E\{\mathbf{w} \mathbf{v}^T\} = E\{\mathbf{v} \mathbf{w}^T\} = 0$

$$E \{ (\hat{\mathbf{x}} - \mathbf{x})(\hat{\mathbf{x}} - \mathbf{x})^T \} = MRM^T + NQN^T$$

- So the loss function becomes

$$J = \frac{1}{2} \text{Tr}(MRM^T + NQN^T) + \text{Tr}[\Lambda(I - MH - N)]$$

- The necessary conditions are

$$\nabla_M J = MR - \Lambda^T H^T = 0 \quad (1)$$

$$\nabla_N J = NQ - \Lambda^T = 0 \quad (2)$$

$$\nabla_\Lambda J = I - MH - N = 0 \quad (3)$$

- Solving Eq. (1) for M gives

$$M = \Lambda^T H^T R^{-1}$$

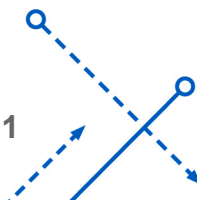
- Solving Eq. (2) for N gives

$$N = \Lambda^T Q^{-1}$$

- Substituting these into Eq. (3) and solving for Λ^T gives

$$\Lambda^T = (H^T R^{-1} H + Q^{-1})^{-1}$$

- This is the covariance for the *a priori* estimates



- Substituting Λ^T into Eqs. (1) and (2) gives

$$M = (H^T R^{-1} H + Q^{-1})^{-1} H^T R^{-1}$$

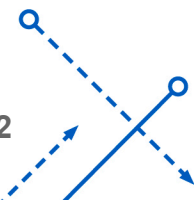
$$N = (H^T R^{-1} H + Q^{-1})^{-1} Q^{-1}$$

- Therefore the *a priori* estimate equation is

$$\hat{\mathbf{x}} = (H^T R^{-1} H + Q^{-1})^{-1} (H^T R^{-1} \tilde{\mathbf{y}} + Q^{-1} \hat{\mathbf{x}}_a)$$

- Some observations

- With poor *a priori* knowledge we have $Q \rightarrow \infty$ and $Q^{-1} \rightarrow 0$, which reduces down to the minimum variance estimator! $\hat{\mathbf{x}} = (H^T R^{-1} H)^{-1} (H^T R^{-1} \tilde{\mathbf{y}})$
- With poor measurements we have $R \rightarrow \infty$ and $R^{-1} \rightarrow 0$, which gives the result $\hat{\mathbf{x}} = \hat{\mathbf{x}}_a$, an intuitively pleasing result!



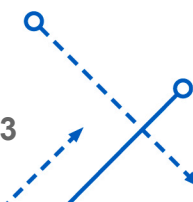
- One of the most useful and important concepts in estimation theory
 - The Cramér-Rao inequality can be used to give us a lower bound on the expected errors between the estimated quantities and the true values from the known statistical properties of the measurement errors
 - Consider the conditional density $p(\tilde{\mathbf{y}}|\mathbf{x})$
 - The Cramér-Rao inequality is given by

$$P \equiv E \left\{ (\hat{\mathbf{x}} - \mathbf{x}) (\hat{\mathbf{x}} - \mathbf{x})^T \right\} \geq F^{-1}$$

where the *Fisher information matrix*, F , is given by

$$F = E \left\{ \left[\frac{\partial}{\partial \mathbf{x}} \ln p(\tilde{\mathbf{y}}|\mathbf{x}) \right] \left[\frac{\partial}{\partial \mathbf{x}} \ln p(\tilde{\mathbf{y}}|\mathbf{x}) \right]^T \right\} = -E \left\{ \frac{\partial^2}{\partial \mathbf{x} \partial \mathbf{x}^T} \ln p(\tilde{\mathbf{y}}|\mathbf{x}) \right\}$$

- Note, Cramér-Rao inequality is only valid for **unbiased estimates** $E(\hat{\mathbf{x}}) = \mathbf{x}$



- Proof begins by using

$$\int_{-\infty}^{\infty} p(\tilde{\mathbf{y}}|\mathbf{x}) d\tilde{\mathbf{y}} = 1$$

- Taking the partial with respect to \mathbf{x} gives

$$\frac{\partial}{\partial \mathbf{x}} \int_{-\infty}^{\infty} p(\tilde{\mathbf{y}}|\mathbf{x}) d\tilde{\mathbf{y}} = \int_{-\infty}^{\infty} \left[\frac{\partial p(\tilde{\mathbf{y}}|\mathbf{x})}{\partial \mathbf{x}} \right] d\tilde{\mathbf{y}} = \mathbf{0}$$

- Since the estimate is assumed to be unbiased we have

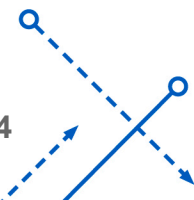
$$E \{ \hat{\mathbf{x}} - \mathbf{x} \} = \int_{-\infty}^{\infty} (\hat{\mathbf{x}} - \mathbf{x}) p(\tilde{\mathbf{y}}|\mathbf{x}) d\tilde{\mathbf{y}} = \mathbf{0} \quad \leftarrow \text{unbiased}$$

Verification

- Differentiating both sides with respect to \mathbf{x} gives

$$\int_{-\infty}^{\infty} (\hat{\mathbf{x}} - \mathbf{x}) \left[\frac{\partial p(\tilde{\mathbf{y}}|\mathbf{x})}{\partial \mathbf{x}} \right]^T d\tilde{\mathbf{y}} - I \int_{-\infty}^{\infty} p(\tilde{\mathbf{y}}|\mathbf{x}) d\tilde{\mathbf{y}} = \mathbf{0}$$

$$\int_{-\infty}^{\infty} (\hat{\mathbf{x}} - \mathbf{x}) \left[\frac{\partial p(\tilde{\mathbf{y}}|\mathbf{x})}{\partial \mathbf{x}} \right]^T d\tilde{\mathbf{y}} - I = \mathbf{0}$$



- Next, we use the following logarithmic differentiation rule

$$\frac{\partial p(\tilde{\mathbf{y}}|\mathbf{x})}{\partial \mathbf{x}} = \left[\frac{\partial}{\partial \mathbf{x}} \ln[p(\tilde{\mathbf{y}}|\mathbf{x})] \right] p(\tilde{\mathbf{y}}|\mathbf{x})$$

- Substitute this into the previous equation to give

$$I = \int_{-\infty}^{\infty} (\mathbf{a} \mathbf{b}^T) d\tilde{\mathbf{y}} \quad (1)$$

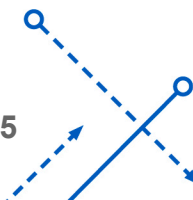
where

$$\mathbf{a} \equiv p(\tilde{\mathbf{y}}|\mathbf{x})^{1/2} (\hat{\mathbf{x}} - \mathbf{x})$$

$$\mathbf{b} \equiv p(\tilde{\mathbf{y}}|\mathbf{x})^{1/2} \left[\frac{\partial}{\partial \mathbf{x}} \ln[p(\tilde{\mathbf{y}}|\mathbf{x})] \right]$$

- Note that P and F can be written now as

$$P = \int_{-\infty}^{\infty} (\mathbf{a} \mathbf{a}^T) d\tilde{\mathbf{y}}, \quad F = \int_{-\infty}^{\infty} (\mathbf{b} \mathbf{b}^T) d\tilde{\mathbf{y}}$$



- Multiply Eq. (1) on the left by an arbitrary row vector α^T and on the right by an arbitrary column vector β

$$\alpha^T \beta = \int_{-\infty}^{\infty} \alpha^T (\mathbf{a} \mathbf{b}^T) \beta d\tilde{\mathbf{y}}$$

- Next, we make use of the *Schwartz inequality*

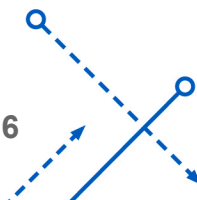
$$\left[\int_{-\infty}^{\infty} g(\tilde{\mathbf{y}}|\mathbf{x}) h(\tilde{\mathbf{y}}|\mathbf{x}) d\tilde{\mathbf{y}} \right]^2 \leq \int_{-\infty}^{\infty} g^2(\tilde{\mathbf{y}}|\mathbf{x}) d\tilde{\mathbf{y}} \int_{-\infty}^{\infty} h^2(\tilde{\mathbf{y}}|\mathbf{x}) d\tilde{\mathbf{y}}$$

If $\int_{-\infty}^{\infty} a(\mathbf{x})b(\mathbf{x}) d\mathbf{x} = 1$ then $\int_{-\infty}^{\infty} a^2(\mathbf{x}) d\mathbf{x} \int_{-\infty}^{\infty} b^2(\mathbf{x}) d\mathbf{x} \geq 1$; the equality holds if $a(\mathbf{x}) = c b(\mathbf{x})$ where c is not a function of \mathbf{x} .

- Define the following quantities

$$g(\tilde{\mathbf{y}}|\mathbf{x}) = \alpha^T \mathbf{a}$$

$$h(\tilde{\mathbf{y}}|\mathbf{x}) = \mathbf{b}^T \beta$$



- Then the Schwartz inequality becomes

$$\left[\int_{-\infty}^{\infty} \alpha^T (\mathbf{a} \mathbf{b}^T) \beta d\tilde{\mathbf{y}} \right]^2 \leq \int_{-\infty}^{\infty} \alpha^T (\mathbf{a} \mathbf{a}^T) \alpha d\tilde{\mathbf{y}} \int_{-\infty}^{\infty} \beta^T (\mathbf{b} \mathbf{b}^T) \beta d\tilde{\mathbf{y}}$$

- Using the definitions of P and F and

$$\int_{-\infty}^{\infty} \mathbf{a} \mathbf{b}^T d\tilde{\mathbf{y}} = \int_{-\infty}^{\infty} (\hat{\mathbf{x}} - \mathbf{x}) \left[\frac{\partial p(\tilde{\mathbf{y}}|\mathbf{x})}{\partial \mathbf{x}} \right]^T d\tilde{\mathbf{y}} = I$$

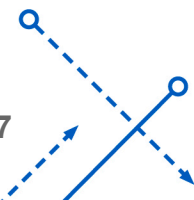
gives

$$(\alpha^T \beta)^2 \leq (\alpha^T P \alpha) (\beta^T F \beta)$$

- Finally, using the particular choice $\beta = F^{-1} \alpha$ gives

$$\alpha^T (P - F^{-1}) \alpha \geq 0$$

- Since α is arbitrary then $P \geq F^{-1}$ must be true, which proves the Cramér-Rao Inequality



- Consider the measurement model

$$\tilde{\mathbf{y}} = H\mathbf{x} + \mathbf{v}, \quad \text{with} \quad E\{\mathbf{v}\} = \mathbf{0} \quad \text{and} \quad E\{\mathbf{v}\mathbf{v}^T\} = R$$

- To determine the mean of the observation model, we take the expectation of both sides

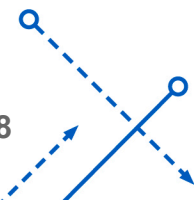
$$\boldsymbol{\mu} \equiv E\{\tilde{\mathbf{y}}\} = E\{H\mathbf{x}\} + E\{\mathbf{v}\} = H\mathbf{x}$$

- The covariance is then given by

$$\begin{aligned} \text{cov}\{\tilde{\mathbf{y}}\} &\equiv E\left\{(\tilde{\mathbf{y}} - \boldsymbol{\mu})(\tilde{\mathbf{y}} - \boldsymbol{\mu})^T\right\} \\ &= E\{\mathbf{v}\mathbf{v}^T\} = R \end{aligned}$$

- The conditional density is then given by

$$p(\tilde{\mathbf{y}}|\mathbf{x}) = \frac{1}{(2\pi)^{m/2} [\det(R)]^{1/2}} \exp\left\{-\frac{1}{2} [\tilde{\mathbf{y}} - H\mathbf{x}]^T R^{-1} [\tilde{\mathbf{y}} - H\mathbf{x}]\right\}$$



- Taking the natural log gives

$$\ln [p(\tilde{\mathbf{y}}|\mathbf{x})] = -\frac{1}{2} [\tilde{\mathbf{y}} - H\mathbf{x}]^T R^{-1} [\tilde{\mathbf{y}} - H\mathbf{x}] - \frac{m}{2} \ln (2\pi) - \frac{1}{2} \ln [\det (R)]$$

- Carry out the computations for the Fisher Information Matrix

$$F = -E \left\{ \frac{\partial^2}{\partial \mathbf{x} \partial \mathbf{x}^T} \ln p(\tilde{\mathbf{y}}|\mathbf{x}) \right\} = (H^T R^{-1} H)$$

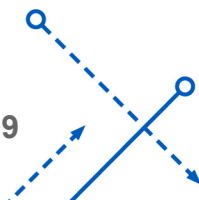
- Hence, the Cramér-Rao inequality is $P \geq (H^T R^{-1} H)^{-1}$
- Let us now find an expression for the estimate covariance P
- Estimate and measurement models

$$\hat{\mathbf{x}} = (H^T R^{-1} H)^{-1} H^T R^{-1} \tilde{\mathbf{y}}$$

$$\tilde{\mathbf{y}} = H\mathbf{x} + \mathbf{v}$$

- Substituting the measurement model into the estimate gives

$$\begin{aligned} \hat{\mathbf{x}} &= (H^T R^{-1} H)^{-1} H^T R^{-1} H\mathbf{x} + (H^T R^{-1} H)^{-1} H^T R^{-1} \mathbf{v} \\ &= \mathbf{x} + (H^T R^{-1} H)^{-1} H^T R^{-1} \mathbf{v} \end{aligned}$$



- The expectation of the estimate is given by

$$E\{\hat{\mathbf{x}}\} = \mathbf{x} + (H^T R^{-1} H)^{-1} H^T R^{-1} E\{\mathbf{v}\} = \mathbf{x}$$

since $E\{\mathbf{v}\} = \mathbf{0}$

- The covariance is

$$\begin{aligned} P &\equiv E\{(\hat{\mathbf{x}} - E\{\hat{\mathbf{x}}\})(\hat{\mathbf{x}} - E\{\hat{\mathbf{x}}\})^T\} \\ &= E\{(\hat{\mathbf{x}} - \mathbf{x})(\hat{\mathbf{x}} - \mathbf{x})^T\} \\ &= (H^T R^{-1} H)^{-1} H^T R^{-1} E\{\mathbf{v} \mathbf{v}^T\} R^{-1} H (H^T R^{-1} H)^{-1} \end{aligned}$$

- From $E\{\mathbf{v} \mathbf{v}^T\} = R$ we have

$$\begin{aligned} P &= (H^T R^{-1} H)^{-1} H^T R^{-1} R R^{-1} H (H^T R^{-1} H)^{-1} \\ &= (H^T R^{-1} H)^{-1} \end{aligned}$$

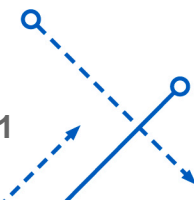
- Therefore, the equality is satisfied, so the least squares estimate from the Gauss-Markov Theorem is the most efficient possible estimate!
- Estimator is thus called **efficient**

- Important result

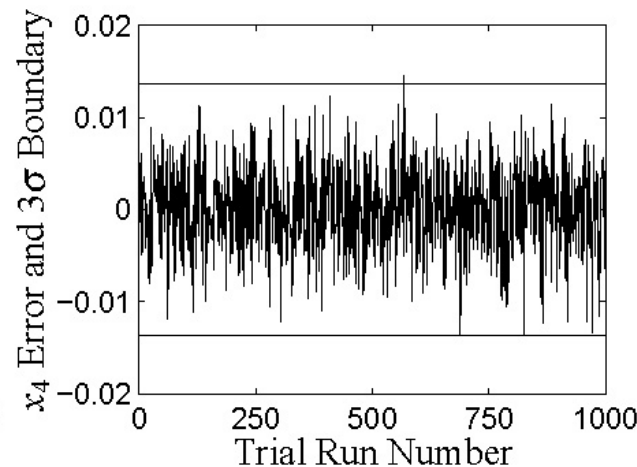
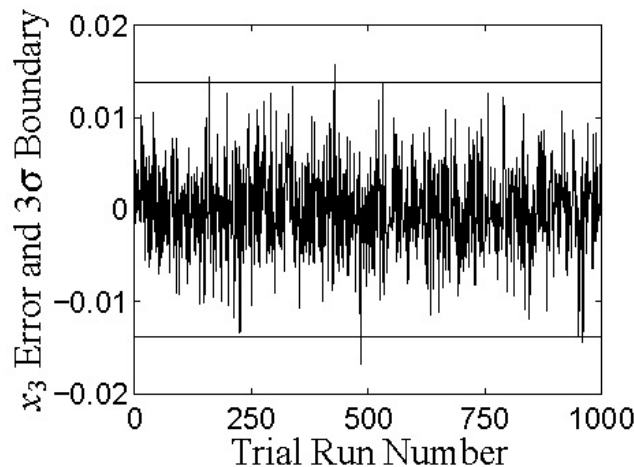
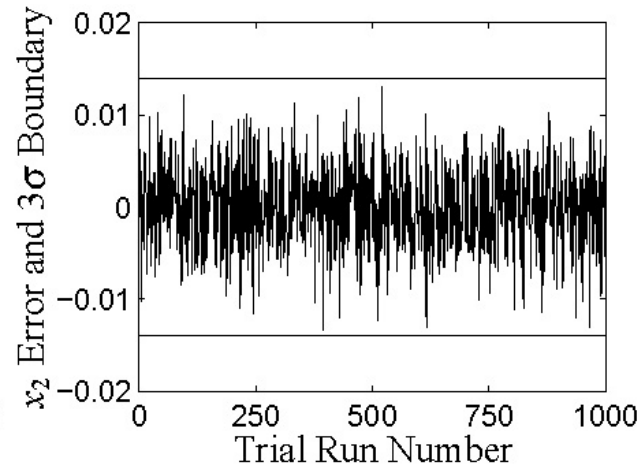
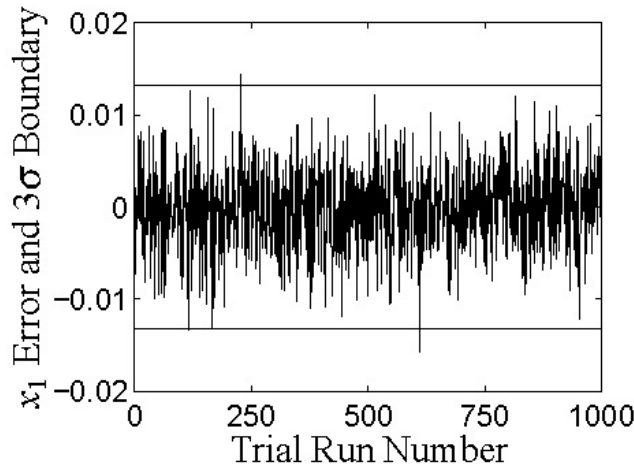
- Note again that the covariance of the estimation errors is given by

$$P = (H^T R^{-1} H)^{-1}$$

- We never know the truth in the real world
- But if we know the characteristics of the measurement errors (zero-mean with known covariance R) then we can determine a bound on the estimation errors from a statistical point of view
- This is certainly useful information!
- Note that the covariance of the estimation errors can be computed without ever computing the estimate
- Helps to assess the performance of the estimator
 - For example, useful to develop an error budget for the total attitude errors in a spacecraft attitude control design



$$\tilde{y}(t) = \cos(t) + 2 \sin(t) + \cos(2t) + 2 \sin(3t) + v(t), \quad R = 0.01I$$



Ran 1,000 Monte Carlo runs

3σ boundaries found from taking the square root of the diagonal elements of P and multiplying the result by 3

Bounds actual errors well




```
% True System
dt=0.01;tf=10;
t=[0:dt:tf]';
m=length(t);
y=cos(t)+2*sin(t)+cos(2*t)+2*sin(3*t);
```

```
% Pre-allocate Space
xe=zeros(1000,4);
pcov=zeros(1000,4);
```

```
% Monte Carlo Simulation
for i=1:1000,
    ym=y+0.1*randn(m,1);w=1/0.01;
    h=[cos(t) sin(t) cos(2*t) sin(3*t)];
    p=inv(h'*w*h);
    xe(i,:)=(p*h'*w*ym)';
    pcov(i,:)=diag(p)';
end
```

```
% Plot Results
subplot(221)
plot([1:1000],xe(:,1)-1,[1:1000],pcov(:,1).^(0.5)*3,[1:1000],-pcov(:,1).^(0.5)*3);
axis([0 1000 -0.02 0.02]);
set(gca,'fontsize',12);
set(gca,'xtick',[0 250 500 750 1000]);
set(gca,'ytick',[-0.02 -0.01 0 0.01 0.02]);
xlabel('Trial Run Number')
ylabel('x_1 Error and 3 sigma Outlier')
```

```
subplot(222)
plot([1:1000],xe(:,2)-2,[1:1000],pcov(:,2).^(0.5)*3,[1:1000],-pcov(:,2).^(0.5)*3);
axis([0 1000 -0.02 0.02]);
set(gca,'fontsize',12);
set(gca,'xtick',[0 250 500 750 1000]);
set(gca,'ytick',[-0.02 -0.01 0 0.01 0.02]);
xlabel('Trial Run Number')
ylabel('x_2 Error and 3 sigma Outlier')
```

```
subplot(223)
plot([1:1000],xe(:,3)-1,[1:1000],pcov(:,3).^(0.5)*3,[1:1000],-pcov(:,3).^(0.5)*3);
axis([0 1000 -0.02 0.02]);
set(gca,'fontsize',12);
set(gca,'xtick',[0 250 500 750 1000]);
set(gca,'ytick',[-0.02 -0.01 0 0.01 0.02]);
xlabel('Trial Run Number')
ylabel('x_3 Error and 3 \sigma Outlier')
```

```
subplot(224)
plot([1:1000],xe(:,4)-2,[1:1000],pcov(:,4).^(0.5)*3,[1:1000],-pcov(:,4).^(0.5)*3);
axis([0 1000 -0.02 0.02]);
set(gca,'fontsize',12);
set(gca,'xtick',[0 250 500 750 1000]);
set(gca,'ytick',[-0.02 -0.01 0 0.01 0.02]);
xlabel('Trial Run Number')
ylabel('x_4 Error and 3 \sigma Outlier')
```

- Suppose we wish to estimate a nonlinear appearing parameter, $a > 0$, of the following exponential model

$$\tilde{y}_k = B e^{a t_k} + v_k, \quad k = 1, 2, \dots, m$$

where v_k is a zero-mean Gaussian white-noise process with variance given by σ^2

- We can choose to employ nonlinear least squares to iteratively determine the parameter a , given the measurements and a known $B > 0$ coefficient
- The covariance of the estimate error is given by $P = \sigma^2 (H^T H)^{-1}$ (this is the Cramér-Rao bound too) with

$$H = \begin{bmatrix} B t_1 e^{a t_1} & B t_2 e^{a t_2} & \dots & B t_m e^{a t_m} \end{bmatrix}^T$$

- Note that H is a function of the true parameter a now
- This can be replaced by the final estimate after the nonlinear least squares iteration is complete (errors are second-order in nature)



- Let's instead employ linear least squares by using a change of variables, as shown before, with $\tilde{z}_k \equiv \ln \tilde{y}_k$
 - Question: How optimal is this approach?
 - Expanding \tilde{z}_k in a first-order series gives

$$\ln \tilde{y}_k - \ln B \approx a t_k + \frac{2 v_k}{2 B e^{a t_k} + v_k}$$

- The least squares “ H matrix” is now simply given by

$$\mathcal{H} = [t_1 \quad t_2 \quad \cdots \quad t_m]^T$$

- A first-order expansion using the binomial series of the new measurement noise is given by

$$\varepsilon_k \equiv 2 v_k (2 B e^{a t_k} + v_k)^{-1} \approx \frac{v_k}{B e^{a t_k}} \left(1 - \frac{v_k}{2 B e^{a t_k}} \right)$$

- The variance can be shown to be given by

$$\varsigma_k^2 = E\{\varepsilon_k^2\} - E\{\varepsilon_k\}^2 = E \left\{ \left(\frac{v_k}{B e^{a t_k}} - \frac{v_k^2}{2 B^2 e^{2 a t_k}} \right)^2 \right\} - \frac{\sigma^4}{4 B^2 e^{4 a t_k}}$$

- This leads to

$$\varsigma_k^2 = \frac{\sigma^2}{B^2 e^{2 a t_k}} + \frac{\sigma^4}{2 B^4 e^{4 a t_k}}$$

- Contains both Gaussian and χ^2 components
- The covariance of the linear approach is given by

$$\mathcal{P} = (\mathcal{H}^T \text{diag} [\varsigma_1^{-2} \quad \varsigma_2^{-2} \quad \cdots \quad \varsigma_m^{-2}] \mathcal{H})^{-1}$$

- Both covariances are equivalent if $\sigma^4 / (2 B^4 e^{4 a t_k})$ is negligible
- If this is not the case, then the Cramér-Rao lower bound is not achieved and the linear approach does not lead to an efficient estimator
- This clearly shows how the Cramér-Rao inequality can be particularly useful to help quantify the errors introduced by using an approximate solution instead of the optimal approach

