

# 7.6 MLE for Transformed Parameters

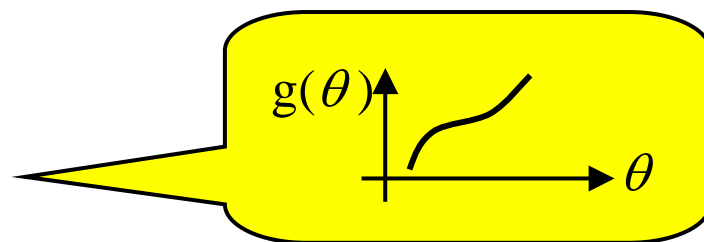
Given PDF  $p(\mathbf{x}; \theta)$  but want an estimate of  $\alpha = g(\theta)$

What is the MLE for  $\alpha$ ??

Two cases:

1.  $\alpha = g(\theta)$  is a one-to-one function

$$\hat{\alpha}_{ML} \text{ maximizes } p(\mathbf{x}; g^{-1}(\alpha))$$



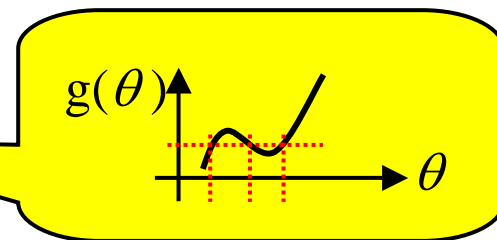
2.  $\alpha = g(\theta)$  is not a one-to-one function

Need to define modified likelihood function:

$$\bar{p}_T(\mathbf{x}; \alpha) = \max_{\{\theta: \alpha = g(\theta)\}} p(\mathbf{x}; \theta)$$

$$\hat{\alpha}_{ML} \text{ maximizes } \bar{p}_T(\mathbf{x}; \alpha)$$

- For each  $\alpha$ , find all  $\theta$ 's that map to it
- Extract largest value of  $p(\mathbf{x}; \theta)$  over this set of  $\theta$ 's



# Invariance Property of MLE

Another Big  
Advantage of MLE!

## Theorem 7.2: Invariance Property of MLE

If parameter  $\theta$  is mapped according to  $\alpha = g(\theta)$  then the MLE of  $\alpha$  is given by

$$\hat{\alpha} = g(\hat{\theta})$$

where  $\hat{\theta}$  is the MLE for  $\theta$  found by maximizing  $p(\mathbf{x}; \theta)$

Note: when  $g(\theta)$  is not one-to-one the MLE for  $\alpha$  maximizes the modified likelihood function

“Proof”:

Easy to see when  $g(\theta)$  is one-to-one

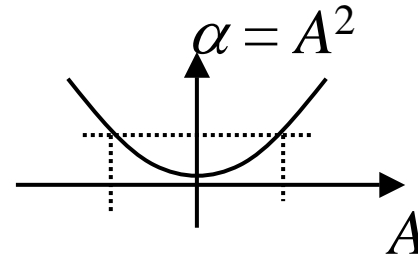
Otherwise... can “argue” that maximization over  $\theta$  inside definition for modified LF ensures the result.

# Ex. 7.9: Estimate Power of DC Level in AWGN

$$x[n] = A + w[n]$$

noise is  $N(0, \sigma^2)$  & White

Want to Est. Power:  $\alpha = A^2 \Rightarrow$



$\Rightarrow$  For each  $\alpha$  value there are 2 PDF's to consider

$$p_{T_1}(\mathbf{x}; \alpha) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp \left[ -\frac{1}{2\sigma^2} \sum_n (x[n] - \sqrt{\alpha})^2 \right]$$

$$p_{T_2}(\mathbf{x}; \alpha) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp \left[ -\frac{1}{2\sigma^2} \sum_n (x[n] + \sqrt{\alpha})^2 \right]$$

Then:

$$\begin{aligned} \hat{\alpha}_{ML} &= \left[ \arg \max_{\sqrt{\alpha} \geq 0} \left\{ p(\mathbf{x}; \sqrt{\alpha}), p(\mathbf{x}; -\sqrt{\alpha}) \right\} \right]^2 \\ &= \left[ \arg \max_{-\infty < A < \infty} p(\mathbf{x}; A) \right]^2 \\ &= \left[ \hat{A}_{ML} \right]^2 \end{aligned}$$

Demonstration that  
Invariance Result  
Holds for this  
Example

## Ex. 7.10: Estimate Power of WGN in dB

$$x[n] = w[n] \quad \text{WGN} \quad \text{w/ var} = \sigma^2 \text{ unknown}$$

$$\text{Recall: } P_{\text{noise}} = \sigma^2$$

Can show that the MLE for variance is:  $\hat{P}_{\text{noise}} = \frac{1}{N} \sum_{n=0}^{N-1} x^2[n]$

To get the dB version of the power estimate:

$$\hat{P}_{dB} = 10 \log_{10} \left[ \frac{1}{N} \sum_{n=0}^{N-1} x^2[n] \right]$$

Using  
Invariance Property !

Note: You may recall a result for estimating variance that divides by  $N-1$  rather than by  $N$  ... that estimator is unbiased, this estimate is biased (but asymptotically unbiased)

## 7.7: Numerical Determination of MLE

**Note:** In all previous examples we ended up with a closed-form expression for the MLE:  $\hat{\theta}_{ML} = f(\mathbf{x})$

**Ex. 7.11:**  $x[n] = r^n + w[n]$

**Estimate  $r$**

To find MLE:

$$\frac{\partial \ln p(\mathbf{x}; \theta)}{\partial \theta} = 0$$

$$\Rightarrow \sum_{n=0}^{N-1} (x[n] - r^n) n r^{n-1} = 0$$

noise is  $N(0, \sigma^2)$  & white

If  $-1 < r < 0$  then this signal is a decaying oscillation that might be used to model:

- A Ship's "Hull Ping"
- A Vibrating String, Etc.

No closed-form solution for the MLE

So...we can't always find a closed-form MLE!

**But a main advantage of MLE is:**

**We can always find it numerically!!!**

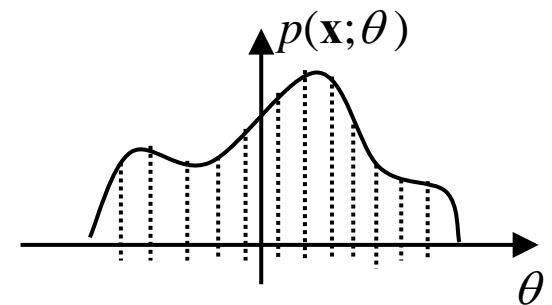
**(Not always computationally efficiently, though)**

### **Brute Force Method**

Compute  $p(\mathbf{x};\theta)$  on a fine grid of  $\theta$  values

Advantage: Sure to Find maximum  
(if grid is fine enough)

Disadvantage: Lots of Computation  
(especially w/ a fine grid)



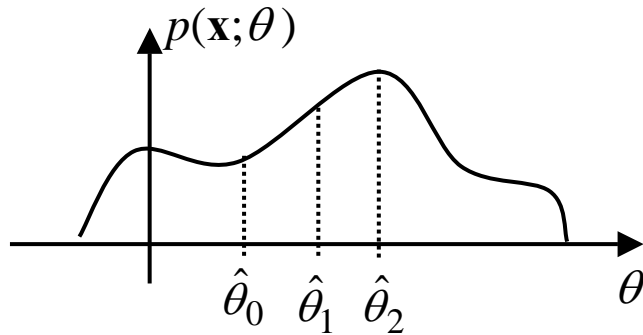
# Iterative Methods for Numerical MLE

Step #1: Pick some “initial estimate”  $\hat{\theta}_0$

Step #2: Iteratively improve it using

$$\hat{\theta}_{i+1} = f(\hat{\theta}_i, \mathbf{x}) \quad \text{such that} \quad \lim_{i \rightarrow \infty} p(\mathbf{x}; \theta_i) = \max_{\theta} p(\mathbf{x}; \theta)$$

## “Hill Climbing in the Fog”



**Note:** A so-called “Greedy” maximization algorithm will always move up even though taking an occasional step downward may be the better global strategy!

## Convergence Issues:

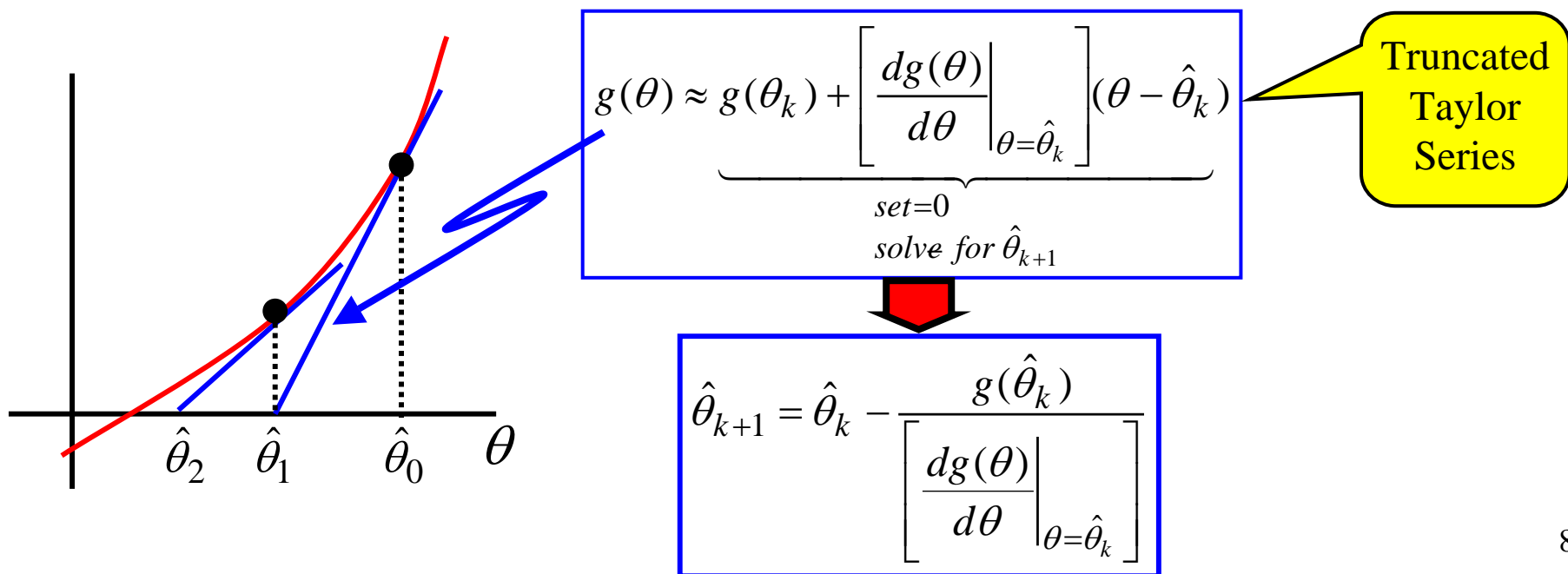
1. May not converge
2. May converge, but to local maximum
  - good initial guess is needed !!
  - can use rough grid search to initialize
  - can use multiple initializations

# Iterative Method: Newton-Raphson MLE

The MLE is the maximum of the LF... so set derivative to 0:

$$\underbrace{\frac{\partial \ln p(\mathbf{x}; \theta)}{\partial \theta}}_{\triangleq g(\theta)} = 0 \quad \text{So... MLE is a zero of } g(\theta)$$

Newton-Raphson is a numerical method for finding the zero of a function... so it can be applied here... **Linearize  $g(\theta)$**





Now... using our “definition of convenience”:  $g(\theta) = \frac{\partial \ln p(\mathbf{x}; \theta)}{\partial \theta}$

So then the Newton-Raphson MLE iteration is:

$$\hat{\theta}_{k+1} = \hat{\theta}_k - \left\{ \left[ \frac{\partial^2 \ln p(\mathbf{x}; \theta)}{\partial \theta^2} \right]^{-1} \frac{\partial \ln p(\mathbf{x}; \theta)}{\partial \theta} \right\} \Big|_{\theta = \hat{\theta}_k}$$

Iterate until  
convergence  
criterion is met:

$$|\hat{\theta}_{k+1} - \hat{\theta}_k| < \varepsilon$$

Look Familiar???

Looks like  $I(\theta)$ , except:  $I(\theta)$  is evaluated at the true  $\theta$ , and has an expected value

You get to choose!

## Generally:

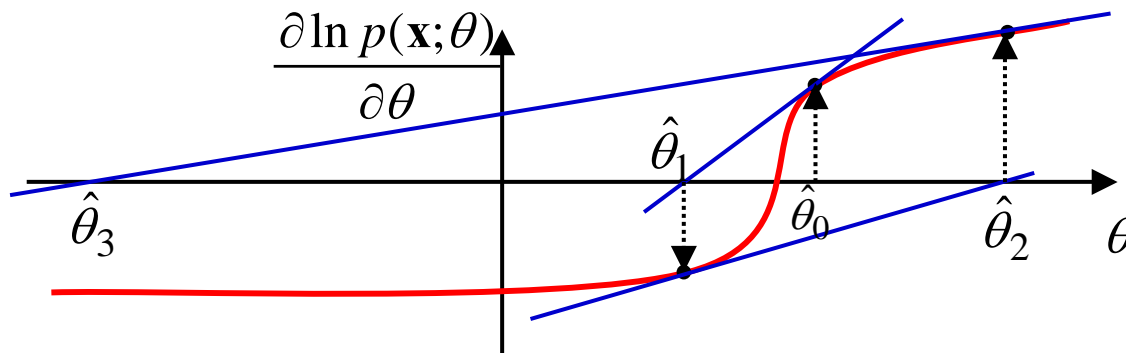
For a given PDF model, compute derivatives analytically...

or... compute derivatives numerically:

$$\frac{\partial \ln p(\mathbf{x}; \theta)}{\partial \theta} \Big|_{\hat{\theta}_k} \approx \frac{\ln p(\mathbf{x}; \hat{\theta}_k + \Delta\theta) - \ln p(\mathbf{x}; \hat{\theta}_k)}{\Delta\theta}$$

## Convergence Issues of Newton-Raphson:

1. May not converge
2. May converge, but to local maximum
  - good initial guess is needed !!
  - can use rough grid search to initialize
  - can use multiple initializations



## Some Other Iterative MLE Methods

1. Scoring Method
  - Replaces second-partial term by  $I(\theta)$
2. Expectation-Maximization (EM) Method
  - Guarantees convergence to at least a local maximum
  - Good for complicated multi-parameter cases

## 7.8 MLE for Vector Parameter

Another nice property of MLE is how easily it carries over to the vector parameter case.

The vector parameter is:  $\boldsymbol{\theta} = [\theta_1 \quad \theta_2 \quad \dots \quad \theta_p]^T$

$\hat{\boldsymbol{\theta}}_{ML}$  is the vector that satisfies:  $\underbrace{\frac{\partial \ln p(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}}_{\text{Derivative w.r.t. a vector}} = 0$

$$\frac{\partial f(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \begin{bmatrix} \frac{\partial f(\boldsymbol{\theta})}{\partial \theta_1} \\ \frac{\partial f(\boldsymbol{\theta})}{\partial \theta_2} \\ \vdots \\ \frac{\partial f(\boldsymbol{\theta})}{\partial \theta_p} \end{bmatrix}$$

Derivative w.r.t.  
a vector

## Ex. 7.12: Estimate DC Level and Variance

$x[n] = A + w[n]$  noise is  $N(0, \sigma^2)$  and white

Estimate: DC level  $A$  and Noise Variance  $\sigma^2 \Rightarrow \boldsymbol{\theta} = \begin{bmatrix} A \\ \sigma^2 \end{bmatrix}$

LF is: 
$$p(\mathbf{x}; A, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{\frac{N}{2}}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{n=0}^{N-1} [x[n] - A]^2\right\}$$

Solve: 
$$\frac{\partial \ln p(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \stackrel{\text{set}}{=} \mathbf{0}$$

$$\frac{\partial \ln p(\mathbf{x}; \boldsymbol{\theta})}{\partial A} = \frac{1}{\sigma^2} \sum_{n=0}^{N-1} (x[n] - A) = \frac{N}{\sigma^2} (\bar{x} - A) = 0$$

$$\frac{\partial \ln p(\mathbf{x}; \boldsymbol{\theta})}{\partial \sigma^2} = -\frac{N}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{n=0}^{N-1} (x[n] - A)^2 = 0$$

$$\hat{\boldsymbol{\theta}}_{ML} = \begin{bmatrix} \bar{x} \\ \frac{1}{N} \sum_n (x[n] - \bar{x})^2 \end{bmatrix}$$

Interesting: For this problem...

First estimate  $A$  just like scalar case

The subtract it off and then estimate variance like scalar case

# Properties of Vector ML

The asymptotic properties are captured in Theorem 7.3:

If  $p(\mathbf{x};\boldsymbol{\theta})$  satisfies some “regularity” conditions, then the MLE is asymptotically distributed according to

$$\hat{\boldsymbol{\theta}}_{ML} \overset{a}{\sim} N(\boldsymbol{\theta}, \mathbf{I}^{-1}(\boldsymbol{\theta}))$$

where  $\mathbf{I}(\boldsymbol{\theta}) =$  Fisher Information Matrix

So the vector ML is asymptotically:

- unbiased
- efficient

## Invariance Property Holds for Vector Case

If  $\boldsymbol{\alpha} = g(\boldsymbol{\theta})$ , then  $\hat{\boldsymbol{\alpha}}_{ML} = g(\hat{\boldsymbol{\theta}}_{ML})$

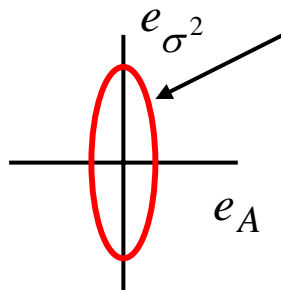
## Ex. 7.12 Revisited

It can be shown that:  $E\{\hat{\boldsymbol{\theta}}\} = \begin{bmatrix} A \\ \frac{(N-1)}{N}\sigma^2 \end{bmatrix}$        $\text{cov}\{\hat{\boldsymbol{\theta}}\} = \begin{bmatrix} \frac{\sigma^2}{N} & 0 \\ 0 & \frac{2(N-1)}{N^2}\sigma^4 \end{bmatrix}$

For large  $N$  then :  $E\{\hat{\boldsymbol{\theta}}\} \approx \begin{bmatrix} A \\ \sigma^2 \end{bmatrix} = \boldsymbol{\theta}$        $\text{cov}\{\hat{\boldsymbol{\theta}}\} \approx \begin{bmatrix} \frac{\sigma^2}{N} & 0 \\ 0 & \frac{2}{N}\sigma^4 \end{bmatrix} = \mathbf{I}^{-1}(\boldsymbol{\theta})$

which we see satisfies the asymptotic property.

Diagonal covariance matrix shows estimates are uncorrelated:



**Error Ellipse is aligned with axes**

This is why we could “decouple” the estimates

# MLE for the General Gaussian Case

Let the data be general Gaussian:  $\mathbf{x} \sim N(\boldsymbol{\mu}(\boldsymbol{\theta}), \mathbf{C}(\boldsymbol{\theta}))$

Thus  $\partial \ln p(\mathbf{x}; \boldsymbol{\theta}) / \partial \boldsymbol{\theta}$  will depend in general on  $\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$  and  $\frac{\partial \mathbf{C}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$

For each  $k = 1, 2, \dots, p$  set:  $\frac{\partial \ln p(\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_k} = 0$

This gives  $p$  simultaneous equations, the  $k^{\text{th}}$  one being:

$$\underbrace{-\frac{1}{2} \text{tr} \left( \mathbf{C}^{-1}(\boldsymbol{\theta}) \frac{\partial \mathbf{C}(\boldsymbol{\theta})}{\partial \theta_k} \right)}_{\text{Term \#1}} + \underbrace{\left[ \frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \theta_k} \right]^T \mathbf{C}^{-1}(\boldsymbol{\theta}) [\mathbf{x} - \boldsymbol{\mu}(\boldsymbol{\theta})]}_{\text{Term \#2}} - \underbrace{\frac{1}{2} [\mathbf{x} - \boldsymbol{\mu}(\boldsymbol{\theta})]^T \left[ \frac{\partial \mathbf{C}^{-1}(\boldsymbol{\theta})}{\partial \theta_k} \right] [\mathbf{x} - \boldsymbol{\mu}(\boldsymbol{\theta})]}_{\text{Term \#3}} = 0$$

**Note:** for the “deterministic signal + noise” case: Terms #1 & #3 are zero

This gives general conditions to find the MLE...

but can't always solve it!!!

# MLE for Linear Model Case

For this case we can solve these equations!

The signal model is:  $\mathbf{x} = \mathbf{H}\boldsymbol{\theta} + \mathbf{w}$  with the noise  $\mathbf{w} \sim N(\mathbf{0}, \mathbf{C})$

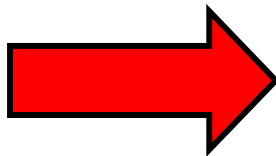
So terms #1 & #3 are zero and term #2 gives: 
$$\underbrace{\left[ \frac{\partial(\mathbf{H}\boldsymbol{\theta})^T}{\partial\boldsymbol{\theta}} \right]}_{=\mathbf{H}} \mathbf{C}^{-1} [\mathbf{x} - \mathbf{H}\boldsymbol{\theta}] = \mathbf{0}$$

Solving this gives: 
$$\hat{\boldsymbol{\theta}}_{ML} = \left( \mathbf{H}^T \mathbf{C}^{-1} \mathbf{H} \right)^{-1} \mathbf{H}^T \mathbf{C}^{-1} \mathbf{x}$$

**Hey!** Same as chapter 4's MVU for linear model

**For Linear Model: ML = MVU**

*Recall: the Linear Model is specified to have Gaussian noise*



$$\hat{\boldsymbol{\theta}}_{ML} \sim N(\boldsymbol{\theta}, (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1})$$

**EXACT...**  
Not Asymptotic!!



# Numerical Solutions for Vector Case

Obvious generalizations... see p. 187

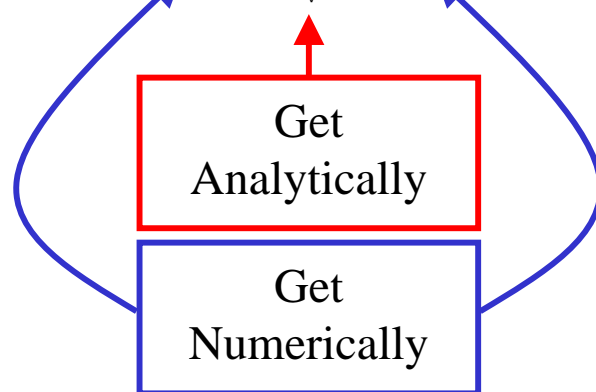
There is one issue to be aware of, though:

The numerical implementation needs  $\partial \ln p(\mathbf{x}; \boldsymbol{\theta}) / \partial \boldsymbol{\theta}$

For the general Gaussian case this requires:  $\frac{\partial \mathbf{C}^{-1}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$

So... we use (3C.2):

$$\frac{\partial \mathbf{C}^{-1}(\boldsymbol{\theta})}{\partial \theta_k} = - \underbrace{\mathbf{C}^{-1}(\boldsymbol{\theta})}_{\text{Get Analytically}} \underbrace{\frac{\partial \mathbf{C}(\boldsymbol{\theta})}{\partial \theta_k}}_{\text{Get Numerically}} \underbrace{\mathbf{C}^{-1}(\boldsymbol{\theta})}_{\text{Get Analytically}}$$



...often hard to analytically:  
get  $\mathbf{C}^{-1}(\boldsymbol{\theta})$   
& then  
differentiate!

## 7.9 Asymptotic MLE

Useful when data samples  $x[n]$  come from a WSS process

**Reading Assignment Only**