

# LINEAR MODELS

**Polynomial Curve Fitting Example.** Continuous signal  $x(t)$  is modeled as a polynomial of degree  $p - 1$  in additive noise:

$$x(t) = \theta_1 + \theta_2 t + \cdots + \theta_p t^{p-1} + w(t).$$

Suppose that we are given  $\{x(t_n)\}_{n=0}^{N-1}$ . Define

$$\begin{aligned} \mathbf{x} &= [x(t_0), \dots, x(t_{N-1})]^T \\ \mathbf{w} &= [w(t_0), \dots, w(t_{N-1})]^T \\ \boldsymbol{\theta} &= [\theta_1, \dots, \theta_p]^T \\ \mathbf{H} &= \begin{bmatrix} 1 & t_0 & \cdots & t_0^{p-1} \\ 1 & t_1 & \cdots & t_1^{p-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & t_{N-1} & \cdots & t_{N-1}^{p-1} \end{bmatrix} \quad (\text{an } N \times p \text{ matrix}). \end{aligned}$$

The data model is then

$$\mathbf{x} = \mathbf{H}\boldsymbol{\theta} + \mathbf{w}$$

where  $\mathbf{H}$  is known and  $\boldsymbol{\theta}$  is the parameter vector to be estimated.

## Sinusoidal Amplitude and Phase Estimation

Measured signal  $x(t)$  is modeled as a superposition of  $p/2$  sinusoids (having known frequencies but unknown amplitudes and phases):

$$x(t) = \sum_{k=1}^{p/2} r_k \sin(\omega_k t + \phi_k) + w(t).$$

This model is linear in  $r_k$  but nonlinear in  $\phi_k$ . However, we can rewrite it as

$$x(t) = \sum_{k=1}^{p/2} [A_k \cos(\omega_k t) + B_k \sin(\omega_k t)] + w(t).$$

Given  $\mathbf{x} = [x(t_0), \dots, x(t_{N-1})]^T$ , we get the following model:

$$\mathbf{x} = \mathbf{H}\boldsymbol{\theta} + \mathbf{w}.$$

## Linear Models (cont.)

For  $p/2 = 2$  sinusoids:

$$\mathbf{H} = \begin{bmatrix} \cos(\omega_1 t_0) & \cos(\omega_2 t_0) & \sin(\omega_1 t_0) & \sin(\omega_2 t_0) \\ \cos(\omega_1 t_1) & \cos(\omega_2 t_1) & \sin(\omega_1 t_1) & \sin(\omega_2 t_1) \\ \vdots & \vdots & \vdots & \vdots \\ \cos(\omega_1 t_{N-1}) & \cos(\omega_2 t_{N-1}) & \sin(\omega_1 t_{N-1}) & \sin(\omega_2 t_{N-1}) \end{bmatrix}$$

and

$$\boldsymbol{\theta} = [A_1, \dots, A_{p/2}, B_1, \dots, B_{p/2}]^T.$$

Once we compute an estimate  $\hat{\boldsymbol{\theta}}$  of  $\boldsymbol{\theta}$ ,  $\hat{r}_k$  and  $\hat{\phi}_k$  are obtained using the simple conversion from rectangular to polar coordinates.

**Note:** Even if  $\hat{\boldsymbol{\theta}}$  is a minimum variance unbiased (MVU) estimator,  $\{\hat{r}_k\}$  and  $\{\hat{\phi}_k\}$  will only be *asymptotically* MVU (for large  $N$ ), as we will see later.

# General Problem Formulation

Consider the model

$$\mathbf{x} = \mathbf{H}\boldsymbol{\theta} + \underbrace{\mathbf{w}}_{\text{noise}}$$

where  $\mathbf{x}$  is a measured  $N \times 1$  vector and  $\mathbf{H}$  is a known *deterministic*  $N \times p$  matrix, with  $N \geq p$ . We wish to estimate the *unknown* parameter vector  $\boldsymbol{\theta}$ .

Assume that  $\mathbf{w}$  is distributed as  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ .

Recall the *identifiability condition*:

$$p(\mathbf{x}; \boldsymbol{\theta}_1) = p(\mathbf{x}; \boldsymbol{\theta}_2) \quad \Leftrightarrow \quad \boldsymbol{\theta}_1 = \boldsymbol{\theta}_2$$

which, in this case, reduces to

$$\mathbf{H}\boldsymbol{\theta}_1 = \mathbf{H}\boldsymbol{\theta}_2 \quad \Leftrightarrow \quad \boldsymbol{\theta}_1 = \boldsymbol{\theta}_2.$$

To satisfy this condition, we assume that  $\mathbf{H}$  has full rank  $p$ .

# Minimum Variance Unbiased Estimator for the Linear Model

**Theorem 1.** *For the model*

$$\mathbf{x} = \mathbf{H}\boldsymbol{\theta} + \mathbf{w}$$

where  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ , the MVU estimator of  $\boldsymbol{\theta}$  is given by

$$\hat{\boldsymbol{\theta}} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{x}. \quad (1)$$

The covariance matrix of  $\hat{\boldsymbol{\theta}}$  attains the Cramér-Rao bound (CRB) for all  $\boldsymbol{\theta} \in \mathbb{R}^p$  and is given by

$$\mathbf{C}_{\hat{\boldsymbol{\theta}}} = \mathbb{E} [(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^T] = \sigma^2 (\mathbf{H}^T \mathbf{H})^{-1}.$$

**Proof.** Verifying the unbiasedness of  $\hat{\boldsymbol{\theta}}$  and the covariance matrix expression  $\text{cov}(\hat{\boldsymbol{\theta}}) = \mathbb{E} [(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^T] = \sigma^2 (\mathbf{H}^T \mathbf{H})^{-1}$  proves the theorem.

For the above model,

$$\text{CRB}(\boldsymbol{\theta}) = \mathcal{I}(\boldsymbol{\theta})^{-1}$$

and the Fisher information matrix (FIM) for  $\boldsymbol{\theta}$ ,  $\mathcal{I}(\boldsymbol{\theta})$ , is computed using the general Gaussian FIM expression in handout # 2:

$$[\mathcal{I}(\boldsymbol{\theta})]_{i,k} = \frac{1}{\sigma^2} \cdot \frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})^T}{\partial \theta_i} \frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \theta_k}$$

where  $\boldsymbol{\mu}(\boldsymbol{\theta}) = \mathbf{H}\boldsymbol{\theta}$ . Now

$$\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \theta_i} = \frac{\partial (\mathbf{H}\boldsymbol{\theta})}{\partial \theta_i} = \textit{i}^{\text{th}} \text{ column of } \mathbf{H}$$

implying that

$$\mathcal{I}(\boldsymbol{\theta}) = \frac{1}{\sigma^2} \cdot \mathbf{H}^T \mathbf{H} \quad \Longrightarrow \quad \text{CRB}(\boldsymbol{\theta}) = \sigma^2 (\mathbf{H}^T \mathbf{H})^{-1}. \quad (2)$$

□

## Comments:

- Since the joint FIM and CRB for  $[\boldsymbol{\theta}^T, \sigma^2]^T$  are block-diagonal matrices,  $\boldsymbol{\theta}$  and  $\sigma^2$  are decoupled  $\Longrightarrow$  **CRB( $\boldsymbol{\theta}$ ) is the same regardless of whether  $\sigma^2$  is known or not.** To be more precise, **CRB( $\boldsymbol{\theta}$ )** in (2) is the *CRB for  $\boldsymbol{\theta}$  assuming that  $\sigma^2$  is known* and here is the full CRB for  $\underbrace{\boldsymbol{\theta} \text{ and } \sigma^2}$  for the case

$$\boldsymbol{\rho} = \begin{bmatrix} \boldsymbol{\theta} \\ \sigma^2 \end{bmatrix}$$

where both  $\boldsymbol{\theta}$  and  $\sigma^2$  are unknown:

$$\mathbf{CRB}_{\rho,\rho}(\boldsymbol{\theta}, \sigma^2) = \begin{bmatrix} \overbrace{\mathbf{CRB}_{\boldsymbol{\theta}\boldsymbol{\theta}}(\boldsymbol{\theta}, \sigma^2)}^{\text{same as (2)}} & \mathbf{0} \\ \mathbf{0} & \mathbf{CRB}_{\sigma^2,\sigma^2}(\sigma^2) \end{bmatrix}.$$

Therefore,  $\hat{\boldsymbol{\theta}}$  in (1) is the MVU estimator of  $\boldsymbol{\theta}$  regardless of whether  $\sigma^2$  is known or not.

- $\hat{\boldsymbol{\theta}}$  in (1) coincides with the *least-squares (LS) estimate* of  $\boldsymbol{\theta}$ :

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \|\mathbf{x} - \mathbf{H}\boldsymbol{\theta}\|^2$$

which can be shown by differentiating  $\|\mathbf{x} - \mathbf{H}\boldsymbol{\theta}\|^2$  with respect to  $\boldsymbol{\theta}$  and setting the result to zero or by completing the squares. Later in this handout, we will see a geometric interpretation of the LS approach.

# Minimum Variance Unbiased Estimator for the Linear Model (cont.)

The solution from the above theorem is numerically not sound as given. It is better to use a  $QR$  factorization, say, briefly outlined below. Suppose that the  $N \times p$  matrix  $\mathbf{H}$  is factored as

$$\mathbf{H} = \mathbf{QR} = [\mathbf{Q}_1 \ \mathbf{Q}_2] \begin{bmatrix} \mathbf{R}_1 \\ 0 \end{bmatrix} = \mathbf{Q}_1 \mathbf{R}_1$$

where  $\mathbf{Q}$  is orthonormal and  $\mathbf{R}_1$  is upper triangular  $p \times p$  (MATLAB: `qr`). Then

$$(\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T = \mathbf{R}_1^{-1} \mathbf{Q}_1^T.$$

Thus,  $\hat{\boldsymbol{\theta}}$  can be obtained by solving the triangular system of equations

$$\mathbf{R}_1 \hat{\boldsymbol{\theta}} = \mathbf{Q}_1^T \mathbf{x}.$$

MATLAB has the “backslash” command for computing the LS solution:

$$\boldsymbol{\theta} = \mathbf{H} \backslash \mathbf{x};$$



# Minimum Variance Unbiased Estimator for the Linear Model, Colored Noise

Suppose that we have colored noise, so that  $w \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{C})$ , where  $\mathbf{C} \neq \mathbf{I}$  is known and positive definite.

We can use *prewhitening* to get back to the old problem (i.e. the white-noise case). We compute the Cholesky factorization of  $\mathbf{C}^{-1}$ :

$$\mathbf{C}^{-1} = \mathbf{D}^T \mathbf{D} \quad \text{Matlab: } \mathbf{D} = \text{inv}(\text{chol}(\mathbf{C}))';$$

(Any other square-root factorization could be used as well.)

Now, define the transformed measurement model:

$$\underbrace{\mathbf{D} \mathbf{x}}_{\mathbf{x}^{\text{transf}}} = \underbrace{\mathbf{D} \mathbf{H}}_{\mathbf{H}^{\text{transf}}} \boldsymbol{\theta} + \underbrace{\mathbf{D} \mathbf{w}}_{\mathbf{w}^{\text{transf}}}.$$

Clearly,  $\mathbf{w}^{\text{transf}} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$  and the problem is reduced to the white-noise case.

## MVU Estimation, Colored Noise (cont.)

**Theorem 2.** For colored Gaussian noise with known covariance  $\mathbf{C}$ , the MVU estimate of  $\boldsymbol{\theta}$  is

$$\hat{\boldsymbol{\theta}} = (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}^{-1} \mathbf{x}.$$

The covariance matrix of  $\hat{\boldsymbol{\theta}}$  attains the CRB and is given by

$$\mathbf{C}_{\hat{\boldsymbol{\theta}}} = (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1}.$$

**Note:**  $\hat{\boldsymbol{\theta}}$  is a weighted LS estimate,

$$\begin{aligned} \hat{\boldsymbol{\theta}} &= \arg \min_{\boldsymbol{\theta}} \|\mathbf{x} - \mathbf{H}\boldsymbol{\theta}\|_{\mathbf{W}}^2 \\ &= \arg \min_{\boldsymbol{\theta}} (\mathbf{x} - \mathbf{H}\boldsymbol{\theta})^T \mathbf{W} (\mathbf{x} - \mathbf{H}\boldsymbol{\theta}). \end{aligned}$$

The “optimal weight matrix,”  $\mathbf{W} = \mathbf{C}^{-1}$ , prewhitens the residuals.

# Best Linear Unbiased Estimator

Given the model

$$\mathbf{x} = \mathbf{H}\boldsymbol{\theta} + \mathbf{w} \quad (3)$$

where  $\mathbf{w}$  has zero mean and covariance matrix  $\mathbb{E}[\mathbf{w}\mathbf{w}^T] = \mathbf{C}$ , we look for the *best linear unbiased estimator (BLUE)*. Hence, we restrict our estimator to be

- linear (i.e. of the form  $\hat{\boldsymbol{\theta}} = \mathbf{A}^T \mathbf{x}$ ) and
- unbiased

and minimize its variance.

**Theorem 3.** (*Gauss-Markov*) *The BLUE of  $\boldsymbol{\theta}$  is*

$$\hat{\boldsymbol{\theta}} = (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}^{-1} \mathbf{x}$$

*and its covariance matrix is*

$$\mathbf{C}_{\hat{\boldsymbol{\theta}}} = (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1}.$$

The expression for  $\mathbf{C}_{\hat{\boldsymbol{\theta}}}$  holds *independently* of the distribution of  $\mathbf{w}$  — all we impose on  $\mathbf{w}$  is that it has known mean vector and covariance matrix, equal to  $\mathbf{0}$  and  $\mathbf{C}$  (respectively).

The estimate  $\hat{\theta}$  is (statistically) efficient if  $w$  is Gaussian (i.e. it attains the CRB), but it is not efficient in general. For non-Gaussian measurement models, there might be a better nonlinear estimate. (Most likely, there exists a better nonlinear estimate.)

**Proof. (of Theorem 3).** For simplicity, consider first the case where  $\theta$  is scalar. Then, our measurement model is

$$x[n] = h[n] \theta + w[n] \quad \Longleftrightarrow \quad \mathbf{x} = \mathbf{h} \theta + \mathbf{w}.$$

The candidate linear estimates of  $\theta$  have the following form:

$$\hat{\theta} = \sum_{n=0}^{N-1} a_n x[n] = \mathbf{a}^T \mathbf{x}.$$

First, the bias is computed:

$$\mathbb{E}[\hat{\theta}] = \mathbf{a}^T \mathbb{E}[\mathbf{x}] = \mathbf{a}^T \mathbf{h} \theta.$$

Thus,  $\hat{\theta}$  is unbiased if and only if  $\mathbf{a}^T \mathbf{h} = 1$ . Next, compute the variance of  $\hat{\theta}$ . We have

$$\hat{\theta} - \theta = \mathbf{a}^T (\underbrace{\mathbf{h} \theta + \mathbf{w}}_{\mathbf{x}}) - \theta = \mathbf{a}^T \mathbf{w}$$

where we have used the unbiasedness condition:  $\mathbf{a}^T \mathbf{h} = 1$ .  
Therefore, the variance is

$$\mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}[(\mathbf{a}^T \mathbf{w})^2] = \mathbb{E}[\mathbf{a}^T \mathbf{w} \mathbf{w}^T \mathbf{a}] = \mathbf{a}^T \mathbf{C} \mathbf{a}.$$

**Note:** The variance of  $\hat{\theta}$  depends only on the second-order properties of the noise. This result holds for any noise distribution that has second-order moments.

Thus, the BLUE problem is

$$\min_{\mathbf{a}} \mathbf{a}^T \mathbf{C} \mathbf{a} \quad \text{such that} \quad \mathbf{a}^T \mathbf{h} = 1.$$

Note the equivalence with MVDR beamforming. To read more about MVDR beamforming, see

H.L. Van Trees, *Detection, Estimation and Modulation Theory*, New York: Wiley, 2002, pt. IV.

Lagrange-multiplier formulation:

$$L(\mathbf{a}) = \mathbf{a}^T \mathbf{C} \mathbf{a} + \lambda \cdot (\mathbf{a}^T \mathbf{h} - 1) \quad \xrightarrow{\text{differentiate}} \quad 2 \mathbf{C} \mathbf{a} + \lambda \mathbf{h} = \mathbf{0}.$$

Hence

$$\mathbf{a} = -\frac{\lambda}{2} \cdot \mathbf{C}^{-1} \mathbf{h}$$

and then

$$\mathbf{a}^T \mathbf{h} = -\frac{\lambda}{2} \mathbf{h}^T \mathbf{C}^{-1} \mathbf{h} = 1 \quad \Rightarrow \quad \lambda = -\frac{2}{\mathbf{h}^T \mathbf{C}^{-1} \mathbf{h}}$$

and optimal  $\mathbf{a}$  follows:  $\mathbf{a} = (\mathbf{h}^T \mathbf{C}^{-1} \mathbf{h})^{-1} \mathbf{C}^{-1} \mathbf{h}$ . Returning to our estimator, we find the BLUE to be

$$\hat{\theta} = (\mathbf{h}^T \mathbf{C}^{-1} \mathbf{h})^{-1} \mathbf{h}^T \mathbf{C}^{-1} \mathbf{x}.$$

and its variance is given by

$$\mathbb{E} [(\hat{\theta} - \theta)^2] = (\mathbf{h}^T \mathbf{C}^{-1} \mathbf{h})^{-1}.$$

□

Consider the vector case. Linear unbiased estimates of  $\boldsymbol{\theta}$ :

$$\hat{\boldsymbol{\theta}} = \mathbf{A}^T \mathbf{x}, \quad \text{where } \mathbf{A} \text{ is independent of } \mathbf{x}. \quad (4)$$

**Remark:** For LS estimate  $\hat{\boldsymbol{\theta}}_{\text{LS}}$ ,  $\mathbf{A}^T = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T$ .

$$\boldsymbol{\theta} = \mathbb{E} [\hat{\boldsymbol{\theta}}] = \mathbb{E} [\mathbf{A}^T \mathbf{x}] = \mathbb{E} [\mathbf{A}^T (\mathbf{H}\boldsymbol{\theta} + \mathbf{w})] = \mathbf{A}^T \mathbf{H}\boldsymbol{\theta}.$$

$$\Rightarrow \mathbf{A}^T \mathbf{H} = \mathbf{I}.$$

**Remark:** For BLUE  $\hat{\boldsymbol{\theta}}_{\text{BLUE}}$ ,  $\mathbf{A}_{\text{BLUE}}^T = (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}^{-1}$ .  
Since  $\mathbf{A}_{\text{BLUE}}^T \mathbf{H} = \mathbf{I} \Rightarrow \mathbb{E} [\hat{\boldsymbol{\theta}}_{\text{BLUE}}] = \boldsymbol{\theta}$ .

$$\text{cov}(\hat{\boldsymbol{\theta}}) = \mathbb{E} \{ [\mathbf{A}^T (\underbrace{\mathbf{H}\boldsymbol{\theta} + \mathbf{w}}_{\mathbf{x}}) - \boldsymbol{\theta}] [\mathbf{A}^T (\mathbf{H}\boldsymbol{\theta} + \mathbf{w}) - \boldsymbol{\theta}]^T \} = \mathbf{A}^T \mathbf{C} \mathbf{A}.$$

and

$$\begin{aligned}\text{COV}(\hat{\boldsymbol{\theta}}_{\text{BLUE}}) &= (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}^{-1} \mathbf{C} \mathbf{C}^{-1} \mathbf{H} (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1} \\ &= (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1}.\end{aligned}$$

To prove that  $\hat{\boldsymbol{\theta}}_{\text{BLUE}}$  has the smallest variance [within the family of linear unbiased estimators  $\hat{\boldsymbol{\theta}}$  in (4)], we show that

$$\text{COV}(\hat{\boldsymbol{\theta}}_{\text{BLUE}}) \leq \text{COV}(\hat{\boldsymbol{\theta}})$$

as follows:

$$\begin{aligned}\text{COV}(\hat{\boldsymbol{\theta}}) - \text{COV}(\hat{\boldsymbol{\theta}}_{\text{BLUE}}) &= \mathbf{A}^T \mathbf{C} \mathbf{A} - (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1} \\ &\quad \mathbf{A}^T \mathbf{H} \mathbf{H}^T \mathbf{A} \\ &\stackrel{\mathbf{A}^T \mathbf{H} \mathbf{H}^T \mathbf{A} = \mathbf{I}}{=} \mathbf{A}^T \mathbf{C} \mathbf{A} - \mathbf{A}^T \mathbf{H} (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{A} \\ &= \mathbf{A}^T [\mathbf{C} - \mathbf{H} (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^T] \mathbf{A} \\ &= \mathbf{A}^T [\mathbf{C} - \mathbf{H} (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^T] \mathbf{C}^{-1} [\mathbf{C} - \mathbf{H} (\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^T] \mathbf{A}\end{aligned}$$

which is always positive semidefinite.

## Examples

**Example 4.4 in Kay-I.** Estimate DC level in colored noise:

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

for  $n = 0, 1, \dots, N - 1$ , where  $\mathbf{w} = [w[0], w[1], \dots, w[N - 1]]^T$  is the colored noise with zero mean and covariance matrix  $E[\mathbf{w}\mathbf{w}^T] = \mathbf{C}$ . Hence,  $\mathbf{H} = \mathbf{h} = \mathbf{1} = [1, 1, \dots, 1]^T$  in (3). The BLUE is

$$\hat{A} = (\mathbf{h}^T \mathbf{C}^{-1} \mathbf{h})^{-1} \mathbf{h}^T \mathbf{C}^{-1} \mathbf{x} = \frac{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{x}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}}$$

and its variance is

$$\text{var}(\hat{A}) = \frac{1}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}}.$$

Consider the Cholesky factorization  $(\mathbf{C})^{-1} = \mathbf{D}^T \mathbf{D}$ ; then the BLUE of  $A$  becomes

$$\hat{A} = \frac{\mathbf{1}^T \mathbf{D}^T \mathbf{D} \mathbf{x}}{\mathbf{1}^T \mathbf{D}^T \mathbf{D} \mathbf{1}} = \frac{(\mathbf{D}\mathbf{1})^T \overbrace{\mathbf{D}\mathbf{x}}^{\mathbf{x}^{\text{transf}}}}{\mathbf{1}^T \mathbf{D}^T \mathbf{D} \mathbf{1}} = \sum_{n=0}^{N-1} d_n x^{\text{transf}}[n]$$

where

$$d_n = [\mathbf{D}\mathbf{1}]_n / \mathbf{1}^T \mathbf{D}^T \mathbf{D} \mathbf{1}.$$



## Examples (cont.)

Sometimes, BLUE is completely wrong. For example,  $x[n] = w[n]$ ,  $n = 1, 2, \dots, N$ , white Gaussian noise with variance  $\sigma^2$ . The MVU estimator is  $\hat{\sigma}^2 = (1/N) \cdot \sum_{n=0}^{N-1} x^2[n]$ . On the other hand,

$$\hat{\sigma}_{\text{BLUE}}^2 = \sum_{n=1}^N a_n x[n].$$

For an estimator  $\hat{\sigma}^2$  to be unbiased, we need  $\mathbb{E}[\hat{\sigma}^2] = \sigma^2$ , but

$$\mathbb{E}[\hat{\sigma}_{\text{BLUE}}^2] = \sum_{n=1}^N a_n \mathbb{E}(x[n]) = 0!$$

It is impossible to find  $a_n$ s to make  $\hat{\sigma}_{\text{BLUE}}^2$  unbiased.

**Note:** Although the BLUE is not suitable for this problem, utilizing the *transformed data*  $y[n] = x^2[n]$  would produce a viable estimator.

# General MVU Estimation

What is the MVU estimate in general?

**Theorem 4. (Rao-Blackwell)** *If  $\tilde{\theta}(\mathbf{x})$  is any unbiased estimator and  $T(\mathbf{x})$  is a sufficient statistic, then*

$$\hat{\theta}(\mathbf{x}) = \mathbb{E}_{p(\mathbf{x} | T(\mathbf{x}))} \left[ \tilde{\theta}(\mathbf{X}) | T(\mathbf{X}) = T(\mathbf{x}) \right] \quad (5)$$

*is no worse than  $\tilde{\theta}(\mathbf{x})$  (in terms of MSE).*

**Problem:** Computing  $\mathbb{E}[\tilde{\theta}(\mathbf{X}) | T(\mathbf{X}) = T(\mathbf{x})]$  may be difficult! Recall that this type of expectation occurs when proving sufficiency, but luckily, in the case of sufficiency, our efforts were greatly simplified by the factorization theorem.

**Definition.**  $T(\mathbf{x})$  is *complete sufficient statistic* if only one estimator  $\hat{\theta} = \mathbf{g}(T(\mathbf{x}))$  is unbiased.

**Corollary:** If  $T(\mathbf{x})$  is a complete sufficient statistic, then the unique unbiased estimate  $\hat{\theta} = \mathbf{g}(T(\mathbf{x}))$  is the MVU estimate.

**Comments:**

- Conditioning always decreases the variance (does not increase, to be more precise).

- To get a realizable estimator, we need to condition on the sufficient statistics. The definition of sufficient statistic [denoted by  $\mathbf{T}(\mathbf{x})$ ] implies that conditioning on it leads to a distribution that is not a function of the unknown parameters  $\theta$ . Hence, (5) is a statistic, i.e. realizable.

**Example:** Suppose that  $x[n]$ ,  $n = 1, 2, \dots, N$  are independent, identically distributed (i.i.d.)  $\mathcal{N}(A, \sigma^2)$  with  $\boldsymbol{\theta} = [A, \sigma^2]^T$ . Then,

$$\begin{aligned} p(\mathbf{x}; \boldsymbol{\theta}) &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left[ -\frac{1}{2\sigma^2} \sum_{n=1}^N (x[n] - A)^2 \right] \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \exp \left\{ -\frac{1}{2\sigma^2} \left[ (N-1) \cdot \frac{1}{N-1} \cdot \sum_{n=1}^N (x[n] - \bar{x})^2 \right. \right. \\ &\quad \left. \left. + N(\bar{x} - A)^2 \right] \right\}. \end{aligned}$$

Therefore, the jointly sufficient statistics are

$$T_1(\mathbf{x}) = \bar{x}, \quad T_2(\mathbf{x}) = \hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^N (x[n] - \bar{x})^2.$$

where  $\bar{x} = \frac{1}{N} \sum_{n=1}^N x[n]$ . It can be shown that  $\hat{A} = \bar{x}$  and  $\hat{\sigma}^2 = \frac{1}{N-1} \sum_{n=1}^N (x[n] - \bar{x})^2$  are the only unbiased functions of  $\mathbf{T}(\mathbf{x}) = [T_1(\mathbf{x}), T_2(\mathbf{x})]^T$ . Hence, the corollary at the previous page implies that they are the MVU estimates (although, in this case, the MVU estimates are *not efficient* and, therefore, could not have been found using the efficiency argument). Indeed, for  $\hat{\boldsymbol{\theta}} = [\hat{A}, \hat{\sigma}^2]^T$ ,

$$\text{cov}(\hat{\boldsymbol{\theta}}) = \mathbf{C}_{\hat{\boldsymbol{\theta}}} = \begin{bmatrix} \sigma^2/N & 0 \\ 0 & 2\sigma^4/(N-1) \end{bmatrix}$$

but, recall the CRB for this case (e.g. p. 24 of handout # 2):

$$\mathbf{CRB}(\boldsymbol{\theta}) = \begin{bmatrix} \sigma^2/N & 0 \\ 0 & 2\sigma^4/N \end{bmatrix}.$$

# MAXIMUM LIKELIHOOD (ML) ESTIMATION

$$\hat{\theta} = \arg \max_{\theta} p(\mathbf{x}; \theta).$$

The pdf  $p(\mathbf{x}; \theta)$ , viewed as function of  $\theta$ , is the *likelihood function*.

**Comments on the likelihood function:** For given  $\theta$  and discrete case,  $p(\mathbf{x}; \theta)$  is the probability of observing the point  $\mathbf{x}$ . In the continuous case, it is approximately proportional to probability of observing a point in a small rectangle around  $\mathbf{x}$ . However, when we think of  $p(\mathbf{x}; \theta)$  as a function of  $\theta$ , it gives, for a given observed  $\mathbf{x}$ , the “likelihood” or “plausibility” of various  $\theta$ .

ML estimate  $\equiv$  value of the parameter  $\theta$  that “makes the probability of the data as great as it can be under the assumed model.”

## ML Estimation (cont.)

**Theorem 5.** Assume that certain regularity conditions hold and let  $\hat{\theta}$  be the ML estimate. Then, as  $N \rightarrow \infty$ ,

$$\hat{\theta} \rightarrow \theta_0 \quad (\text{with probability 1}) \quad (\text{consistency}) \quad (6)$$

$$\sqrt{N} (\hat{\theta} - \theta_0) \xrightarrow{d} \mathcal{N}(0, N \mathcal{I}^{-1}(\theta_0)) \quad (\text{asymptotic efficiency}) \quad (7)$$

where  $\theta_0$  is the true value of the parameter and  $\mathcal{I}(\theta_0)$  is the Fisher information [and  $\mathcal{I}^{-1}(\theta_0)$  the CRB]. Moreover, if an efficient (in finite samples) estimate exists, it is given by the ML estimate.

**Proof.** See e.g. Rao, Chapter 5f.2 at pp. 364–366 for the case of independent observations.  $\square$

**Note:** At lower signal-to-noise ratios (SNRs), a threshold effect occurs — outliers give rise to increased variance (more than predicted by the CRB). This behavior is characteristic of practically all nonlinear estimators.

**Example:**  $x[n]$  i.i.d.  $\mathcal{N}(\theta, \sigma^2)$ ,  $n = 0, 1, \dots, N - 1$ , for  $\sigma^2$

known. Maximizing  $p(\mathbf{x}; \theta)$  is equivalent to

$$\max_{\theta} \log p(\mathbf{x}; \theta) = \text{const} - \frac{1}{2\sigma^2} \sum_{n=0}^{N-1} (x[n] - \theta)^2.$$

Thus, the ML estimate is the sample mean

$$\hat{\theta} = \frac{1}{N} \sum_{n=0}^{N-1} x[n] \sim \mathcal{N}(\theta_0, \sigma^2/N).$$

In this example, ML estimator = MVU estimator = BLUE.

**Note:** When estimation error cannot be made small as  $N \rightarrow \infty$ , the asymptotic pdf in (7) is invalid. For asymptotics to work, there has to be an averaging effect!

**Example 7.7 in Kay-I:** Estimation of the DC level in fully dependent non-Gaussian noise:

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

We observe  $x[0], x[1], \dots, x[N-1]$  but  $w[0] = w[1] = \dots = w[N-1]$ , i.e. all noise samples are the same. Hence, we discard  $x[1], x[2], \dots, x[N-1]$ . Then,  $\hat{A} = x[0]$ , say. The pdf of  $\hat{A}$  remains non-Gaussian as  $N \rightarrow \infty$ . Also,  $\hat{A}$  is not consistent since  $\text{var}(\hat{A}) = \text{var}(x[0]) \not\rightarrow 0$  as  $N \rightarrow \infty$ .



# ML Estimation: Vector Parameters

Nothing really changes. The ML estimate is

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} p(\mathbf{x}; \boldsymbol{\theta}).$$

Under appropriate regularity conditions, this estimate is consistent and

$$\sqrt{N} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} \mathcal{N}(\mathbf{0}, N \mathcal{I}^{-1}(\boldsymbol{\theta}_0))$$

where  $\mathcal{I}(\boldsymbol{\theta}_0)$  is the Fisher information *matrix* now.

# ML Estimation: Properties

**Theorem 6. (ML Invariance Principle)** *The ML estimate of  $\alpha = g(\theta)$  where the pdf/pmf  $p(\mathbf{x}; \theta)$  is parametrized by  $\theta$ , is given by*

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

*where  $\hat{\theta}$  is the ML estimate of  $\theta$  [obtained by maximizing  $p(\mathbf{x}; \theta)$  with respect to  $\theta$ ].*

## Comments:

- For a more precise formulation, see Theorems 7.2 and 7.4 in Kay-I.
- Invariance is often combined with the *delta method* which we introduce later in this handout.

## More properties:

- If a given scalar parameter  $\theta$  has a single sufficient statistic  $T(\mathbf{x})$ , say, then the ML estimate of  $\theta$  must be a function of  $T(\mathbf{x})$ . Furthermore, if  $T(\mathbf{x})$  is minimal and complete, then the ML estimate is unique.
- **(Connection between ML and MVU estimation)** If the ML estimate is unbiased, then it is MVU.

## Statistical Motivation

ML has a nice intuitive interpretation, but is it justifiable statistically? Now we try to add to the answer to this question, focusing on the case of i.i.d. observations.

In Ch. 6.2 of *Theory of Point Estimation*, Lehmann shows the following result.

**Theorem 7.** *Suppose that the random observations  $X_i$  are i.i.d. with common pdf/pmf  $p(x_i; \theta_0)$  where  $\theta_0$  is in the interior of the parameter space. Then, as  $N \rightarrow \infty$*

$$P\left\{ \prod_{i=1}^N p(X_i; \theta_0) > \prod_{i=1}^N p(X_i; \theta) \right\} \rightarrow 1$$

for any fixed  $\theta \neq \theta_0$ .

**Comment:** This theorem states that, for large number of i.i.d. samples (i.e. large  $N$ ), the joint pdf/pmf of  $X_1, X_2, \dots, X_N$  at the true parameter value

$$\prod_{i=1}^N p(X_i; \theta_0)$$

exceeds the joint pdf/pmf of  $X_1, X_2, \dots, X_N$  at any other parameter value (with probability one). Consequently, as the

number of observations increases, the parameter estimate that maximizes the joint distribution of the measurements (i.e. the ML estimate) must become close to the true value.

# Regularity Conditions for I.I.D. Observations

Not one set of regularity conditions applies to all scenarios.

Here are some typical regularity conditions for the i.i.d. case. Suppose  $X_1, \dots, X_n$  are i.i.d. with pdf

$$p(x_i; \boldsymbol{\theta}), \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix}.$$

## Regularity conditions:

- (i)  $p(x; \boldsymbol{\theta})$  is identifiable for  $\boldsymbol{\theta}$  and the support of  $p(x_i; \boldsymbol{\theta})$  is not a function of  $\boldsymbol{\theta}$ ;
- (ii) The true value of the parameter, say  $\boldsymbol{\theta}_0$ , lies in an open subset of the parameter space  $\Theta$ ;
- (iii) For almost all  $x$ , the pdf  $p(x; \boldsymbol{\theta})$  has continuous derivatives to order three with respect to all elements of  $\boldsymbol{\theta}$  and all values in the open subset of (ii);
- (iv) The following are satisfied:

$$\mathbb{E}_{p(x;\theta)} \left[ \frac{\partial}{\partial \theta_k} \log p(X; \boldsymbol{\theta}) \right] = 0, \quad k = 1, 2, \dots, p$$

and

$$\begin{aligned} \mathcal{I}_{i,k}(\boldsymbol{\theta}) &= \mathbb{E}_{p(x;\boldsymbol{\theta})} \left[ \frac{\partial}{\partial \theta_i} \log p(X; \boldsymbol{\theta}) \cdot \frac{\partial}{\partial \theta_k} \log p(X; \boldsymbol{\theta}) \right] \\ &= -\mathbb{E} \left[ \frac{\partial^2}{\partial \theta_i \partial \theta_k} \log p(X; \boldsymbol{\theta}) \right], \quad i, k = 1, 2, \dots, p. \end{aligned}$$

**(v)** The FIM  $\mathcal{I}(\boldsymbol{\theta}) = [\mathcal{I}(\boldsymbol{\theta})]_{i,k}$  is positive definite;

**(vi)** Bounding functions  $m_{i,k,l}(\cdot)$  exist such that

$$\left| \frac{\partial^3}{\partial \theta_i \partial \theta_k \partial \theta_l} \log p(x; \boldsymbol{\theta}) \right| \leq m_{i,k,l}(x)$$

for all  $\boldsymbol{\theta}$  in the open subset of (ii), and

$$\mathbb{E}_{p(x;\boldsymbol{\theta})} [m_{i,k,l}(X)] < \infty.$$

**Theorem 8.** (*≈ same as Theorem 5*) *If  $X_1, X_2, \dots, X_N$  are i.i.d. with pdf  $p(x_i; \boldsymbol{\theta})$  such that the conditions (i)–(vi) hold, then there exists a sequence of solutions  $\{\hat{\boldsymbol{\theta}}_N\}$  to the likelihood equations such that*

**(i)**  $\hat{\boldsymbol{\theta}}_N$  is consistent for  $\boldsymbol{\theta}$ ;

**(ii)**  $\sqrt{N} (\hat{\boldsymbol{\theta}}_N - \boldsymbol{\theta})$  is asymptotically Gaussian with mean  $\mathbf{0}$  and covariance matrix  $N \mathcal{I}^{-1}(\boldsymbol{\theta}) = N \mathbf{CRB}(\boldsymbol{\theta})$ ;

**(iii)**  $\sqrt{N} ([\hat{\theta}_N]_i - \theta_i)$  is asymptotically Gaussian with mean 0 and variance  $N [\mathcal{I}^{-1}(\boldsymbol{\theta})]_{i,i}$ ,  $i = 1, 2, \dots, p$ .

**Comments:** What we are *not* given:

**(i)** uniqueness of  $\hat{\boldsymbol{\theta}}_N$ ;

**(ii)** existence for all  $x_1, \dots, x_N$ ;

**(iii)** even if the solution exists and is unique, that we can find it.

## An Array Processing Example

$$\mathbf{x}[n] = \mathbf{A}(\boldsymbol{\phi})\mathbf{s}[n] + \mathbf{w}[n], \quad n = 0, 1, \dots, N - 1$$

where  $\boldsymbol{\theta} = [\boldsymbol{\phi}^T, \mathbf{s}[0]^T, \mathbf{s}[1]^T, \dots, \mathbf{s}[N - 1]^T]^T$  is the vector of unknown parameters and  $\mathbf{w}[n]$  is complex WGN.

### Note:

- $\mathbf{x}[n]$  are *not i.i.d.* (conditions that we stated are not enough);
- $\boldsymbol{\theta}$  grows with  $N$ .

It is well known that CRB cannot be attained asymptotically in this case, see

P. Stoica and A. Nehorai, "MUSIC, maximum likelihood and Cramér-Rao bound," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 37, pp. 720-741, May 1989.

**What if  $\mathbf{s}[n]$  are random  $\sim \mathcal{N}_c(\mathbf{0}, \boldsymbol{\Gamma})$ ?** Then,  $\mathbf{x}[n]$ ,  $n = 0, 1, \dots, N - 1$  are i.i.d. with

$$\mathbf{x}[n] \sim \mathcal{N}(\mathbf{0}, \mathbf{A}(\boldsymbol{\theta})\boldsymbol{\Gamma}\mathbf{A}(\boldsymbol{\theta})^H + \sigma^2\mathbf{I}).$$

Here, the number of parameters *does not grow*. If the regularity conditions that we stated for the i.i.d. case hold, the CRB will be attained asymptotically! Also, the CRB for this case will be different (smaller) than the CRB for deterministic  $\mathbf{s}[n]$ .



## Digression: Delta Method

**Theorem 9.** (*Gauss Approximation Formula, Delta Method*)  
Assume  $\alpha = g(\boldsymbol{\theta})$  has bounded derivatives up to the 2nd order.  
Then, if  $\hat{\boldsymbol{\theta}}$  is consistent, so is  $\hat{\alpha}$ . Moreover, the asymptotic MSE matrices  $C_{\hat{\boldsymbol{\theta}}}$  and  $C_{\hat{\alpha}}$  are related by

$$C_{\hat{\alpha}} = \frac{\partial \mathbf{g}}{\partial \boldsymbol{\theta}^T} C_{\hat{\boldsymbol{\theta}}} \frac{\partial \mathbf{g}^T}{\partial \boldsymbol{\theta}}.$$

**Proof.** Follows from Taylor expansion [around the true value  $\alpha_0 = g(\boldsymbol{\theta}_0)$ ]

$$\hat{\alpha} = g(\boldsymbol{\theta}_0) + \frac{\partial \mathbf{g}}{\partial \boldsymbol{\theta}^T} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) + o(\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\|).$$

□

## Example: Amplitude and Phase Estimation

Assume  $x[n] = A \cos(\omega_0 n + \phi) + e[n]$ ,  $n = 0, 1, \dots, N - 1$ , where  $\omega_0$  is known and  $e[n]$  is additive white Gaussian noise (AWGN). We wish to estimate  $A$  and  $\phi$ .

We rewrite this model as a linear model:

$$\mathbf{x} = \begin{bmatrix} x[0] \\ \dots \\ x[N-1] \end{bmatrix} = \mathbf{H}\boldsymbol{\theta} + \mathbf{w}$$

where  $\boldsymbol{\theta} = [\theta_1, \theta_2]^T$  and

$$H_{i,1} = \cos[\omega_0(i-1)], \quad i = 1, 2, \dots, N$$

$$H_{i,2} = \sin[\omega_0(i-1)], \quad i = 1, 2, \dots, N$$

$$(A \cos \phi, -A \sin \phi) \leftrightarrow (\theta_1, \theta_2)$$

We have

$$\hat{\boldsymbol{\theta}} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{x} \sim \mathcal{N}(\boldsymbol{\theta}_0, \sigma^2 (\mathbf{H}^T \mathbf{H})^{-1}).$$

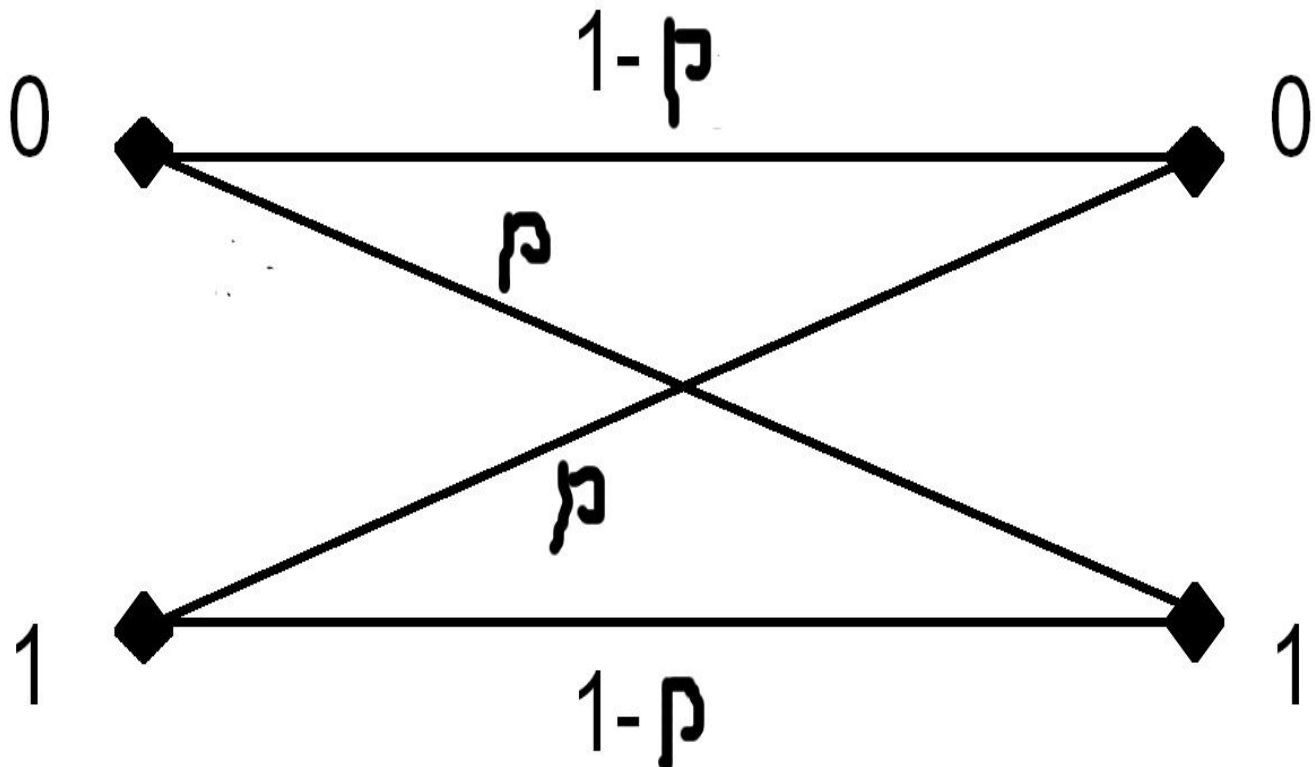
By the ML invariance principle,  $\hat{A}$  and  $\hat{\phi}$  can be found from  $\hat{\boldsymbol{\theta}} = [\hat{\theta}_1, \hat{\theta}_2]^T$  via rectangular-to-polar coordinate conversion:

$$(\hat{\theta}_1, \hat{\theta}_2) \leftrightarrow (\hat{A} \cos \hat{\phi}, -\hat{A} \sin \hat{\phi}).$$

Define  $\boldsymbol{\alpha} = [A, \phi]^T = \mathbf{g}(\boldsymbol{\theta})$ . Then, the delta method yields

$$\mathbf{C}_{\hat{\boldsymbol{\alpha}}} = \frac{\partial \mathbf{g}}{\partial \boldsymbol{\theta}^T} \mathbf{C}_{\hat{\boldsymbol{\theta}}} \frac{\partial \mathbf{g}^T}{\partial \boldsymbol{\theta}}.$$

## Example: ML Decoding



For a symmetric channel, the ML decoder is the minimum Hamming distance decoder.

**Proof.** Let  $x$  and  $\theta$  be the received and transmitted vectors from a binary symmetric channel (i.e. the elements of  $x$  and  $\theta$  are zeros or ones). Note that  $\theta$  belongs to a finite set of codewords. We wish to find which  $\theta$  was transmitted based on the received  $x$ . We have

$$x = \theta + w \pmod{2} \triangleq \theta \oplus w$$

where  $w = [w_1, \dots, w_N]^T$  and  $w_i$  are i.i.d. Bernoulli( $p$ ). The

likelihood function is given by

$$\begin{aligned}
 p(\mathbf{x}; \boldsymbol{\theta}) &= P\{\mathbf{X} = \mathbf{x}\} = P\{\boldsymbol{\theta} \oplus \underbrace{\mathbf{W}}_{\text{i.i.d. Bernoulli}} = \mathbf{x}\} \\
 &= P\{\mathbf{W} = \underbrace{\mathbf{x} \oplus \boldsymbol{\theta}}_{\mathbf{w}}\} \\
 &\quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_N \end{bmatrix} \\
 &= p^{\sum_{i=1}^N w_i} \cdot (1-p)^{N - \sum_{i=1}^N w_i} \\
 &= \left(\frac{p}{1-p}\right)^{d_H(\mathbf{x}, \boldsymbol{\theta})} (1-p)^N
 \end{aligned}$$

where

$$d_H(\mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^N x_i \oplus \theta_i$$

is the Hamming distance between  $\mathbf{x}$  and  $\boldsymbol{\theta}$  (i.e. the number of bits that are different between the two vectors). Hence, if  $p < 0.5$ , then

$$\max_{\mathbf{x}} p(\mathbf{x}; \boldsymbol{\theta}) \iff \min_{\boldsymbol{\theta}} d_H(\mathbf{x}, \boldsymbol{\theta}).$$

□

# Asymptotic ML for WSS Processes

Consider data  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{C}(\boldsymbol{\theta}))$ . To find the ML estimate of  $\boldsymbol{\theta}$ , maximize

$$p(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{(2\pi)^{N/2} |\mathbf{C}(\boldsymbol{\theta})|^{1/2}} \exp \left[ -\frac{1}{2} \cdot \mathbf{x}^T \mathbf{C}(\boldsymbol{\theta})^{-1} \mathbf{x} \right]$$

over  $\boldsymbol{\theta}$ .

If  $x[n]$  is WSS, then  $\mathbf{C}(\boldsymbol{\theta})$  is Töplitz, so, as  $N \rightarrow \infty$ , we can approximate the log likelihood as:

$$\begin{aligned} \log p(\mathbf{x}; \boldsymbol{\theta}) &= -\frac{N}{2} \log 2\pi \\ &\quad -\frac{N}{2} \int_{-1/2}^{1/2} \left( \log P_{xx}(f; \boldsymbol{\theta}) + \frac{I_x(f)}{P_{xx}(f; \boldsymbol{\theta})} \right) df \end{aligned}$$

where  $I_x(f)$  is the periodogram:

$$I_x(f) = \frac{1}{N} \left| \sum_{n=0}^{N-1} x[n] e^{-j2\pi f n} \right|^2$$

and  $P_{xx}(f; \boldsymbol{\theta})$  is the PSD of  $x[n]$ .

This result is based on the Whittle approximation, see e.g.

P. Whittle, "The analysis of multiple stationary time series," *J. R. Stat. Soc., Ser. B* vol. 15, pp. 125–139, 1953.

**Proof.** See e.g. Ch. 7.9 in Kay-I.  $\square$

**Note:** Kay calls the Whittle approximation “asymptotic ML”.

The discrete-frequency version of the above expression is also useful:

$$-\frac{N}{2} \log 2\pi - \frac{1}{2} \sum_{k=0}^{N-1} \left\{ \log[P_{yy}(f_k; \boldsymbol{\theta})] + \frac{I_x(f_k)}{P_{yy}(f_k; \boldsymbol{\theta})} \right\} \quad (8)$$

where

$$f_k = k/N, \quad k = 0, 1, \dots, N - 1.$$

For example, (8) may exist even when the integral form does not. An example of such a case that we mentioned earlier is the Doppler PSD (which goes to infinity, causing an integrability problem), see

A. Dogandžić and B. Zhang, “Estimating Jakes’ Doppler power spectrum parameters using the Whittle approximation,” *IEEE Trans. Signal Processing*, vol. 53, pp. 987–1005, Mar. 2005.

**Example.** Autoregressive (AR) parameter estimation:

$$P_{xx}(f; \boldsymbol{\theta}) = \frac{\sigma_u^2}{|A(f; \mathbf{a})|^2}$$

where  $\boldsymbol{\theta} = \underbrace{[a[1], a[2], \dots, a[p]]}_{\mathbf{a}^T}, \sigma_u^2]^T$ , and

$$A(f; \mathbf{a}) = 1 + \sum_{m=1}^p a[m] \exp(-j2\pi fm).$$

So

$$\begin{aligned} \log p(\mathbf{x}; \mathbf{a}, \sigma_u^2) &= -\frac{N}{2} \log 2\pi \\ &\quad - \frac{N}{2} \int_{-1/2}^{1/2} \left( \log \frac{\sigma_u^2}{|A(f; \mathbf{a})|^2} + \frac{I_x(f)}{\frac{\sigma_u^2}{|A(f; \mathbf{a})|^2}} \right) df. \end{aligned}$$

Assuming  $\mathcal{A}(z) = 1 + \sum_{m=1}^p a[m]z^{-m}$  to be minimum-phase [typically required for stability of  $1/\mathcal{A}(z)$ ], then

$$\int_{-1/2}^{1/2} \log |A(f; \mathbf{a})|^2 df = 0$$

see Problem 7.22 in Kay-I. Therefore

$$\begin{aligned} \log p(\mathbf{x}; \mathbf{a}, \sigma_u^2) &= -\frac{N}{2} \log 2\pi - \frac{N}{2} \log \sigma_u^2 \\ &\quad - \frac{N}{2\sigma_u^2} \int_{-1/2}^{1/2} |A(f; \mathbf{a})|^2 I_x(f) df. \end{aligned}$$



Differentiating with respect to  $\sigma_u^2$  and setting the result to zero, we obtain, for a fixed  $\mathbf{a}$ ,

$$\hat{\sigma}_u^2(\mathbf{a}) = \int_{-1/2}^{1/2} |A(f; \mathbf{a})|^2 I_x(f) df.$$

Then, the concentrated Whittle log-likelihood function of  $\mathbf{a}$  is obtained by substituting  $\hat{\sigma}_u^2(\mathbf{a})$  into the Whittle log-likelihood  $\log p(\mathbf{x}; \mathbf{a}, \sigma_u^2)$ :

$$\log p(\mathbf{x}; \mathbf{a}, \hat{\sigma}_u^2(\mathbf{a})) = -\frac{N}{2} \log 2\pi - \frac{N}{2} \log \hat{\sigma}_u^2(\mathbf{a}) - N/2.$$

We will discuss the concentrated likelihood approach later, see p. 49. To find  $\hat{\mathbf{a}}$ , we need to maximize the above concentrated likelihood or, equivalently, minimize  $\hat{\sigma}_u^2$ :

$$\min_{\mathbf{a}} J(\mathbf{a}) = \min_{\mathbf{a}} \int_{-1/2}^{1/2} |A(f; \mathbf{a})|^2 I_x(f) df.$$

The above function is quadratic in  $\mathbf{a}$ , resulting in the global minimum upon differentiation. For  $k = 1, 2, \dots, p$ , we have

$$\frac{\partial J(\mathbf{a})}{\partial a[k]} = \int_{-1/2}^{1/2} \left[ A(f; \mathbf{a}) \underbrace{\frac{\partial A^*(f; \mathbf{a})}{\partial a[k]}}_{\exp(j2\pi f k)} + \underbrace{\frac{\partial A(f; \mathbf{a})}{\partial a[k]}}_{\exp(-j2\pi f k)} A^*(f; \mathbf{a}) \right] I_x(f) df.$$

Since  $A(-f; \mathbf{a}) = A^*(f; \mathbf{a})$  and  $I_x(-f) = I_x(f)$ , we have

$$\frac{\partial J(\mathbf{a})}{\partial a[k]} = 2 \int_{-1/2}^{1/2} A(f; \mathbf{a}) I_x(f) \exp(j2\pi f k) df.$$

Setting the above expression to zero, we get

$$\int_{-1/2}^{1/2} \left[ 1 + \sum_{m=1}^p a[m] \exp(-j2\pi f m) \right] I_x(f) \exp(j2\pi f k) df = 0$$

or

$$\begin{aligned} & \sum_{m=1}^p a[m] \int_{-1/2}^{1/2} I_x(f) \exp[j2\pi f(k - m)] df \\ &= - \int_{-1/2}^{1/2} I_x(f) \exp(j2\pi f k) df. \end{aligned}$$

But,  $\int_{-1/2}^{1/2} I_x(f) \exp(j2\pi f k) df$  is just the inverse discrete-time Fourier transform (DTFT) of the periodogram evaluated at  $k$ , which is equal to the biased sample estimate of the autocorrelation function:

$$\hat{r}_{xx}[k] = \begin{cases} \frac{1}{N} \sum_{i=k}^{N-1} x[i]x[i - |k|], & |k| \leq N - 1, \\ 0, & k \geq N - 1 \end{cases}$$

Hence, the Whittle (asymptotic) ML estimate of the AR filter parameter vector  $\mathbf{a}$  solves

$$\sum_{m=1}^p \hat{a}[m] \hat{r}_{xx}[k-m] = -\hat{r}_{xx}[k], \quad k = 1, 2, \dots, p$$

which are the *estimated Yule-Walker equations*.

## Computing the Estimates

Typically, finding the ML estimate requires a nonlinear  $p$ -dimensional optimization (for  $\boldsymbol{\theta}$  of size  $p$ ). More generally,

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} V(\boldsymbol{\theta})$$

where, for ML estimation, we have  $V(\boldsymbol{\theta}) = -\log p(\mathbf{x}; \boldsymbol{\theta})$ .

**Newton-Raphson Iteration:** Assume that a guess  $\boldsymbol{\theta}^{(i)}$  is available. We wish to improve  $\boldsymbol{\theta}^{(i)}$ , yielding  $\boldsymbol{\theta}^{(i+1)}$ . Let us apply a quadratic Taylor expansion:

$$V(\boldsymbol{\theta}) \approx V(\boldsymbol{\theta}^{(i)}) + \mathbf{g}_i^T \bar{\boldsymbol{\theta}}^{(i)} + \frac{1}{2} \cdot (\bar{\boldsymbol{\theta}}^{(i)})^T \mathbf{H}_i \bar{\boldsymbol{\theta}}^{(i)}$$

where

$$\begin{aligned} \bar{\boldsymbol{\theta}}^{(i)} &= \boldsymbol{\theta} - \boldsymbol{\theta}^{(i)} \\ \mathbf{g}_i &= \left. \frac{\partial V(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(i)}} \\ \mathbf{H}_i &= \left. \frac{\partial^2 V(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(i)}}. \end{aligned}$$

# Newton-Raphson Iteration

## (Ch. 7.7 in Kay-I)

Complete the squares:

$$V(\boldsymbol{\theta}) \approx (\bar{\boldsymbol{\theta}}^{(i)} + \mathbf{H}_i^{-1} \mathbf{g}_i)^T \frac{1}{2} \mathbf{H}_i (\bar{\boldsymbol{\theta}}^{(i)} + \mathbf{H}_i^{-1} \mathbf{g}_i) + \text{const.}$$

We assume that  $\mathbf{H}_i > 0$

Hessian matrix of  $V(\boldsymbol{\theta})$

(i.e. the second derivative of  $V(\boldsymbol{\theta})$  with respect to  $\boldsymbol{\theta}$ , computed at  $\boldsymbol{\theta}_i$ , is positive definite)

and thus choose

$$\boldsymbol{\theta}^{(i+1)} = \boldsymbol{\theta}^{(i)} - \mathbf{H}_i^{-1} \mathbf{g}_i.$$

Newton-Raphson iteration usually has *quadratic convergence* near the optimum, i.e.

$$\|\boldsymbol{\theta}^{(i+1)} - \boldsymbol{\theta}\| \leq c \|\boldsymbol{\theta}^{(i)} - \boldsymbol{\theta}\|^2.$$

where  $c$  is a positive constant. Therefore, we gain approximately one significant digit per iteration.

However, the algorithm can diverge if we start too far from the optimum. To facilitate convergence (to a *local optimum*, in general), we can apply a damped Newton-Raphson algorithm.

Here is one such damped algorithm:

$$\boldsymbol{\theta}^{(i+1)} = \boldsymbol{\theta}^{(i)} - \mu_i \cdot \mathbf{H}_i^{-1} \mathbf{g}_i \quad (9)$$

where the step length  $\mu_i$  is  $\mu_i = 1, 1/2, 1/4, \dots$ . In particular, in the  $i$ th iteration, start with the step length  $\mu_i = 1$ , compute  $\boldsymbol{\theta}^{(i+1)}$  using (9), and check if

$$V(\boldsymbol{\theta}^{(i+1)}) < V(\boldsymbol{\theta}^{(i)})$$

holds; if yes, go to the  $(i + 1)$ st iteration. If no, keep halving  $\mu_i$  and recomputing  $\boldsymbol{\theta}^{(i+1)}$  using (9) until  $V(\boldsymbol{\theta}^{(i+1)}) < V(\boldsymbol{\theta}^{(i)})$  holds — then go to the  $(i + 1)$ st iteration. Once in the  $(i + 1)$ st iteration, reset  $\mu^{(i+1)}$  to 1 and continue in the same manner.

**Modification:** Use an approximate form of the Hessian matrix of  $V(\boldsymbol{\theta})$

$$\frac{\partial^2 V(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T}.$$

In the case of ML estimation, use the FIM instead of this Hessian:

$$\mathbf{H}_i = \mathcal{I}(\boldsymbol{\theta}^{(i)})$$

and the resulting algorithm is called *Fisher scoring*. This choice of  $\mathbf{H}_i$  usually guarantees positive definiteness of the Hessian matrix:  $\mathbf{H}_i = \mathcal{I}(\boldsymbol{\theta}^{(i)}) > 0$ .

**Note:** The convergence point is a *local* minimum of  $V(\boldsymbol{\theta})$ . It is a global minimum if  $V(\boldsymbol{\theta})$  is a unimodal function of  $\boldsymbol{\theta}$

or if the initial estimate is sufficiently good. (If we suspect that) there are multiple local minima of  $V(\boldsymbol{\theta})$  (i.e. multiple local maxima of the likelihood function), we should try many (wide-spread/different) starting values for our Newton-Raphson or Fisher scoring iterations.

If the parameter space  $\Theta$  is not  $\mathbf{R}^p$ , we should also examine the boundary of the parameter space to see if a global maximum of the likelihood function (i.e. a global minimum of  $V(\boldsymbol{\theta})$ ) lies on this boundary. **Pay attention to this issue in your homework assignments and exams, as well as in general.**

## Newton-Raphson Iteration: Example

Suppose  $x[n] = s[n; \boldsymbol{\theta}] + e[n]$  where  $e[n]$  is white Gaussian noise with known variance  $\sigma^2$ ,  $s[n; \boldsymbol{\theta}] = \sin(\omega_1 n) + \sin(\omega_2 n)$ ,  $\boldsymbol{\theta} = [\omega_1, \omega_2]^T$ , and  $n = 0, 1, \dots, N - 1$ .

Ignoring constants, we obtain the negative log likelihood:

$$V(\boldsymbol{\theta}) = -\log p(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{2\sigma^2} \sum_{n=0}^{N-1} (x[n] - s[n; \boldsymbol{\theta}])^2 + \text{const}$$

its gradient:

$$V'(\boldsymbol{\theta}) = -\frac{1}{\sigma^2} \sum_{n=0}^{N-1} (x[n] - s[n; \boldsymbol{\theta}]) \cdot \frac{\partial s[n; \boldsymbol{\theta}]}{\partial \boldsymbol{\theta}}$$

and the FIM:

$$\mathcal{I}(\boldsymbol{\theta}) = \frac{1}{\sigma^2} \sum_{n=0}^{N-1} \frac{\partial s[n; \boldsymbol{\theta}]}{\partial \boldsymbol{\theta}} \frac{\partial s[n; \boldsymbol{\theta}]}{\partial \boldsymbol{\theta}^T}.$$

We can use  $\mathbf{H} = \mathcal{I}(\boldsymbol{\theta})$  and the damped Fisher scoring iteration



becomes

$$\begin{aligned}
 \boldsymbol{\theta}^{(i+1)} &= \boldsymbol{\theta}^{(i)} - \mu_i \cdot \mathcal{I}(\boldsymbol{\theta}^{(i)})^{-1} \cdot V'(\boldsymbol{\theta}^{(i)}) \\
 &= \boldsymbol{\theta}^{(i)} \\
 &+ \mu_i \cdot \left\{ \sum_{n=0}^{N-1} n^2 \begin{bmatrix} \cos^2(\omega_1^{(i)} n) & \cos(\omega_1^{(i)} n) \cos(\omega_2^{(i)} n) \\ \cos(\omega_1^{(i)} n) \cos(\omega_2^{(i)} n) & \cos^2(\omega_2^{(i)} n) \end{bmatrix} \right\}^{-1} \\
 &\cdot \sum_{n=0}^{N-1} (x[n] - s[n; \boldsymbol{\theta}^{(i)}]) \begin{bmatrix} n \cos(\omega_1^{(i)} n) \\ n \cos(\omega_2^{(i)} n) \end{bmatrix}.
 \end{aligned}$$

## Example: Concentrated Likelihood

Consider a situation in which a small-scale disease epidemic has been observed, with individuals exposed to the disease (e.g. virus) at a common place and time. Or, in a similar computer-analogous scenario, consider computers infected by a virus. We assume that a time interval is known for exposure, but not the exact time.

We collect times at which infection was detected at various computers ('incubation times'), say, with time 0 corresponding to the start of a known interval in which exposure occurred. Let  $x_1, x_2, \dots, x_n$  be the collected infection times after the exposure. Assume that  $x_1, x_2, \dots, x_n$  are i.i.d. following some distribution. Here, we adopt the following lognormal model:

$$p(x_i; \boldsymbol{\theta}) = \begin{cases} \frac{1}{(x_i - \alpha)\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma^2}[\log(x_i - \alpha) - \mu]^2\right\}, & x_i > \alpha \\ 0, & \text{otherwise} \end{cases}$$

with parameters

$$\boldsymbol{\theta} = [\alpha, \mu, \sigma]^T.$$

where the parameter  $\alpha > 0$  represents the time at which the exposure took place. Since the support of the above distribution depends on the parameter  $\alpha$ , the regularity condition (i) on p. 29 does not hold.

Here are some references related to the above model:

H.L. Harter and A.H. Moore, “Local-maximum-likelihood estimation of the parameters of three-parameter lognormal populations from complete and censored samples,” *J. Amer. Stat. Assoc.*, vol. 61, pp. 842–851, Sept. 1966.

B.M. Hill, “The three-parameter lognormal distribution and Bayesian analysis of a point-source epidemic,” *J. Amer. Stat. Assoc.*, vol. 68, pp. 72–84, Mar. 1963.

**Note:** this model is equivalent to having  $x_i$ ,  $i = 1, 2, \dots, N$  such that  $\log(x_i - \alpha) \sim \mathcal{N}(\mu, \sigma^2)$ .

The log-likelihood function is

$$\begin{aligned} l(\boldsymbol{\theta}) &= \sum_{i=1}^N \log p(x_i; \boldsymbol{\theta}) \\ &= -\frac{N}{2} \log(2\pi\sigma^2) - \sum_{i=1}^N \log(x_i - \alpha) \\ &\quad - \frac{1}{2\sigma^2} \sum_{i=1}^N [\log(x_i - \alpha) - \mu]^2 \end{aligned}$$

where  $x_i > \alpha$ ,  $\forall i = 1, 2, \dots, N$ .

For a fixed  $\alpha$ , we can easily find  $\mu$  and  $\sigma^2$  that maximize the

likelihood:

$$\hat{\mu}(\alpha) = \frac{1}{N} \sum_{t=1}^N \log(x_i - \alpha)$$

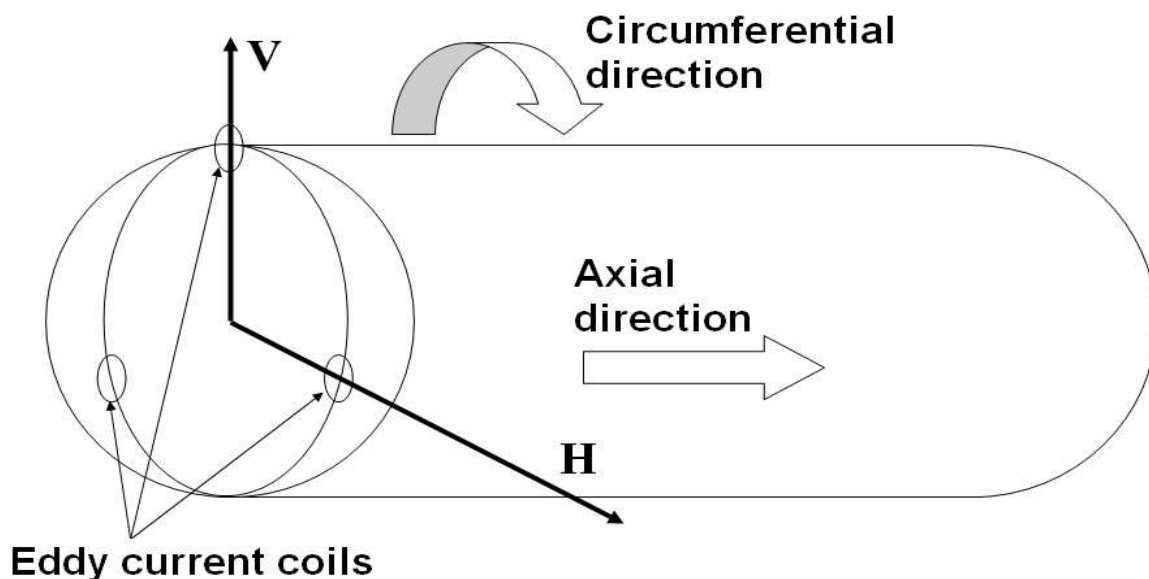
$$\hat{\sigma}^2(\alpha) = \frac{1}{N} \sum_{t=1}^N [\log(x_i - \alpha) - \hat{\mu}(\alpha)]^2$$

which follows from the above relationship with the normal pdf. Now, we can write the log-likelihood function as a function of  $\alpha$  alone, by substituting  $\hat{\mu}(\alpha)$  and  $\hat{\sigma}^2(\alpha)$  into  $l(\boldsymbol{\theta})$ :

$$\begin{aligned} l([\alpha, \hat{\mu}(\alpha), \hat{\sigma}^2(\alpha)]^T) \\ = -\frac{N}{2} \log[2\pi\hat{\sigma}^2(\alpha)] - \sum_{i=1}^N \log(x_i - \alpha) - \frac{N}{2}. \end{aligned}$$

## Example: ML Estimation for Eddy-Current Data

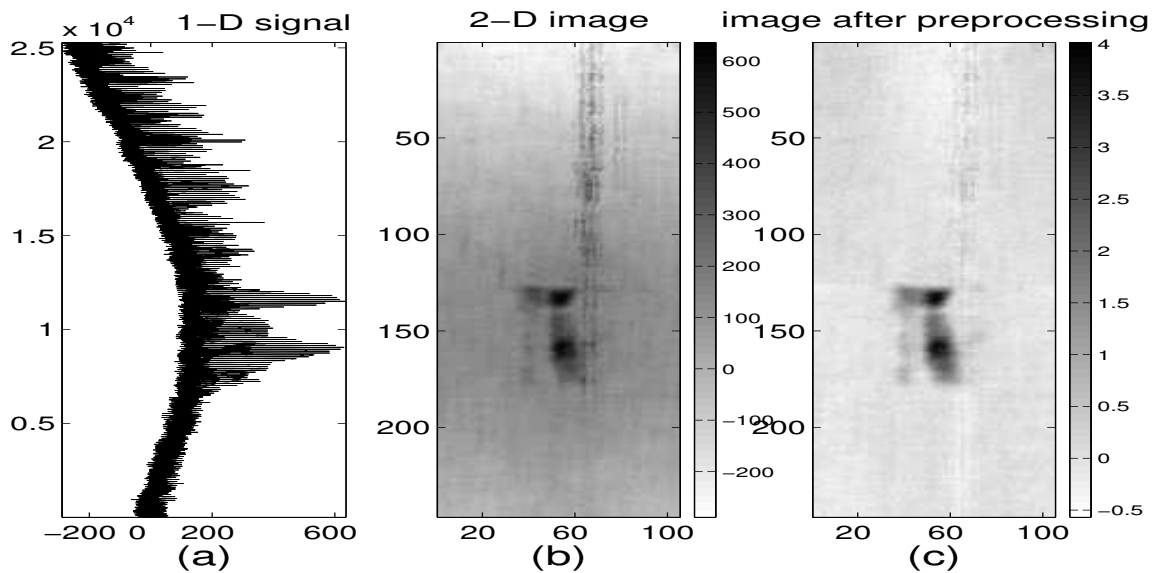
Data collected with a rotating probe, consisting of three coils spaced  $2\pi/3$  rad ( $120^\circ$ ) apart. Each coil scans the inner surface of the tube by moving along a helical path.



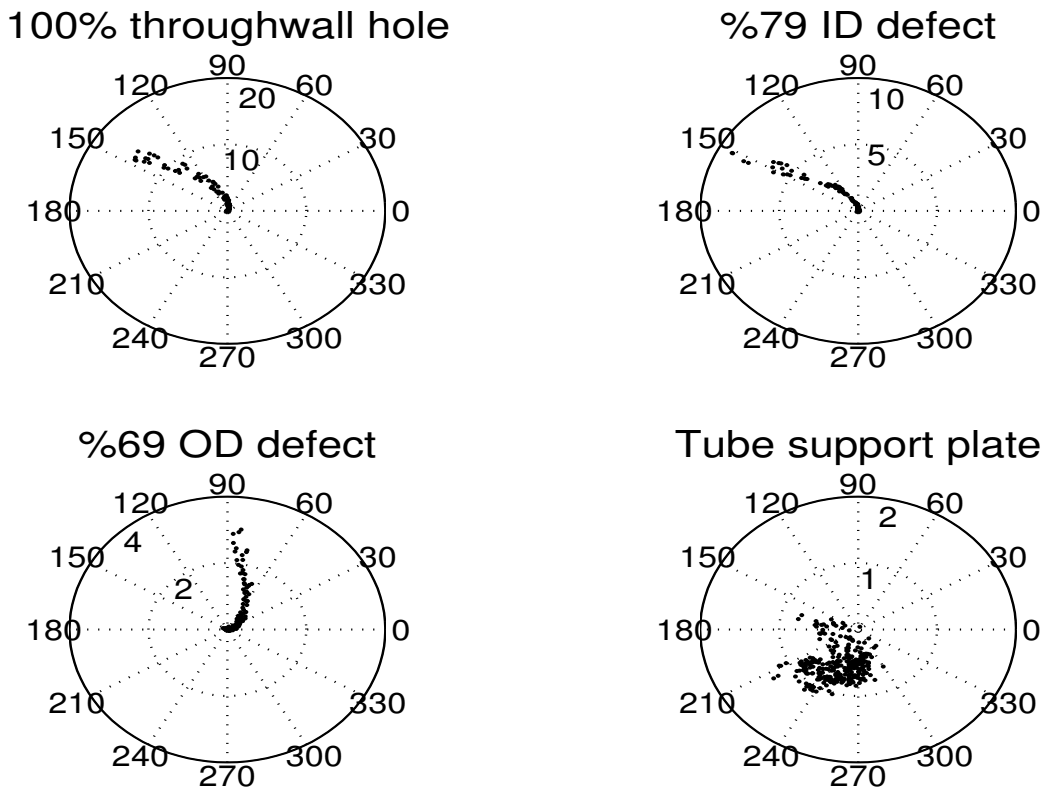
Data acquisition scheme.

A preprocessing step is performed first, with goal to extract meaningful information from the rotating-probe data, see

P. Xiang, S. Ramakrishnan, X. Cai, P. Ramuhalli, R. Polikar, S.S. Udpa, and L. Udpa, "Automated analysis of rotating probe multi-frequency eddy current data from steam generator tubes," *Intl. J. Applied Electrom. Mech.*, vol. 12, pp. 151–164, 2000.



Signal preprocessing: (a) 1-D raw data, (b) 2-D image, and (c) 2-D image after preprocessing.



Impedance-plane plots of typical signals measured by the rotating-probe eddy-current system.

**Objective:** Characterize the amplitude and phase probability distributions of the potential defects. The estimated distribution parameters can be used for:

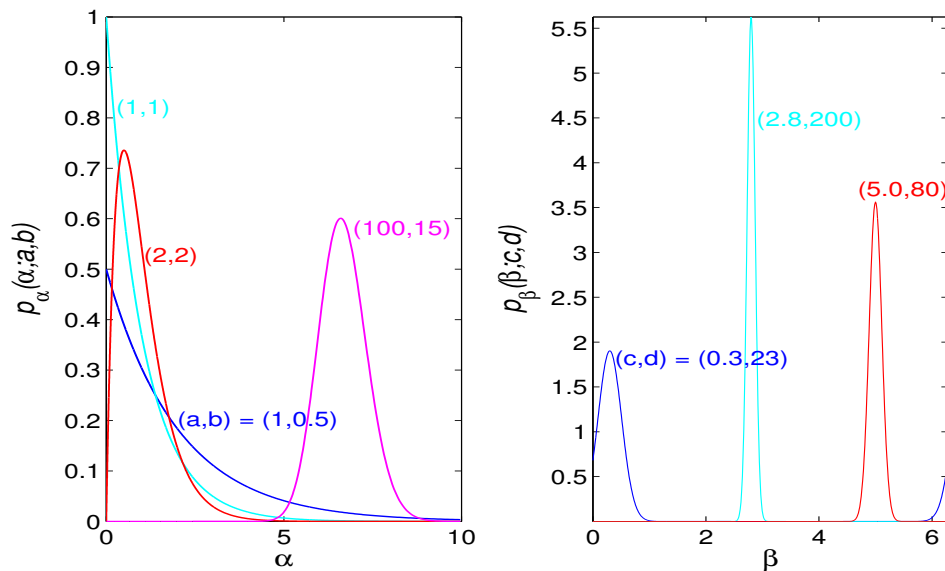
- defect detection,
- defect classification and characterization, e.g. discriminating between inner diameter (ID) and outer diameter (OD) defects,
- denoising.

**Statistical Model:**  $K$  complex measurements  $x[k]$ ,  $k = 0, 1, \dots, K - 1$  at the defect location modeled as

$$x[k] = \sqrt{\alpha_k} \cdot e^{j\beta_k} + e[k]$$

where

- $\alpha_k \equiv$  i.i.d. signal powers following a Gamma( $a, b$ ) distribution (Interestingly, in the special case where  $a = 1$ , the amplitudes  $\sqrt{\alpha_k}$  follow a Rayleigh distribution.),
- $\beta_k \equiv$  i.i.d. signal phases independent of the amplitudes, where  $\beta_k$  follow a von Mises distribution (with parameters  $c$  and  $d$ ),
- $e[k]$  i.i.d. zero-mean complex Gaussian noise independent of the signal amplitudes and phases, having known variance  $\sigma^2$ . [The noise variance  $\sigma^2$  can be estimated from the neighboring pixels that contain only noise.]



Gamma distribution pdf (left)  
and von Mises distribution pdf (right).

### Pdfs of $\alpha_k$ and $\beta_k$ :

$$p_{\alpha}(\alpha_k; a, b) = \frac{b^a}{\Gamma(a)} \alpha_k^{a-1} \exp(-b\alpha_k), \quad \alpha_k > 0, \quad a, b > 0$$

and

$$p_{\beta}(\beta_k; c, d) = \frac{1}{2\pi I_0(d)} \exp[d \cos(\beta_k - c)], \quad 0 < \beta_k \leq 2\pi, \quad d > 0$$

where  $I_0(\cdot)$  denotes the modified Bessel function of the first kind and order zero. Define the unknown parameter vector

$$\boldsymbol{\lambda} = [a, b, c, d]^T$$



and the vectors of signal amplitudes and phases

$$\boldsymbol{\theta}_k = [\alpha_k, \beta_k]^T, \quad k = 0, 1, \dots, K - 1.$$

Marginal distribution of the  $k$ th observation:

$$p_x(x[k]; \boldsymbol{\lambda}) = \int_{\Theta} p_{x|\theta}(x[k] | \boldsymbol{\theta}) p_{\theta}(\boldsymbol{\theta}; \boldsymbol{\lambda}) d\boldsymbol{\theta}, \quad k = 0, 1, \dots, K - 1$$

where  $\boldsymbol{\theta} = [\alpha, \beta]^T$ ,  $\Theta = \{(\alpha, \beta) : 0 < \alpha < \infty, 0 < \beta < 2\pi\}$ , and

$$p_{x|\theta}(x[k] | \boldsymbol{\theta}) = \frac{1}{\pi\sigma^2} \exp \left[ -\frac{|x[k] - \sqrt{\alpha} \cdot e^{j\beta}|^2}{\sigma^2} \right]$$

$$p_{\theta}(\boldsymbol{\theta}; \boldsymbol{\lambda}) = p_{\alpha}(\alpha; a, b) p_{\beta}(\beta; c, d).$$

ML estimate of  $\boldsymbol{\lambda}$  obtained by maximizing the log-likelihood of  $\boldsymbol{\lambda}$  for all measurements  $\boldsymbol{x} = [x[0], x[1], \dots, x[K - 1]]^T$ :

$$L(\boldsymbol{\lambda}, \boldsymbol{y}) = \sum_{k=0}^{K-1} \log p_x(x[k]; \boldsymbol{\lambda}). \quad (10)$$

Newton-Raphson iteration for finding the ML estimates of  $\boldsymbol{\lambda}$  [i.e. maximizing (10)]:

$$\boldsymbol{\lambda}^{(i+1)} = \boldsymbol{\lambda}^{(i)} - \delta^{(i)} \cdot \left[ \underbrace{\frac{\partial^2 L(\boldsymbol{\lambda}^{(i)})}{\partial \boldsymbol{\lambda} \partial \boldsymbol{\lambda}^T}}_{\text{Hessian of the log likelihood}} \right]^{-1} \underbrace{\frac{\partial L(\boldsymbol{\lambda}^{(i)})}{\partial \boldsymbol{\lambda}}}_{\text{gradient}} \quad (11)$$

where the *damping factor*  $0 < \delta^{(i)} \leq 1$  is chosen (at every step  $i$ ) to ensure that the likelihood function (10) increases and the parameter estimates remain in the allowable parameter space ( $a, b, d > 0$ ).

We utilized the following formulas to compute the gradient vector and Hessian matrix in (11):

$$\begin{aligned} \frac{\partial}{\partial \lambda_i} \{\log p_x(x; \boldsymbol{\lambda})\} &= \frac{1}{p_x(x; \boldsymbol{\lambda})} \int_{\Theta} p_{x|\theta}(x|\boldsymbol{\theta}) \frac{\partial p_{\theta}(\boldsymbol{\theta}; \boldsymbol{\lambda})}{\partial \lambda_i} d\boldsymbol{\theta} \\ \frac{\partial^2}{\partial \lambda_i \partial \lambda_m} \{\log p_x(x; \boldsymbol{\lambda})\} &= \frac{1}{p_x(x; \boldsymbol{\lambda})} \int_{\Theta} p_{x|\theta}(x|\boldsymbol{\theta}) \frac{\partial^2 p_{\theta}(\boldsymbol{\theta}; \boldsymbol{\lambda})}{\partial \lambda_i \partial \lambda_m} d\boldsymbol{\theta} \\ &\quad - \frac{1}{[p_x(x; \boldsymbol{\lambda})]^2} \\ &\quad \cdot \int_{\Theta} p_{x|\theta}(x|\boldsymbol{\theta}) \frac{\partial p_{\theta}(\boldsymbol{\theta}; \boldsymbol{\lambda})}{\partial \lambda_i} d\boldsymbol{\theta} \cdot \int_{\Theta} p_{x|\theta}(x|\boldsymbol{\theta}) \frac{\partial p_{\theta}(\boldsymbol{\theta}; \boldsymbol{\lambda})}{\partial \lambda_m} d\boldsymbol{\theta} \end{aligned}$$

for  $i, m = 1, 2, 3, 4$ .

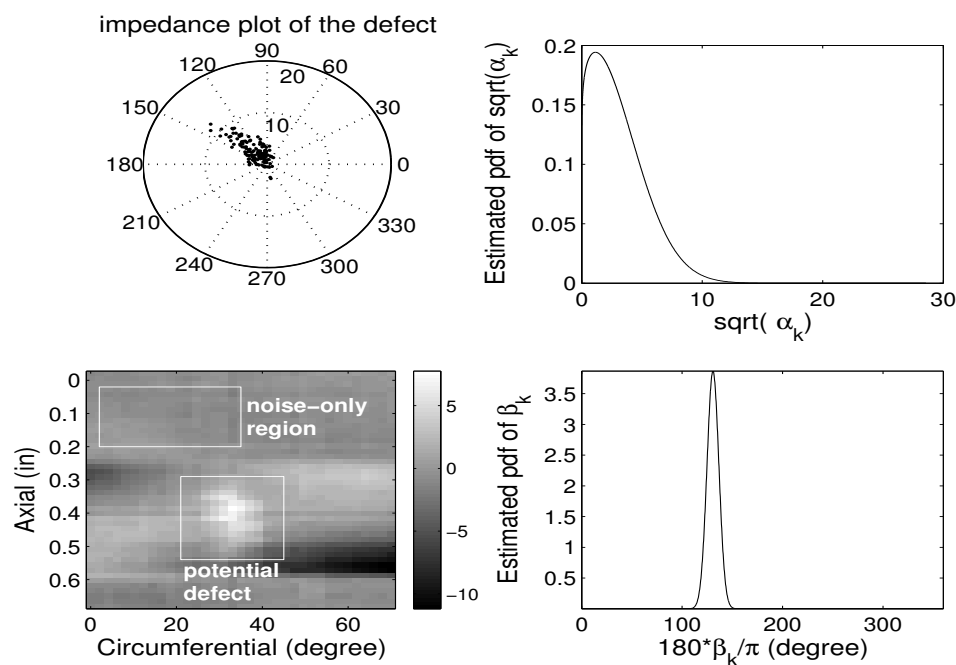
The above integrals with respect to  $\boldsymbol{\theta}$  can be easily computed using Gauss quadratures.

For more details, see

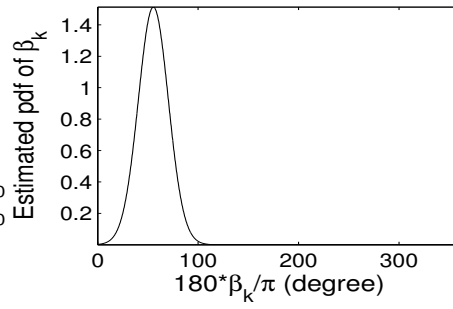
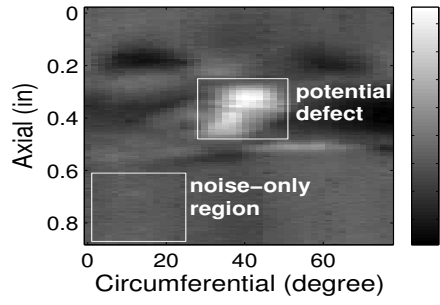
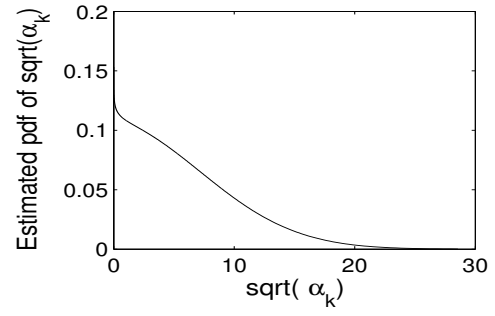
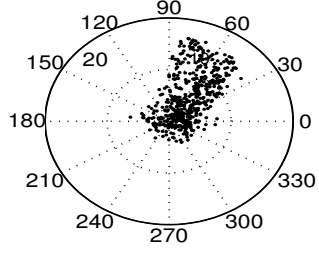
A. Dogandžić and P. Xiang, “A statistical model for eddy-current defect signals from steam generator tubes,” in *Rev. Progress Quantitative Nondestructive Evaluation*, D.O.

Thompson and D.E. Chimenti (Eds.), Melville NY: Amer. Inst. Phys., vol. 23, 2004, pp. 605–612.

We now show the performance of the proposed ML estimation method using data containing two real defects, obtained by inspecting steam-generator tubes. The  $K$  measurements  $x[k]$ ,  $k = 0, 1, \dots, K - 1$  were selected from potential defect regions, and the noise variance  $\sigma^2$  was estimated from the neighboring regions that contain only noise.



impedance plot of the defect



# Least-Squares Approach to Estimation

Suppose that we have a signal model

$$\mathbf{x} = \mathbf{H}\boldsymbol{\theta} + \mathbf{w}$$

where  $\mathbf{x} = [x[0], \dots, x[N-1]]^T$  is the vector of observations,  $\mathbf{H}$  is a known *regression vector matrix*, and  $\mathbf{w}$  is “error” vector.

LS problem formulation:

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \|\mathbf{x} - \mathbf{H}\boldsymbol{\theta}\|^2.$$

Solution:

$$\hat{\boldsymbol{\theta}} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{x}.$$

We can also use weighted least squares, which allows us to assign different weights to measurements. For example, if  $E[\mathbf{w}\mathbf{w}^T] = \mathbf{C}$  is known, we could use

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \|\mathbf{x} - \mathbf{H}\boldsymbol{\theta}\|_{\mathbf{C}^{-1}}^2 = \arg \min_{\boldsymbol{\theta}} (\mathbf{x} - \mathbf{H}\boldsymbol{\theta})^H \mathbf{C}^{-1} (\mathbf{x} - \mathbf{H}\boldsymbol{\theta})$$

Let  $\mathbf{H} = [\mathbf{h}_1 \cdots \mathbf{h}_p]$ :

$$\mathbf{x} = \sum_{k=1}^p \theta_k \mathbf{h}_k + \mathbf{w}.$$

The “signal part”  $\mathbf{H}\boldsymbol{\theta}$  of the  $N$ -vector  $\mathbf{x}$  is confined to the  $p$ -dimensional subspace spanned by  $[\mathbf{h}_1 \cdots \mathbf{h}_p]$  — the “signal subspace”!

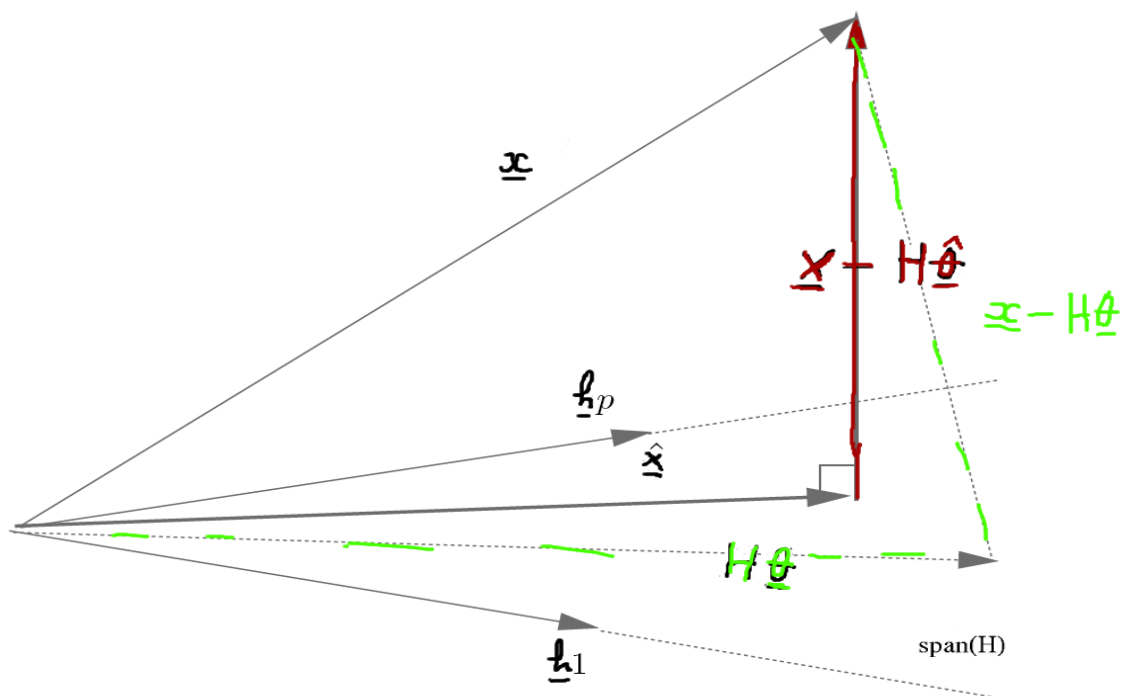
The signal estimate

$$\hat{\mathbf{x}} = \mathbf{H}\hat{\boldsymbol{\theta}} = \underbrace{\mathbf{H}(\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T}_{\mathbf{P}} \mathbf{x} = \mathbf{P}\mathbf{x}$$

is the *orthogonal projection* of  $\mathbf{x}$  onto  $\text{span}(\mathbf{H})$  (the column space of  $\mathbf{H}$ ). The error

$$\hat{\mathbf{w}} = \mathbf{x} - \mathbf{H}(\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{x}$$

is the *orthogonal projection* of  $\mathbf{x}$  onto the *orthogonal complement* of  $\text{span}(\mathbf{H})$ .



**Note:**  $P = H(H^T H)^{-1}H^T$  is the projection matrix onto the column space of  $H$ , and  $P^\perp = I - P$  is the complementary projection matrix.

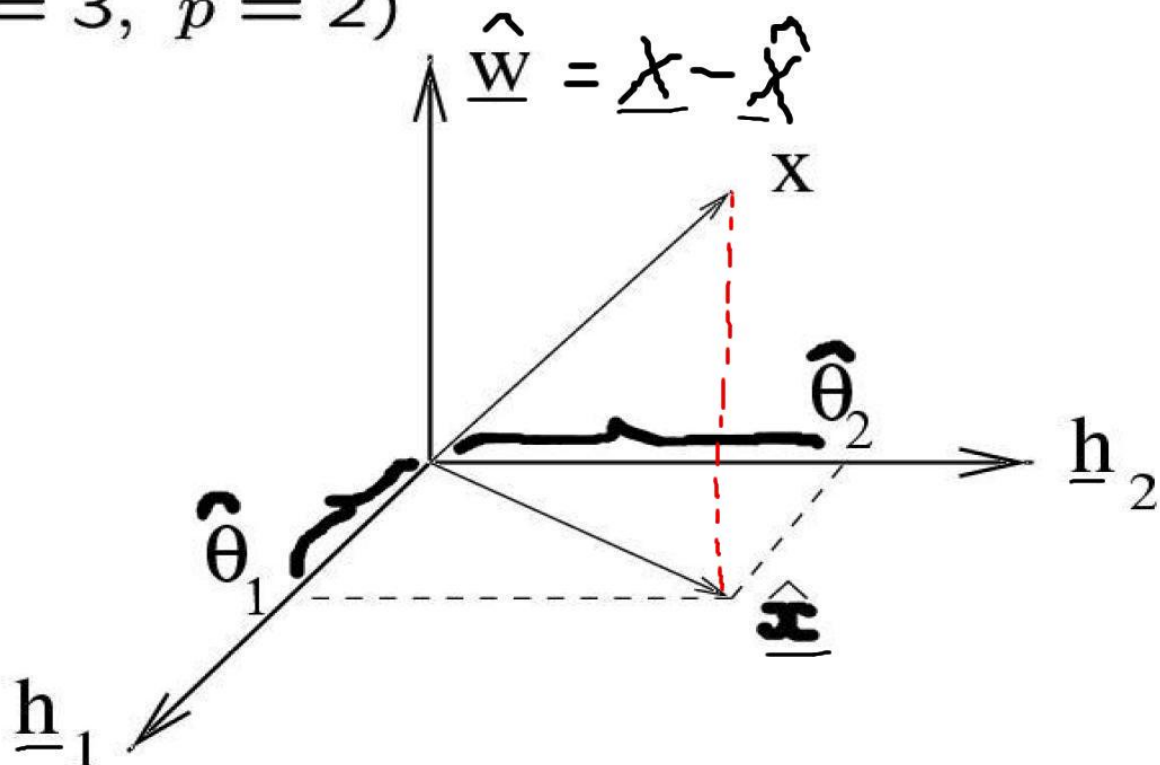
**Recall:** Projection matrices are symmetric and idempotent:  $P = P^T = P^2$ .

Also

$$(x - \underbrace{H\hat{\theta}}_{\hat{x}})^T H = 0.$$

**Example:**

$(N = 3, p = 2)$



Obviously,  $(x - \hat{x}) \perp \{h_1, h_2\}$ . This is the *orthogonality principle*: the minimum error is orthogonal to the columns of  $H$  (called regressors in statistics).

In general,

$$\mathbf{x} - \hat{\mathbf{x}} \perp \text{span}(\mathbf{H}) \Leftrightarrow \mathbf{x} - \hat{\mathbf{x}} \perp \mathbf{h}_j, \forall \mathbf{h}_j \Leftrightarrow \mathbf{H}^T (\mathbf{x} - \mathbf{H}\hat{\boldsymbol{\theta}}) = \mathbf{0}.$$



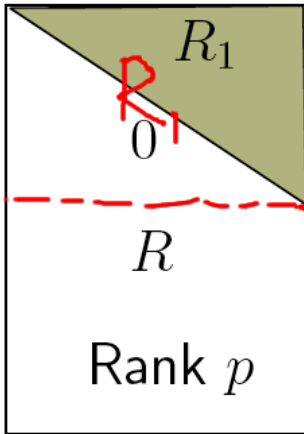
# Computational Aspects of Least Squares

QR decomposition of  $H$ :

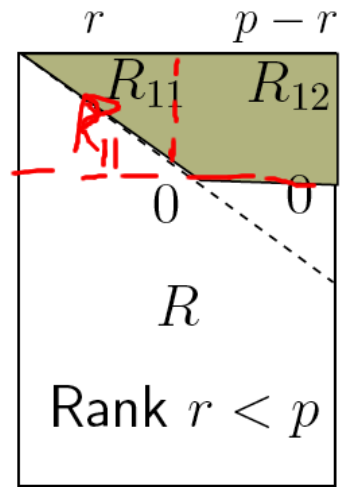
$$H_{N \times p} = Q_{N \times N} R_{N \times p}$$
$$= \begin{array}{|c|c|} \hline Q_1 & Q_2 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline R \\ \hline \end{array}$$

where  $Q$  has orthonormal columns:  $Q^T Q = I$  (and rows, i.e.  $Q Q^T = I$ ).

$R$  is upper triangular, and may not have full rank:



or



For the full-rank case,

$$\begin{aligned}\|x - H\theta\|^2 &= \|Q^T x - R\theta\|^2 \\ &= \|Q_1^T x - R_1 \theta\|^2 + \|Q_2^T x\|^2 \\ \implies \hat{\theta} &= R_1^{-1} Q_1^T x.\end{aligned}$$

### Comments:

- $Q^T x$  yields coordinates of  $x$  on columns of  $Q$ .
- $\hat{x} = Q_1 Q_1^T x = P x = H(H^T H)^{-1} H^T x$ . Here, the projection matrix  $P$  is also known as the *hat matrix* (because it puts the hat on  $x$ ).
- Non full rank case:  $\text{rank}(H) = r < p$ . We need to solve  $Q_1^T x = R_{11}\theta_1 + R_{12}\theta_2$ , where  $Q_1$  has  $r$  columns. There are infinitely many solutions — to get one, arbitrarily set  $\theta_2 = \mathbf{0}_{(p-r) \times 1}$  and solve for  $\theta_1$ . Here,  $\hat{x} = Q_1 Q_1^T x$  is still well defined, and unique.

# Nonlinear Least Squares (NLLS)

Often, the signal is *not* a linear function of  $\boldsymbol{\theta}$ , say  $\mathbf{f}(\boldsymbol{\theta})$ . Then, we obtain a NLLS estimate of  $\boldsymbol{\theta}$  as follows:

$$\begin{aligned}\hat{\boldsymbol{\theta}} &= \arg \min_{\boldsymbol{\theta}} V(\boldsymbol{\theta}) \\ V(\boldsymbol{\theta}) &= \|\mathbf{x} - \mathbf{f}(\boldsymbol{\theta})\|^2.\end{aligned}$$

**Example:**  $s[n] = r \cos(\omega n + \phi)$ ,  $n = 0, 1, 2, \dots, N - 1$  gives

$$\mathbf{f}(r, \omega, \phi) = \left[ r \cos(\phi), \dots, r \cos((N - 1)\omega + \phi) \right]^T.$$

Nonlinear problem  $\implies$  we usually need iterative optimization.

Recall the damped Newton-Raphson's method:

$$\hat{\boldsymbol{\theta}}^{(k+1)} = \hat{\boldsymbol{\theta}}^{(k)} - \mu_k \cdot \mathbf{H}_k^{-1} \mathbf{g}_k$$

where  $\mu_k$  is the step length and  $\mathbf{H}_k$ ,  $\mathbf{g}_k$  are the Hessian and gradient of  $V(\boldsymbol{\theta})$ , evaluated at  $\boldsymbol{\theta}^{(k)}$ .

# Nonlinear Least Squares

## Newton-Raphson Iteration

Define

$$\mathbf{f}_\theta^{(k)} = \left. \frac{\partial \mathbf{f}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^T} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(k)}}, \quad \mathbf{f}^{(k)} = \mathbf{f}(\boldsymbol{\theta}^{(k)}).$$

The partial derivatives are then

$$\mathbf{g}_k = \left. \frac{\partial (\mathbf{x} - \mathbf{f})^T (\mathbf{x} - \mathbf{f})}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(k)}} = -2(\mathbf{f}_\theta^{(k)})^T (\mathbf{x} - \mathbf{f}^{(k)})$$
$$\mathbf{H}_k = \left. \frac{\partial^2 (\mathbf{x} - \mathbf{f})^T (\mathbf{x} - \mathbf{f})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(k)}} = 2(\mathbf{f}_\theta^{(k)})^T \mathbf{f}_\theta^{(k)} - 2\mathbf{G}^{(k)}$$

where  $[\mathbf{G}^{(k)}]_{i,l} = \left. \frac{\partial^2 \mathbf{f}^T}{\partial \theta_i \partial \theta_l} (\mathbf{x} - \mathbf{f}) \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(k)}}$ .

Assuming that we have a “small residual” problem, such that  $\mathbf{x} - \mathbf{f}^{(k)} \approx \mathbf{0}$  (close to the optimum), the Hessian is approximated by

$$\mathbf{H}_k = 2(\mathbf{f}_\theta^{(k)})^T \mathbf{f}_\theta^{(k)}.$$

Recall:  $(\mathbf{f}_\theta^{(k)})^T \mathbf{f}_\theta^{(k)}$  is the FIM for  $\boldsymbol{\theta}$  (under the AWGN measurement model), hence this approach is equivalent to Fisher scoring when the noise is AWGN. It is also known as the *Gauss-Newton algorithm*.

# Nonlinear Least Squares (cont.)

## (Damped) Gauss-Newton:

$$\hat{\boldsymbol{\theta}}^{(k+1)} = \hat{\boldsymbol{\theta}}^{(k)} + \mu_k [(\mathbf{f}_{\boldsymbol{\theta}}^{(k)})^T (\mathbf{f}_{\boldsymbol{\theta}}^{(k)})]^{-1} (\mathbf{f}_{\boldsymbol{\theta}}^{(k)})^T (\mathbf{x} - \mathbf{f}^{(k)}).$$

The search direction  $\boldsymbol{\gamma} = [(\mathbf{f}_{\boldsymbol{\theta}}^{(k)})^T (\mathbf{f}_{\boldsymbol{\theta}}^{(k)})]^{-1} (\mathbf{f}_{\boldsymbol{\theta}}^{(k)})^T (\mathbf{x} - \mathbf{f}^{(k)})$  is the LS solution to

$$\min_{\boldsymbol{\gamma}} \|(\mathbf{x} - \mathbf{f}^{(k)}) - \mathbf{f}_{\boldsymbol{\theta}}^{(k)} \boldsymbol{\gamma}^{(k)}\|^2$$

which is efficiently computed in MATLAB using

$$\boldsymbol{\gamma} = \mathbf{f}_{\boldsymbol{\theta}} \backslash (\mathbf{x} - \mathbf{f}).$$

Note that the approximate Hessian  $\mathbf{f}_{\boldsymbol{\theta}}^T \mathbf{f}_{\boldsymbol{\theta}}$  is always positive (semi)definite, which is generally *not true for the exact Hessian!*

# Separable NLLS

Consider the sinusoid example

$$s[n] = r \cos(\omega n + \phi) = A \sin(\omega n) + B \cos(\omega n).$$

$A$  and  $B$  enter linearly in  $s[n]$ ! We can write

$$\mathbf{f}(\boldsymbol{\theta}) = \mathbf{H}(\underbrace{\boldsymbol{\alpha}}_{\omega}) \underbrace{\begin{bmatrix} \beta \\ A \\ B \end{bmatrix}}_{\boldsymbol{\beta}}$$

where  $\boldsymbol{\theta} = \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix}$ . For a fixed  $\boldsymbol{\alpha}$ , the LS solution for  $\boldsymbol{\beta}$  is

$$\hat{\boldsymbol{\beta}}(\boldsymbol{\alpha}) = [\mathbf{H}^T(\boldsymbol{\alpha})\mathbf{H}(\boldsymbol{\alpha})]^{-1}\mathbf{H}^T(\boldsymbol{\alpha})\mathbf{x}.$$

Substituting into  $V(\boldsymbol{\theta})$  gives the concentrated criterion:

$$V_c(\boldsymbol{\alpha}) = \left\| \mathbf{x} - \underbrace{\mathbf{H}(\boldsymbol{\alpha})[\mathbf{H}^T(\boldsymbol{\alpha})\mathbf{H}(\boldsymbol{\alpha})]^{-1}\mathbf{H}^T(\boldsymbol{\alpha})}_{\mathbf{P}(\boldsymbol{\alpha})} \mathbf{x} \right\|^2$$

where  $\mathbf{P}(\boldsymbol{\alpha})$  is the projection matrix onto the column space of  $\mathbf{H}(\boldsymbol{\alpha})$ . Equivalently

$$\hat{\boldsymbol{\alpha}} = \arg \max_{\boldsymbol{\alpha}} \mathbf{x}^T \underbrace{\mathbf{H}(\boldsymbol{\alpha})[\mathbf{H}^T(\boldsymbol{\alpha})\mathbf{H}(\boldsymbol{\alpha})]^{-1}\mathbf{H}^T(\boldsymbol{\alpha})}_{\mathbf{P}(\boldsymbol{\alpha})} \mathbf{x}.$$

Here,  $\hat{\alpha}$  maximizes the projection of  $x$  onto the signal subspace.

We have used the fact that our cost function can be easily minimized with respect to a subset of parameters ( $\beta$ , in our case) if the rest of the parameters  $\alpha$  are fixed. We have obtained a *concentrated* cost function to be maximized with respect to  $\alpha$  only.

### Comments:

- There is nothing fundamentally statistical about LS: the least squares approach solves a minimization problem in vector spaces.
- In linear problems, LS allows a closed-form solution.
- We need to replace  $T$  with  $H$  (the Hermitian transpose) to obtain the corresponding results for complex data:

$$\hat{\beta}(\alpha) = [\mathbf{H}^H(\alpha)\mathbf{H}(\alpha)]^{-1}\mathbf{H}^H(\alpha)x$$

minimizes  $\|x - \mathbf{H}(\alpha)\beta\|^2 = [x - \mathbf{H}(\alpha)\beta]^H [x - \mathbf{H}(\alpha)\beta]$ ,  
i.e.

$$\hat{\beta}(\alpha) = \arg \min_{\beta} \|x - \mathbf{H}(\alpha)\beta\|^2$$

and  $\alpha$  can be estimated using the concentrated criterion:

$$\hat{\alpha} = \arg \max_{\alpha} x^H \underbrace{\mathbf{H}(\alpha)[\mathbf{H}^H(\alpha)\mathbf{H}(\alpha)]^{-1}\mathbf{H}^H(\alpha)}_{P(\alpha)} x.$$