

Least Squares Estimation, Filtering, and Prediction

- Principle of least squares
- Normal equations
- Weighted least squares
- Statistical properties
- FIR filters
- Windowing
- Combined forward-backward linear prediction
- Narrowband Interference Cancellation

Motivation

- If the second-order statistics are known, the optimum estimator is given by the normal equations
- In many applications, they aren't known
- Alternative approach is to estimate the coefficients from observed data
- Two possible approaches
 - Estimate required moments from available data and build an *approximate* MMSE estimator
 - Build an estimator that minimizes some error functional calculated from the available data

MMSE versus Least Squares

- Recall that MMSE estimators are optimal in expectation across the ensemble of all stochastic processes with the same second order statistics
- Least squares estimators minimize the error on a given **block** of data
 - In signal processing applications, the block of data is a finite-length period of time
- No guarantees about optimality on other data sets or other stochastic processes
- If the process is ergodic and stationary, the LSE estimator approaches the MMSE estimator as the size of the data set grows
 - This is the first time in this class we've discussed estimation from data
 - First time we need to consider ergodicity

Principle of Least Squares

- Will only discuss the sum of squares as the performance criterion
 - Recall our earlier discussion about alternatives
 - Essentially, picking sum of squares will permit us to obtain a closed-form optimal solution
- Requires a data set where both the inputs and desired responses are known
- Recall the range of possible applications
 - Plant modeling for control (system identification)
 - Inverse modeling/deconvolution
 - Interference cancellation
 - Prediction

Recalling the Book's Notation

- $y(n) \in \mathbb{C}^{1 \times 1}$ is the **target** or **desired response**
- $x_k(n)$ represent the inputs
- These may be of several types
 - Multiple sensors, no lags: $\mathbf{x}(n) = [x_1(n), x_2(n), \dots, x_M(n)]^T$
 - Lag window: $\mathbf{x}(n) = [x(n), x(n-1), \dots, x(n-M+1)]^T$
 - Combined
- Data sets consists of values over the time span $0 \leq n \leq N-1$
- Boldface is now used for vectors and matrices
- The coefficients are represented as $\mathbf{c}(n)$

Change in Notation

In a trade of elegance and simplicity for inconsistency, I'm going to break with some of the book's notational conventions

$$\text{Notes: } \hat{y}(n) \triangleq \sum_{k=1}^M c_k(n)x_k(n) = \mathbf{c}^T(n)\mathbf{x}(n)$$

$$\text{Book: } \hat{y}(n) \triangleq \sum_{k=1}^M c_k^*(n)x_k(n) = \mathbf{c}^H(n)\mathbf{x}(n)$$

- In the case that \mathbf{c} is real, they are consistent
- Rationale
 - The inner product notation leads to unnecessary complications in the notation
 - Most books use the same notation that the notes use
 - Leads to a symmetry: $\mathbf{c}^T(n)\mathbf{x}(n) = \mathbf{x}^T(n)\mathbf{c}(n)$

Definitions

$$\text{Estimate: } \hat{y}(n) \triangleq \sum_{k=1}^M c_k(n)x_k(n) = \mathbf{c}^T(n)\mathbf{x}(n)$$

$$\text{Estimation error: } e(n) \triangleq y(n) - \hat{y}(n) = y(n) - \mathbf{c}(n)\mathbf{x}(n)$$

$$\text{Sum of squared errors: } E_e \triangleq \sum_{n=0}^{N-1} |e(n)|^2$$

- Coefficient vector $\mathbf{c}(n)$ is typically held constant over the data window, $0 \leq n \leq N-1$
- Contrast with adaptive filter approach
- The coefficients \mathbf{c} that minimize E_e are called the **linear LSE estimator**

Matrix Formulation

$$\begin{bmatrix} e(0) \\ e(1) \\ \vdots \\ e(N-1) \end{bmatrix} = \begin{bmatrix} y(0) \\ y(1) \\ \vdots \\ y(N-1) \end{bmatrix} - \begin{bmatrix} x_1(0) & x_2(0) & \dots & x_M(0) \\ x_1(1) & x_2(1) & \vdots & x_M(1) \\ \vdots & \vdots & \ddots & \vdots \\ x_1(N-1) & x_2(N-1) & \dots & x_M(N-1) \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_M \end{bmatrix}$$

$$\mathbf{e} = \mathbf{y} - \mathbf{X}\mathbf{c}$$

where

$$\mathbf{e} \triangleq [e(0) \ e(1) \ \dots \ e(N-1)]^T \quad \text{error data vector } (N \times 1)$$

$$\mathbf{y} \triangleq [y(0) \ y(1) \ \dots \ y(N-1)]^T \quad \text{desired response vector } (N \times 1)$$

$$\mathbf{X} \triangleq [\mathbf{x}^T(0) \ \mathbf{x}^T(1) \ \dots \ \mathbf{x}^T(N-1)]^T \quad \text{input data matrix } (N \times M)$$

$$\mathbf{c} \triangleq [c_1 \ c_2 \ \dots \ c_M]^T \quad \text{combiner parameter vector } (M \times 1)$$

Note: my definitions differ from the book by a conjugate factor *

Input Data Matrix Notation

$$\mathbf{X} = \begin{bmatrix} x_1(0) & x_2(0) & \dots & x_M(0) \\ x_1(1) & x_2(1) & \vdots & x_M(1) \\ \vdots & \vdots & \ddots & \vdots \\ x_1(N-1) & x_2(N-1) & \dots & x_M(N-1) \end{bmatrix} = \begin{bmatrix} \mathbf{x}^T(0) \\ \mathbf{x}^T(1) \\ \vdots \\ \mathbf{x}^T(N-1) \end{bmatrix}$$

$$= [\vec{x}_1 \quad \vec{x}_2 \quad \dots \quad \vec{x}_M]$$

- We will need to reference both the row and column vectors of the data matrix \mathbf{X}
- This is always awkward to denote
- Book's notation is redundant
 - Row vectors ("snapshots") indicated by (n) and boldface \mathbf{x}
 - Column vectors ("data records") indicated by k and \vec{x} (book uses \tilde{x})
- I don't know of a more elegant solution

Size of Data Matrix

$$\underset{N \times 1}{\mathbf{e}} = \underset{N \times 1}{\mathbf{y}} - \underset{N \times M}{\mathbf{X}} \underset{M \times 1}{\mathbf{c}}$$

- Suppose we wish to make $|e| = 0$

$$\mathbf{y} = \mathbf{X}\mathbf{c}$$
- Suppose \mathbf{X} has maximum rank
 - Linearly independent rows or columns
 - Always occurs with real data
- N linear equations and M unknowns
 - $N < M$: underdetermined, infinite number of solutions
 - $N = M$: unique solution, $\mathbf{c} = \mathbf{X}^{-1}\mathbf{y}$
 - $N > M$: overdetermined, no solution, in general
- In practical applications we pick $N > M$ (why?)

Block Processing

- LSE estimators *can* be used in **block processing mode**
 - Take a segment of N input-output observations, say $n_1 \leq n \leq n_1 + N - 1$
 - Estimate the coefficients
 - Increment the temporal location of the block to $n_1 + N_0$
- The blocks overlap by $N - N_0$ samples
- Reminiscent of Welch's method of PSD estimation
- Useful for parametric time-frequency analysis
- In most other nonstationary applications, adaptive filters are usually used instead

Normal Equations

$$\begin{aligned} E_e &= \|\mathbf{e}\|^2 = \mathbf{e}^H \mathbf{e} \\ &= (\mathbf{y} - \mathbf{X}\mathbf{c})^H (\mathbf{y} - \mathbf{X}\mathbf{c}) \\ &= (\mathbf{y}^H - \mathbf{c}^H \mathbf{X}^H) (\mathbf{y} - \mathbf{X}\mathbf{c}) \\ &= \mathbf{y}^H \mathbf{y} - \mathbf{c}^H \mathbf{X}^H \mathbf{y} - \mathbf{y}^H \mathbf{X} \mathbf{c} + \mathbf{c}^H \mathbf{X}^H \mathbf{X} \mathbf{c} \end{aligned}$$

- E_e is a nonlinear function of \mathbf{y} , \mathbf{X} , and \mathbf{c}
- Is a quadratic function of each of these components
- E_e is the sum of squared errors
 - Energy of the error signal over the interval $0 \leq n \leq N - 1$
- If we take the average squared errors (ASE), we have an estimate of the mean square error = the estimated power of the error

$$\hat{P}_e = \frac{1}{N} E_e = \frac{1}{N} \sum_{n=0}^{N-1} |e(n)|^2$$

Normal Equation Components

Let us define

$$\text{Average squared error (ASE): } \hat{P}_e \triangleq \frac{1}{N} \|e\|^2 = \frac{1}{N} \sum_{n=0}^{N-1} |e(n)|^2$$

$$\text{Average signal power: } \hat{P}_y \triangleq \frac{1}{N} \|y\|^2 = \frac{1}{N} \sum_{n=0}^{N-1} |y(n)|^2$$

$$\text{Average correlation: } \hat{\mathbf{R}} \triangleq \frac{1}{N} \mathbf{X}^H \mathbf{X} = \frac{1}{N} \sum_{n=0}^{N-1} \mathbf{x}(n) \mathbf{x}^T(n)$$

$$\text{Average cross-correlation: } \hat{\mathbf{d}} \triangleq \frac{1}{N} \mathbf{X}^H \mathbf{y} = \frac{1}{N} \sum_{n=0}^{N-1} \mathbf{x}(n) y(n)$$

If we replaced the sample average $\frac{1}{N} \sum_{n=0}^{N-1} (\cdot)$ with expectation, $E[\cdot]$, each of these terms would be the **mean ensemble value** rather than the **sample average estimate**

Relating LSE and MMSE Estimation

Then

$$\begin{aligned} \hat{P}_e &= \frac{1}{N} (\mathbf{y}^H \mathbf{y} - \mathbf{c}^H \mathbf{X}^H \mathbf{y} - \mathbf{y}^H \mathbf{X} \mathbf{c} + \mathbf{c}^H \mathbf{X}^H \mathbf{X} \mathbf{c}) \\ &= \hat{P}_y - \mathbf{c}^H \hat{\mathbf{d}} - \hat{\mathbf{d}}^H \mathbf{c} + \mathbf{c}^H \hat{\mathbf{R}} \mathbf{c} \end{aligned}$$

This should look familiar. . .

If we assume $\hat{\mathbf{R}} > 0$, we can complete the square (fourth time!),

$$\begin{aligned} \hat{P}_e(\mathbf{c}) &= \hat{P}_y + \mathbf{c}^H \hat{\mathbf{R}} (\mathbf{c} - \hat{\mathbf{R}}^{-1} \hat{\mathbf{d}}) - \hat{\mathbf{d}}^H \mathbf{c} + \hat{\mathbf{d}}^H \hat{\mathbf{R}}^{-1} \hat{\mathbf{d}} - \hat{\mathbf{d}}^H \hat{\mathbf{R}}^{-1} \hat{\mathbf{d}} \\ &= \hat{P}_y + (\mathbf{c}^H) \hat{\mathbf{R}} (\mathbf{c} - \hat{\mathbf{R}}^{-1} \hat{\mathbf{d}}) - (\hat{\mathbf{d}}^H \hat{\mathbf{R}}^{-1}) \hat{\mathbf{R}} (\mathbf{c} - \hat{\mathbf{R}}^{-1} \hat{\mathbf{d}}) - \hat{\mathbf{d}}^H \hat{\mathbf{R}}^{-1} \hat{\mathbf{d}} \\ &= \hat{P}_y + (\mathbf{c}^H - \hat{\mathbf{d}}^H \hat{\mathbf{R}}^{-1}) \hat{\mathbf{R}} (\mathbf{c} - \hat{\mathbf{R}}^{-1} \hat{\mathbf{d}}) - \hat{\mathbf{d}}^H \hat{\mathbf{R}}^{-1} \hat{\mathbf{d}} \\ &= \hat{P}_y - \hat{\mathbf{d}}^H \hat{\mathbf{R}}^{-1} \hat{\mathbf{d}} + (\mathbf{c} - \hat{\mathbf{R}}^{-1} \hat{\mathbf{d}})^H \hat{\mathbf{R}} (\mathbf{c} - \hat{\mathbf{R}}^{-1} \hat{\mathbf{d}}) \\ &= \hat{P}_y - \hat{\mathbf{d}}^H \hat{\mathbf{R}}^{-1} \hat{\mathbf{d}} + (\hat{\mathbf{R}} \mathbf{c} - \hat{\mathbf{d}})^H \hat{\mathbf{R}}^{-1} (\hat{\mathbf{R}} \mathbf{c} - \hat{\mathbf{d}}) \\ &= \hat{P}_{ls} + (\hat{\mathbf{R}} \mathbf{c} - \hat{\mathbf{d}})^H \hat{\mathbf{R}}^{-1} (\hat{\mathbf{R}} \mathbf{c} - \hat{\mathbf{d}}) \end{aligned}$$

Least Squares Estimate

$$\hat{P}_e(\mathbf{c}) = \hat{P}_y - \hat{\mathbf{d}}^H \hat{\mathbf{R}}^{-1} \hat{\mathbf{d}} + (\hat{\mathbf{R}} \mathbf{c} - \hat{\mathbf{d}})^H \hat{\mathbf{R}}^{-1} (\hat{\mathbf{R}} \mathbf{c} - \hat{\mathbf{d}})$$

So the least squares estimate and minimum average squared error are given by

$$\begin{aligned} \mathbf{c}_{ls} &= \hat{\mathbf{R}}^{-1} \hat{\mathbf{d}} \\ &= (\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}^H \mathbf{y} \\ \hat{P}_{ls} &\triangleq \hat{P}_y - \hat{\mathbf{d}}^H \hat{\mathbf{R}}^{-1} \hat{\mathbf{d}} \\ &= \hat{P}_y - \hat{\mathbf{d}}^H \mathbf{c}_{ls} \end{aligned}$$

- Both the LSE and MSE criteria are quadratic functions of the coefficient vector \mathbf{c}

Computational Issues

$$\begin{aligned} \mathcal{O}(M^2 N) & \quad \hat{\mathbf{R}} \triangleq \frac{1}{N} \mathbf{X}^H \mathbf{X} \\ \mathcal{O}(MN) & \quad \hat{\mathbf{d}} \triangleq \frac{1}{N} \mathbf{X}^H \mathbf{y} \\ \mathcal{O}(M^3) & \quad \mathbf{c}_{ls} = \hat{\mathbf{R}}^{-1} \hat{\mathbf{d}} \end{aligned}$$

- Typically $M \ll N$
- The most computationally expensive operation is calculating the input correlation matrix!
- If $\mathbf{x}(n)$ comes from a stationary time series, we can speed this up with the FFT
 - Recall the estimate from ECE 538/638
 - More later

Geometric Derivation

$$\mathbf{e} = \mathbf{y} - \mathbf{X}\mathbf{c}$$

- I was disparaging of the geometric interpretation for the MSE case
- It makes a lot more sense for the LSE case
- We are trying to estimate an N dimensional vector as a linear combination of the M columns of \mathbf{X}
- The orthogonal projection of \mathbf{y} onto the M dimensional subspace minimizes the error

Inner Product Definition

- In this case an inner product can be defined as

$$\langle \vec{\mathbf{x}}_i, \vec{\mathbf{x}}_j \rangle \triangleq \vec{\mathbf{x}}_i^H \vec{\mathbf{x}}_j = \sum_{n=0}^{N-1} x_i(n) x_j^*(n)$$
$$\|\vec{\mathbf{x}}_i\|^2 \triangleq \vec{\mathbf{x}}_i^H \vec{\mathbf{x}}_i = \sum_{n=0}^{N-1} |x_i(n)|^2$$

- Can easily show that this has all of the required properties of an inner product

Projection Theorem: Estimator

- The projection theorem holds
 - Therefore the projection of \mathbf{y} onto the column space of \mathbf{X} minimizes the length of the residual vector

$$\|\mathbf{e}\|^2 = \sum_{n=0}^{N-1} |y(n) - \mathbf{c}^T \mathbf{x}(n)|^2$$

- The error vector is also orthogonal to the observations

$$\langle \mathbf{X}, \mathbf{y} - \mathbf{X}\mathbf{c}_{ls} \rangle = \mathbf{X}^H \mathbf{y} - \mathbf{X}^H \mathbf{X} \mathbf{c}_{ls} = N\hat{\mathbf{d}} - N\hat{\mathbf{R}}\mathbf{c}_{ls} = 0$$

- Thus we have another way of obtaining the normal equations

$$\hat{\mathbf{R}}\mathbf{c}_{ls} = \hat{\mathbf{d}}$$

Projection Theorem: LSE

- Further, by orthogonality we have

$$\|\mathbf{y}\|^2 = \|\mathbf{e}\|^2 + \|\hat{\mathbf{y}}\|^2$$
$$\|\mathbf{e}\|^2 = \|\mathbf{y}\|^2 - \|\hat{\mathbf{y}}\|^2$$
$$P_{ls} = P_y - \mathbf{c}_{ls}^H \mathbf{X}^H \mathbf{X} \mathbf{c}_{ls}$$
$$= P_y - \hat{\mathbf{d}}^H \hat{\mathbf{R}}^{-1} \hat{\mathbf{R}} \mathbf{c}_{ls}$$
$$= P_y - \hat{\mathbf{d}}^H \mathbf{c}_{ls}$$

Uniqueness

- As with the MMSE case we have
 - \mathbf{c}_{ls} is unique if \mathbf{X} has full column rank and $M \leq N$
 - Otherwise there are infinite, equivalent solutions with the same LSE
- Regardless the LSE estimate $\hat{\mathbf{y}}(n) = \mathbf{c}^T(n)\mathbf{x}(n)$ is the same
- Full column rank is equivalent to the requirement that $\hat{\mathbf{R}}$ be positive definite

Weighted Least Squares

- Suppose some errors are more significant than others
- A closely related problem is weighted least squares

$$\begin{aligned} E_e &= \sum_{n=0}^{N-1} w_n^2 |y(n) - \mathbf{c}^T \mathbf{x}(n)|^2 \\ &= (\mathbf{y} - \mathbf{X}\mathbf{c})^H \mathbf{W}^2 (\mathbf{y} - \mathbf{X}\mathbf{c}) \\ &= (\mathbf{W}\mathbf{y} - \mathbf{W}\mathbf{X}\mathbf{c})^H (\mathbf{W}\mathbf{y} - \mathbf{W}\mathbf{X}\mathbf{c}) \end{aligned}$$

If we define

$$\hat{\mathbf{y}} \triangleq \mathbf{W}\mathbf{y} \qquad \hat{\mathbf{X}} \triangleq \mathbf{W}\mathbf{X}$$

we have the original LSE problem

$$E_e = \sum_{n=0}^{N-1} |\hat{\mathbf{y}}(n) - \mathbf{c}^T \hat{\mathbf{x}}(n)|^2 = (\hat{\mathbf{y}} - \hat{\mathbf{X}}\mathbf{c})^H (\hat{\mathbf{y}} - \hat{\mathbf{X}}\mathbf{c})$$

Weighted Least Squares Comments

$$E_e = (\mathbf{y} - \mathbf{X}\mathbf{c})^H \mathbf{W}^2 (\mathbf{y} - \mathbf{X}\mathbf{c})$$

- Weighted least squares can be applied with any weighting matrix \mathbf{W}^2 that is positive definite
 - \mathbf{W}^2 does not have to be diagonal
 - All positive definite matrices have a square root $\mathbf{W}^2 = \mathbf{W}^H \mathbf{W}$ for some matrix \mathbf{W}
 - The square root is not unique
- Useful for
 - Windowing parametric estimators
 - Accounting for correlated observation noise

Statistical Models

$$\mathbf{y} = \mathbf{X}\mathbf{c}_o + \mathbf{v}$$

$$\begin{aligned} \mathbf{c}_{ls} &= (\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}^H \mathbf{y} \\ &= (\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}^H \mathbf{X} \mathbf{c}_o + (\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}^H \mathbf{v} \\ &= \mathbf{c}_o + (\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}^H \mathbf{v} \end{aligned}$$

- Most of the statistical properties require a **statistical model** of how the data was generated
- This idea is similar to the state space model used by the Kalman filter
- The linear statistical model used here is more general (no assumption of state dynamics)
- Some of the properties don't hold when the model is not accurate

Linear Statistical Model

$$\underset{N \times 1}{\mathbf{y}} = \underset{N \times M}{\mathbf{X}} \underset{M \times 1}{\mathbf{c}_o} + \underset{N \times 1}{\mathbf{v}}$$

- \mathbf{y} is a vector of observed measurements
- \mathbf{X} is well conditioned
 - Deterministic: \mathbf{X} has full column rank
 - Stochastic: $E[\mathbf{X}^H \mathbf{X}]$ is positive definite
- \mathbf{c}_o is considered the “true” parameter vector
 - The notational overlap with MSE estimation notes is intentional
- \mathbf{v} is a vector of random noise or “errors”
 - Note my notation differs from the text, which used e_o
 - All of the elements of \mathbf{X} are independent with \mathbf{v}
 - \mathbf{v} has zero mean: $E[\mathbf{v}] = 0$
 - The elements of \mathbf{v} are uncorrelated: $E[\mathbf{v}\mathbf{v}^H] = \sigma_v^2 \mathbf{I}$

Deterministic versus Stochastic Data Matrix

$$\mathbf{y} = \mathbf{X}\mathbf{c}_o + \mathbf{v}$$

- The **data matrix** \mathbf{X} may be deterministic or stochastic
- Appropriate model depends on the application
- Deterministic
 - Appropriate for conditions where the entries in \mathbf{X} are controlled by the user
 - Only source of randomness is then \mathbf{v}
 - Simplifies the analysis of the estimator
- Stochastic
 - Appropriate for conditions where \mathbf{X} are observed and randomly fluctuating
 - Most signal processing application

Deterministic Case: Estimator Properties

$$\mathbf{y} = \mathbf{X}\mathbf{c}_o + \mathbf{v} \quad \mathbf{c}_{ls} = \mathbf{c}_o + (\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}^H \mathbf{v}$$

- The estimator \mathbf{c}_{ls} is unbiased

$$E[\mathbf{c}_{ls}] = \mathbf{c}_o$$

- The estimator covariance matrix is

$$\Lambda_{ls} \triangleq E[(\mathbf{c}_{ls} - \mathbf{c}_o)(\mathbf{c}_{ls} - \mathbf{c}_o)^H] = \sigma_v^2 (\mathbf{X}^H \mathbf{X})^{-1} = \frac{\sigma_v^2}{N} \hat{\mathbf{R}}^{-1}$$

- The diagonal elements of Λ_{ls} contain the variance of the elements of \mathbf{c}_{ls}
- If \mathbf{v} is Gaussian, we can construct exact confidence intervals!
- Difference in definition of $\hat{\mathbf{R}}$ make the dependence on N explicit
- Covariance is inversely proportional to the number of observations

Deterministic Case: Residuals

$$\mathbf{P} \triangleq \mathbf{X}(\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}^H$$

$$\mathbf{c}_{ls} = \mathbf{c}_o + (\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}^H \mathbf{v}$$

$$\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}}$$

$$= \mathbf{y} - \mathbf{X}\mathbf{c}_{ls}$$

$$= \mathbf{y} - \mathbf{X}(\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}^H \mathbf{y}$$

$$= \mathbf{X}\mathbf{c}_o + \mathbf{v} - \mathbf{X}(\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}(\mathbf{X}\mathbf{c}_o + \mathbf{v})$$

$$= \mathbf{v} - \mathbf{P}\mathbf{v}$$

$$= (\mathbf{I} - \mathbf{P})\mathbf{v}$$

Deterministic Case: Error Variance

$$\begin{aligned}
 E_e &= \mathbf{e}^H \mathbf{e} \\
 &= \mathbf{v}^H (\mathbf{I} - \mathbf{P})^H (\mathbf{I} - \mathbf{P}) \mathbf{v} \\
 &= \mathbf{v}^H (\mathbf{I} - \mathbf{P})^H (\mathbf{I} - \mathbf{P}) \mathbf{v} \\
 &= \mathbf{v}^H (\mathbf{I} - \mathbf{P}) \mathbf{v} - \mathbf{v}^H (\mathbf{P} - \mathbf{P}\mathbf{P}) \mathbf{v} \\
 &= \mathbf{v}^H (\mathbf{I} - \mathbf{P}) \mathbf{v} \\
 &= \text{trace}[\mathbf{v}^H (\mathbf{I} - \mathbf{P}) \mathbf{v}] \\
 &= \text{trace}[(\mathbf{I} - \mathbf{P}) \mathbf{v} \mathbf{v}^H] \\
 E[E_e] &= \text{trace}[(\mathbf{I} - \mathbf{P}) \sigma_v^2 \mathbf{I}] \\
 &= \sigma_v^2 \text{trace}[\mathbf{I} - \mathbf{P}]
 \end{aligned}$$

Properties of the trace[.] operator

$$\text{trace}[\mathbf{A}\mathbf{B}] = \text{trace}[\mathbf{B}\mathbf{A}] \quad \text{trace}[\mathbf{A} + \mathbf{B}] = \text{trace}[\mathbf{A}] + \text{trace}[\mathbf{B}]$$

Deterministic Case: Error Variance

$$\begin{aligned}
 E[E_e] &= \sigma_v^2 \text{trace}[\mathbf{I} - \mathbf{P}] \\
 \text{trace}[\mathbf{I} - \mathbf{P}] &= \text{trace}[\mathbf{I} - \mathbf{X}(\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}^H] \\
 &= \text{trace}[\mathbf{I}] - \text{trace}[(\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}^H \mathbf{X}] \\
 &= N - \text{trace}[\mathbf{I}] \\
 &= N - M \\
 E[E_e] &= \sigma_v^2 (N - M) \\
 \hat{\sigma}_v^2 &\triangleq \frac{1}{N - M} E_e \\
 E[\hat{\sigma}_v^2] &= \frac{1}{N - M} E[E_e] \\
 &= E_e
 \end{aligned}$$

- Therefore $\hat{\sigma}_v^2$ is unbiased
- The difference $N - M$ is often called the **degrees of freedom**

Deterministic Case: Other, Unproved Properties

$$\mathbf{y} = \mathbf{X}\mathbf{c}_o + \mathbf{v} \quad \mathbf{c}_{ls} = \mathbf{c}_o + (\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}^H \mathbf{v}$$

- The weighted LSE estimate is also unbiased
- If the error covariance is not diagonal, $R_v \neq \sigma_v^2 \mathbf{I}$, the optimal estimator is obtained by setting $\mathbf{W}^2 = \mathbf{R}_v^{-1}$
- In both cases, the LSE estimator is the **best linear unbiased estimator** (BLUE)
 - Of all the unbiased estimators, this one has the smallest variance
- If \mathbf{v} is normally distributed,, the LSE estimator is also the maximum likelihood estimator

Stochastic Case: Properties

$$\mathbf{y} = \mathbf{X}\mathbf{c}_o + \mathbf{v} \quad \mathbf{c}_{ls} = \mathbf{c}_o + (\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}^H \mathbf{v}$$

- Model assumptions
 - \mathbf{X} and \mathbf{v} are statistically independent
 - Merely uncorrelated is insufficient in this case
- The estimator is still unbiased in this case
- The covariance of \mathbf{c}_{ls} is given by

$$\mathbf{\Lambda}_{ls} \triangleq E[(\mathbf{c}_{ls} - \mathbf{c}_o)(\mathbf{c}_{ls} - \mathbf{c}_o)^H] = \sigma_v^2 E[(\mathbf{X}^H \mathbf{X})^{-1}]$$

- Proofs are in the text

Least Squares Filters

$$e(n) = y(n) - \sum_{k=0}^{M-1} h(k)x(n-k)$$

$$= y(n) - \mathbf{c}^T \mathbf{x}(n)$$

$$\mathbf{x}(n) \triangleq [\mathbf{x}(n) \quad \mathbf{x}(n-1) \quad \dots \quad \mathbf{x}(n-M+1)]^T$$

- Several things change we we consider the lagged window case
 - The observations are stochastic, not deterministic
 - The estimated correlation matrix $\hat{\mathbf{R}}$ has a relationship to the estimated correlation of the stochastic process
 - Edge effects of the signals mean our input vectors are not always complete

Edge Conditions

$$\begin{bmatrix} e(0) \\ e(1) \\ \vdots \\ e(M-2) \\ e(M-1) \\ \vdots \\ e(N-1) \\ e(N) \\ \vdots \\ e(N+M-3) \\ e(N+M-2) \end{bmatrix}_{(N+M-1) \times 1} = \begin{bmatrix} y(0) \\ y(1) \\ \vdots \\ y(M-2) \\ y(M-1) \\ \vdots \\ y(N-1) \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}_{(N+M-1) \times 1} - \begin{bmatrix} x(0) & 0 & \dots & 0 \\ x(1) & x(0) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ x(M-2) & x(M-3) & \dots & 0 \\ x(M-1) & x(M-2) & \dots & x(0) \\ \vdots & \vdots & \ddots & \vdots \\ x(N-1) & x(N-2) & \dots & x(N-M) \\ 0 & x(N-1) & \dots & x(N-M+1) \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & x(N-2) \\ 0 & 0 & \dots & x(N-1) \end{bmatrix}_{(N+M-1) \times M} \begin{bmatrix} c_0 \\ c_1 \\ \vdots \\ c_{M-1} \end{bmatrix}_{M \times 1}$$

Computing the Correlation Matrix

$$\hat{r}_{i+1,j+1} = x^*(N_i - i)x(N_i - j) - x^*(N_f + 1 - i)x(N_f + 1 - j) + \hat{r}_{ij}$$

- The data matrix \mathbf{X} is toeplitz
- This gives a structure to $\mathbf{X}^H \mathbf{X}$ that can be used to increase computational efficiency
- However, $\hat{\mathbf{R}}$ is not necessarily toeplitz
 - Depends on the data matrix

Derivation of Correlation Matrix Recursion

$$\begin{aligned} \hat{r}_{ij} &= \sum_{n=N_i}^{N_f} x^*(n+1-i)x(n+1-j) \\ \hat{r}_{i+1,j+1} &= \sum_{n=N_i}^{N_f} x^*(n+1-(i+1))x(n+1-(j+1)) \\ &= \sum_{n=N_i-1}^{N_f-1} x^*(n+1-i)x(n+1-j) \\ &= x^*(N_i-1+1-i)x(N_i-1+1-j) \\ &\quad - x^*(N_f+1-i)x(N_f+1-j) + \sum_{n=N_i}^{N_f} x^*(n+1-i)x(n+1-j) \\ &= x^*(N_i-i)x(N_i-j) - x^*(N_f+1-i)x(N_f+1-j) + \hat{r}_{ij} \end{aligned}$$

Correlation Matrix Recursions

$$\hat{r}_{i+1,j+1} = x^*(N_i - i)x(N_i - j) - x^*(N_f + 1 - i)x(N_f + 1 - j) + \hat{r}_{ij}$$

- Once the first row of \hat{R} is calculated, the recursion above can be used to fill out the rest of the matrix
- Reduces the computation from $\mathcal{O}(NM^2)$ to $\mathcal{O}(NM)$
- Can reduce even further to $\mathcal{O}(N \log N)$ via the FFT if the first row is equivalent to the signal correlation estimates discussed last term

Windowing

$$E_e = \sum_{n=N_i}^{N_f} |e(n)|^2 = \mathbf{e}^H \mathbf{e}$$

There are four ways to select the range for LSE estimation

- **Prewindowing:** $N_i = 0$ and $N_f = N - 1$
 - Essentially treats $x(-1), \dots, x(-M + 1)$ equal to zero
 - Used in adaptive filters mostly for sake of simplicity
- **Postwindowing:** $N_i = M - 1$ and $N_f = N + M - 2$
 - No one uses this
 - Included for completeness only

Windowing Continued

- **Short/No Windowing:** $N_i = M - 1$ and $N_f = N - 1$
 - Use only available data
 - No artificial data
 - No distortions
 - Unbiased, highest variance
 - Sometimes called the **autocorrelation method**
- **Tall/Full windowing:** $N_i = 0$ and $N_f = N + M - 2$
 - \hat{R} becomes toeplitz (efficient order recursions)
 - Sometimes called the **covariance method**
 - Equivalent to calculating \hat{R} with the biased signal correlation estimate (ECE 5/638)
 - Tip: make sure data is zero mean or detrended

Unbiased Autocorrelation Estimate?

- Full windowing is equivalent to using the biased correlation estimate discussed last term
- Conceptually could also use the unbiased estimate
- Not discussed in text
- Properties
 - \hat{R} is toeplitz (by construction)
 - \hat{R} may not be positive definite
 - Estimated AR process models may be unstable — but does it matter?
 - Uses all of the data to calculate every element of \hat{R}

Forward and Backward Linear Prediction (FBLP)

$$\hat{x}_f(n) = - \sum_{k=1}^M a_k(n)x(n-k) = -\mathbf{a}^T(n)\mathbf{x}(n-1)$$

$$\hat{x}_b(n) = - \sum_{k=0}^{M-1} b_k(n)x(n-k) = -\mathbf{b}^T(n)\mathbf{x}(n)$$

$$e_f = x(n) + \sum_{k=1}^M a_k(n)x(n-k) = x(n) + \mathbf{a}^T(n)\mathbf{x}(n-1)$$

$$e_b = \sum_{k=0}^{M-1} b_k(n)x(n-k) + x(n-M) = \mathbf{b}^T(n)\mathbf{x}(n) + x(n-M)$$

- Recall that for *stationary* stochastic processes, the optimum MMSE forward and backward linear predictors have conjugate symmetry

$$\mathbf{a}_o = J\mathbf{b}_o^*$$

Forward and Backward Linear Prediction Continued

$$\mathbf{a}_o = J\mathbf{b}_o^*$$

- This symmetry stems from the Toeplitz structure of the autocorrelation matrix
- If the estimated autocorrelation matrix doesn't have it, perhaps we could improve performance by minimize the forward and backward sum of squared errors

Forward and Backward Linear Prediction Continued

$$\bar{\mathbf{X}}_{(N+M) \times (M+1)} \triangleq \begin{bmatrix} x(0) & 0 & \dots & 0 \\ x(1) & x(0) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ x(M) & x(M-1) & \dots & x(0) \\ \vdots & \vdots & \ddots & \vdots \\ x(N-1) & x(N-2) & \dots & x(N-M-1) \\ 0 & x(N-1) & \dots & x(N-M) \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & x(N-1) \end{bmatrix}$$

$$\mathbf{e}_f \triangleq \bar{\mathbf{X}} \begin{bmatrix} 1 \\ \mathbf{a}_{fb} \end{bmatrix} \quad \mathbf{e}_b^* \triangleq \bar{\mathbf{X}}^* \begin{bmatrix} \mathbf{b}^* \\ 1 \end{bmatrix} = \bar{\mathbf{X}}^* J \begin{bmatrix} 1 \\ \mathbf{a}_{fb} \end{bmatrix}$$

Forward and Backward Linear Prediction Continued

$$\begin{bmatrix} \mathbf{e}_f \\ \mathbf{e}_b^* \end{bmatrix} = \begin{bmatrix} \bar{\mathbf{X}} \\ \bar{\mathbf{X}}^* J \end{bmatrix} \begin{bmatrix} 1 \\ \mathbf{a}_{fb} \end{bmatrix} \\ = \mathbf{x}_{fb} + \begin{bmatrix} \bar{\mathbf{X}} \\ \bar{\mathbf{X}}^* J \end{bmatrix} \begin{bmatrix} \mathbf{a}_{fb} \end{bmatrix}$$

$$\begin{bmatrix} e_f \\ e_b^* \end{bmatrix} = \begin{bmatrix} x(0) \\ \vdots \\ x(M-1) \\ \hline x(M) \\ \vdots \\ x(N-1) \\ \hline 0 \\ \vdots \\ 0 \\ \hline 0 \\ \vdots \\ 0 \\ \hline x^*(0) \\ \vdots \\ x^*(N-M-1) \\ \hline x^*(N-M) \\ \vdots \\ x^*(N-1) \end{bmatrix} + \begin{bmatrix} 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x(M-2) & x(M-3) & \dots & x(0) & 0 \\ \hline x(M-1) & x(M-2) & \dots & x(1) & x(0) \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x(N-2) & x(N-3) & \dots & x(N-M) & x(N-M-1) \\ \hline x(N-1) & x(N-2) & \dots & x(N-M+1) & x(N-M) \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & x(N-1) \\ \hline 0 & 0 & \dots & 0 & x^*(0) \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x^*(0) & x^*(1) & \dots & x^*(M-2) & x^*(M-1) \\ \hline x^*(1) & x^*(2) & \dots & x^*(M-1) & x^*(M) \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x^*(N-M) & x^*(N-M+1) & \dots & x^*(N-2) & x^*(N-1) \\ \hline x^*(N-M+1) & x^*(N-M+2) & \dots & x^*(N-1) & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 0 \end{bmatrix} \mathbf{a}_{fb}$$

Forward and Backward Linear Prediction Continued

$$\mathbf{e}_{fb} = \mathbf{x}_{fb} + \mathbf{X}_{fb} \mathbf{a}_{fb} \quad \mathbf{x}_{fb} \triangleq \begin{bmatrix} x \\ 0 \\ 0 \\ x^* \end{bmatrix} \quad \mathbf{X}_{fb} \triangleq \begin{bmatrix} 0 \\ \mathbf{X} \\ \mathbf{X}^* \mathbf{J} \\ 0 \end{bmatrix}$$

We've already solved this problem

$$\begin{aligned} \mathbf{a}_{ls} &= -(\mathbf{X}_{fb}^H \mathbf{X}_{fb})^{-1} \mathbf{X}_{fb}^H \mathbf{x}_{fb} \\ \mathbf{X}_{fb}^H \mathbf{X}_{fb} &= \begin{bmatrix} 0 & \mathbf{X}^H & \mathbf{J} \mathbf{X}^T & 0 \end{bmatrix} \begin{bmatrix} 0 \\ \mathbf{X} \\ \mathbf{X}^* \mathbf{J} \\ 0 \end{bmatrix} \\ &= \mathbf{X}^H \mathbf{X} + \mathbf{J} \mathbf{X}^T \mathbf{X}^* \mathbf{J} \\ &\propto \hat{\mathbf{R}} + \mathbf{J} \hat{\mathbf{R}}^* \mathbf{J} \\ \mathbf{X}_{fb}^H \mathbf{x}_{fb} &\propto \hat{\mathbf{r}}_f + \hat{\mathbf{r}}_f^* \end{aligned}$$

Forward and Backward Linear Prediction Continued

$$\hat{\mathbf{R}}_{fb} = \mathbf{X}^H \mathbf{X} + \mathbf{J} \mathbf{X}^T \mathbf{X}^* \mathbf{J}$$

- Makes the autocorrelation matrix more symmetric

$$\hat{\mathbf{R}}_{fb} = \mathbf{J} \hat{\mathbf{R}}_{fb} \mathbf{J}$$

- If no windowing is used, is sometimes called the **modified covariance method**
- With full windowing

$$\mathbf{a}_{fb} = (\mathbf{a} + \mathbf{J} \mathbf{b}^*)/2$$

- Works really well for AR signal modeling and parametric spectral estimation (more later)

Summary of Least Squares Filter Estimation

Technique	PD	Unbiased	FFT	Toeplitz	WLS
Pre-Windowing	✓				✓
Post-Windowing	✓				✓
Full-Windowing	✓		✓	✓	✓
Short-Windowing	✓	✓			✓
Unbiased		✓	✓		
Forward/Backward	✓	*		*	✓

*:Can be, depending on windowing technique used. *:Persymmetric, $\hat{\mathbf{R}} = \mathbf{J} \hat{\mathbf{R}} \mathbf{J}$.

- Forward/backward limited to one-step prediction applications and stationary segments.
 - Works best when these conditions are met
- Otherwise short-windowing or unbiased is probably best, depending on importance of $\hat{\mathbf{R}}$ being positive definite

Narrowband Inteferece

- Many types of signals are corrupted by **narrowband interference**
- Typically electrical noise from electrical lines, radio frequency interference, and electrical equipment
- Sensors can pick up this noise through many different means
 - Closed loops (magnetic coupling)
 - Capacitive (electrostatic) coupling
 - Electromagnetic radiation (wires = antennas)
 - Power line interference
 - Voltage drops on ground lines or power lines

Narrowband Inteferece Properties

$$x(n) = s(n) + y(n) + v(n)$$

where

$s(n)$ = signal of interest

$y(n)$ = narrowband interference

$v(n)$ = white noise

- In many applications signal is broadband and unpredictable
- Cannot be separated from white noise with a linear filter (spectral overlap)
- Narrowband interference is predictable
- Consists of sharp peaks in the frequency domain
- Can be eliminated with notch filters, but must know the frequencies

Narrowband Inteferece Cancellation Concept

$$x(n) = s(n) + y(n) + v(n)$$

Suppose all three components of $x(n)$ are mutually uncorrelated

$$\begin{aligned} r_x(\ell) &\triangleq E[x(n)x^*(n-\ell)] \\ &= E[(s(n) + y(n) + v(n))(s^*(n-\ell) + y^*(n-\ell) + v^*(n-\ell))] \\ &= r_s(\ell) + r_y(\ell) + \sigma_v^2\delta(\ell) \end{aligned}$$

- Suppose $r_s(\ell) \approx 0$ for $\ell > D$ and $r_y(\ell) \neq 0$ for $\ell > D$
- This implies $s(n)$ is broadband and $r_y(\ell)$ is “narrowband”
- This means if we try to predict $x(n)$ D steps ahead, only part of the $y(n)$ component will be (partially) predictable

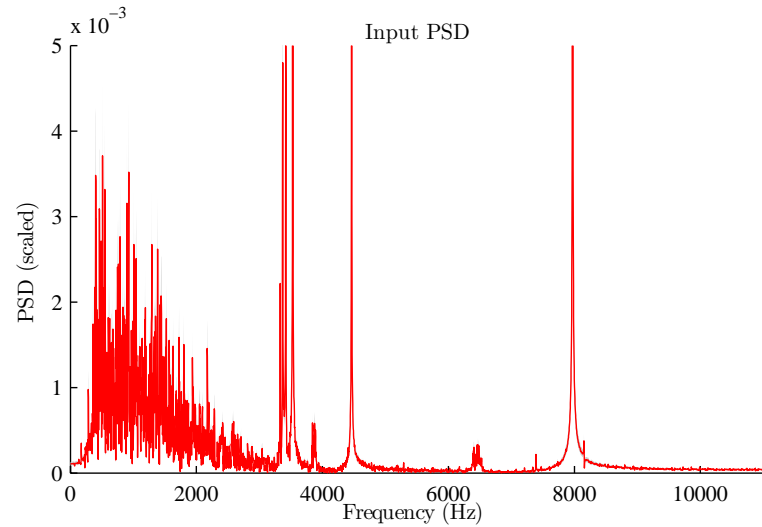
$$\hat{x}_{n-D}(n) = \hat{y}_{n-D}(n)$$

- The difference will then contain less of the interference

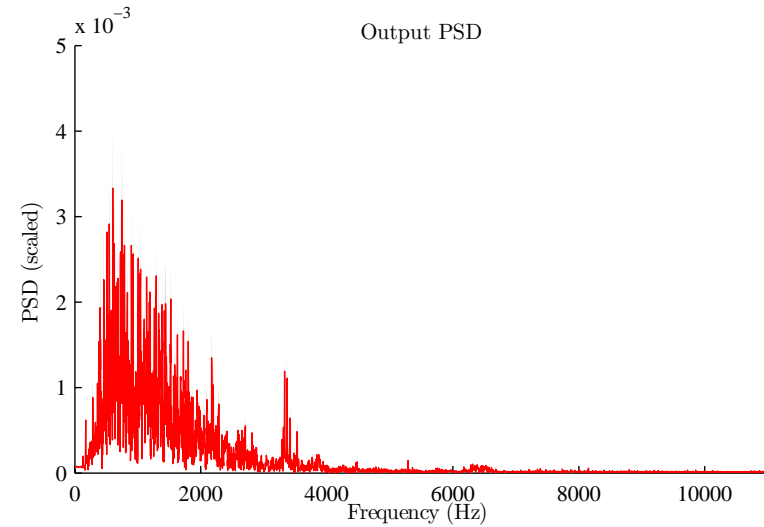
Example 1: Microelectrode Narrowband Interference

- Microelectrode recordings often contain narrowband interference
- The signal of interest are spikes that can be modeled as a point process
- In most cases, the spikes are not predictable with linear filters
- The duration of the spikes is typically 1 ms
- Use a narrowband interference canceller to eliminate the narrowband interference

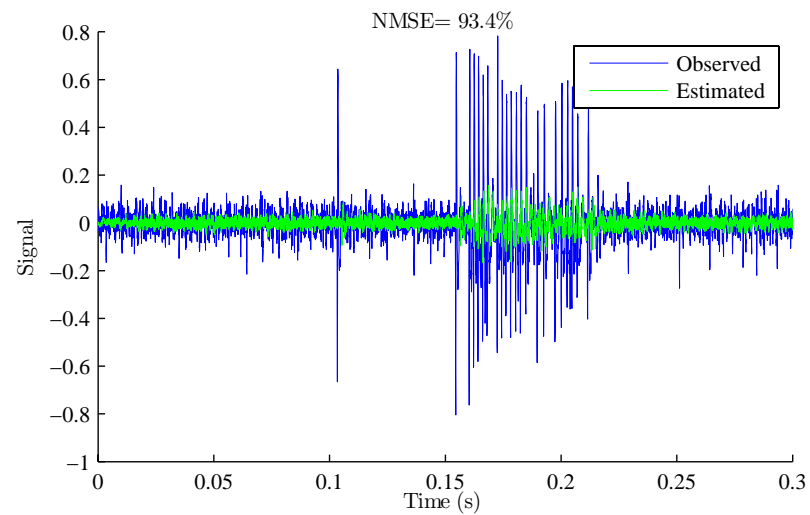
Example 1: Signal PSD



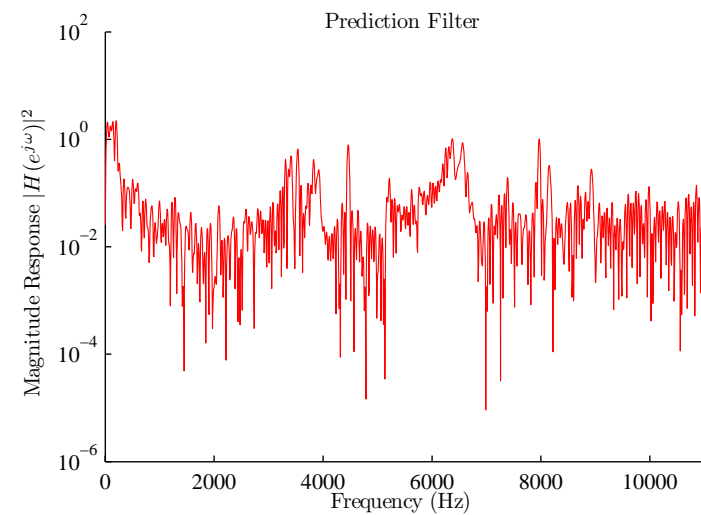
Example 1: Output PSD



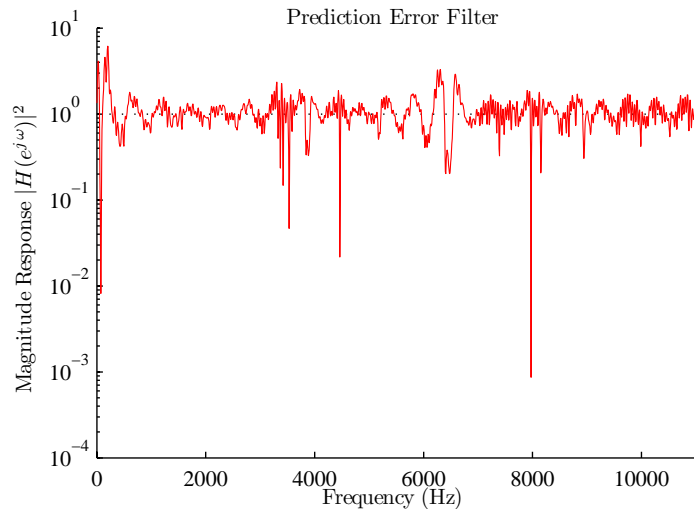
Example 1: Input and Predicted Signal



Example 1: Prediction Filter Frequency Response



Example 1: Prediction Error Filter Frequency Response



Example 1: MATLAB Code

```
clear all;
close all;

%=====
% Parameters
%=====
T = 5; % Signal duration (s)
td = 2e-3; % Prediction delay (s)
tf = 50e-3; % Filter window duration (s)

%=====
% Load the Data
%=====
load MER.mat;

%=====
% Preprocessing
%=====
nx = ceil(fs*T); % Duration of extracted segment
d = round(td*fs); % Delay in units of samples
fo = round(tf*fs);
y = x(d+(1:nx)); % Extract the target output
x = x(1:nx); % Extract a segment of the signal

my = mean(y); % Mean of target signal

%=====
% Estimate the Coefficients
%=====
[c,yh] = LeastSquaresFilter(x,y,fs,500,1,1);

%=====
% Post Processing
%=====
```

```
hpe = [1;zeros(d-1,1);-c]; % Impulse response of the prediction error filter
e = y-yh; % Residuals

%=====
% Figures
%=====
BlackmanTukey(y,fs,1);
FigureSet(1,'Slides');
ylim([0 0.005]);
title('Input PSD');
AxisSet(8);
print('NBISignalPSD','-depsc');

BlackmanTukey(y-yh,fs,1);
FigureSet(1,'Slides');
ylim([0 0.005]);
title('Output PSD');
AxisSet(8);
print('NBIOutputPSD','-depsc');

figure;
FigureSet(1,'Slides');
k = 1:nx;
t = (k-0.5)/fs;
h = plot(t,y,'b',t,yh,'g');
set(h,'LineWidth',0.2);
xlim([0 0.3]);
xlabel('Time (s)');
ylabel('Signal');
legend('Observed','Estimated');
set(get(gca,'Title'),'Interpreter','LaTeX');
title(sprintf('NMSE=%5.1f%%',100*sum((y-yh).^2)/sum((y-my).^2)));
box off;
AxisSet(8);
print('NBISignalEstimate','-depsc');

figure;
```

```
FigureSet(2,'Slides');
[h,f] = freqz(c,1,2^11,fs);
h = semilogy(f,abs(h).^2,'r');
set(h,'LineWidth',0.4);
xlim([0 fs/2]);
xlabel('Frequency (Hz)');
set(get(gca,'YLabel'),'Interpreter','LaTeX');
ylabel('Magnitude Response |H(e^{j\omega})|^2$');
title('Prediction Filter');
box off;
AxisSet(8);
print('NBIPredictionFrequencyResponse','-depsc');

figure;
FigureSet(2,'Slides');
[h,f] = freqz(hpe,1,2^11,fs);
semilogy([0 fs/2],[1 1],'k');
hold on;
h = semilogy(f,abs(h).^2,'r');
set(h,'LineWidth',0.4);
hold off;
xlim([0 fs/2]);
xlabel('Frequency (Hz)');
set(get(gca,'YLabel'),'Interpreter','LaTeX');
ylabel('Magnitude Response |H(e^{j\omega})|^2$');
title('Prediction Error Filter');
box off;
AxisSet(8);
print('NBIPredictionErrorFrequencyResponse','-depsc');
```

Example 1: Least Squares Filter MATLAB Code

```
function [c,yh] = LeastSquaresFilter(x,y,fsa,foa,wta,pfa);
%LeastSquaresFilter: Least squares estimate of FIR filter coefficients
%
% [c,yh] = NonparametricSpectrogram(x,y,fs,wl,fr,nf,ns,pf);
%
% x Input signal.
% y Target signal.
% fs Sample rate (Hz). Default = 1 Hz.
% fo Filter order. Default = 10.
% wt Window type: 0=full (default), 1=none,
% 2=unbiased autocorrelation estimate
% pf Plot flag: 0=none (default), 1=screen, 2=current figure.
%
% c Vector of coefficients.
% yh Estimate of y using the estimator.
%
% Calculates the least squares estimate of the coefficient vector
% c using the input data. Efficiently calculates the
% autocorrelation matrix using the recursive approach described
% in Manolakis.
%
% Example: Estimate the coefficients for doing narrowband
% interference of a microelectrode recording.
%
% load MER.mat;
% d = round(2e-3*fs);
% y = x(d+(1:50e3));
% x = x(1:50e3);
% [c,yh] = LeastSquaresFilter(x,y,fs,500,1,1);
%
% D. G. Manolakis, V. K. Ingle, S. M. Kogon, "Statistical and
% adaptive signal processing," Artech House, 2005.
%
% Version 1.00 JM
```

```
%
% See also signal, armcov, arcov, and aryule.
%
%=====
% Error Checking
%=====
if nargin<2,
    help LeastSquaresFilter;
    return;
end;

if isempty(x) | isempty(y),
    error('Signal is empty.');
```

```
end;

if length(x)~=length(y),
    error('Input signals are different lengths.');
```

```
end;

if var(x)==0 | var(y)==0,
    error('Signal is constant.');
```

```
end;

%=====
% Calculate Basic Signal Statistics
%=====
nx = length(x); % No. samples in x
mx = mean(x); % Input signal mean
my = mean(y); % Target signal mean

%=====
% Process Function Arguments
%=====
fs = 1; % Default sample rate
if exist('fsa','var') & ~isempty(fsa),
    fs = fsa;
end;
```

```
fo = 10; % Default filter order
if exist('foa','var') & ~isempty(foa),
    fo = foa;
end;

wt = 0; % Default filter order
if exist('wta','var') & ~isempty(wta),
    wt = wta;
end;

pf = 0; % Default - no plotting
if nargout==0, % Plot if no output arguments
    pf = 1;
end;
if exist('pfa') & ~isempty(pfa),
    pf = pfa;
end;

%=====
% Preprocessing
%=====
x = x(:); % Make into a column vector
%x = x - mx; % Remove mean

%=====
% Estimate the ACF
%=====
if wt==0 | wt==2,
    np = 2*nextpow2(2*nx-1); % Figure out how much to pad the signal
    X = fft(x,np);
    xc = ifft(abs(X).^2);
    ac = real(xc(1:fo));

    Y = fft(y,np);
    cp = ifft(Y.*conj(X));
    cc = real(cp(1:fo));
```

```
if wt==0, % Biased estimate
    ac = ac./nx;
    cc = cc./nx;
else % Unbiased estimate
    k = (0:fo-1).';
    ac = ac./(nx-k);
    cc = cc./(nx-k);
end;

%=====
% Create the Estimated Autocorrelation (R) and Cross-Correlation (d)
%=====
R = zeros(fo,fo);
d = zeros(fo,1);
switch wt,
case {0,2},
    for c1=1:fo,
        d(c1) = cc(c1);
        for c2=c1:fo,
            R(c1,c2) = ac(c2-c1+1);
            R(c2,c1) = R(c1,c2);
        end;
    end;
case 1,
    %-----
    % Calculate the First Row
    %-----
    for c1=1:fo,
        d(c1) = sum(y(fo:nx).*x(fo-(c1-1):nx-(c1-1)));
        R(1,c1) = sum(x(fo:nx).*x(fo-(c1-1):nx-(c1-1)));
        R(c1,1) = R(1,c1);
    end;
    %-----
    % Fill out the Remainder of the Matrix
    %-----
    for c1=2:fo,
```



```

        for c2=c1:fo,
            R(c1,c2) = R(c1-1,c2-1) + x(fo-(c1-1)).*x(fo-(c2-1)) - x(nx-(c1-2)).*x(nx-(c2-2));
            R(c2,c1) = R(c1,c2);
        end;
    end;
    %X = zeros(nx-fo+1,fo);           % Verification code
    %for c1=1:fo,
    %    X(:,c1) = x(fo-(c1-1):nx-(c1-1));
    %    end;
    %R2 = X'*X;
    %d2 = X'*y(fo:nx);
    %disp([max(max(abs(R-R2))) max(abs(d-d2))]);
end;

%=====
% Calculate the Coefficients
%=====
c = pinv(R)*d;
yh = filter(c,1,x);

%=====
% Plot Results
%=====
if pf>=1,
    if pf~=2,
        figure;
        end;
    FigureSet(1);

    k = 1:nx;
    t = (k-0.5)/fs;
    h = plot(t,y,'r',t,yh,'g');
    set(h,'LineWidth',1.2);
    xlim([0 nx/fs]);
    xlabel('Time (s)');
    ylabel('Signal');
    legend('Observed','Estimated');
end;

```

```

set(get(gca,'Title'),'Interpreter','LaTeX');
title(sprintf('NMSE=%5.1f\\%',100*sum((y-yh).^2)/sum((y-my).^2)));
box off;
AxisSet;

if pf~=2,
    figure;
    end;
FigureSet(2);
[h,f] = freqz(c,1,2^12,fs);
h = semilogy(f,abs(h).^2,'r');
set(h,'LineWidth',1.2);
xlim([0 fs/2]);
xlabel('Frequency (Hz)');
set(get(gca,'YLabel'),'Interpreter','LaTeX');
ylabel('Magnitude Response  $|H(\exp\{j\omega\})|^2$ ');
box off;
AxisSet;

end;

%=====
% Process Return Arguments
%=====
if nargin==0,
    clear('c','yh');
end;

```