

Adaptive Filtering - Theory and Applications

José C. M. Bermudez

Department of Electrical Engineering
Federal University of Santa Catarina
Florianópolis – SC
Brazil

IRIT - INP-ENSEEIH, Toulouse
May 2011

- 1 Introduction
- 2 Adaptive Filtering Applications
- 3 Adaptive Filtering Principles
- 4 Iterative Solutions for the Optimum Filtering Problem
- 5 Stochastic Gradient Algorithms
- 6 Deterministic Algorithms
- 7 Analysis

Introduction

Estimation Techniques

Several techniques to solve estimation problems.

- Classical Estimation

Maximum Likelihood (ML), Least Squares (LS), Moments, etc.

- Bayesian Estimation

Minimum MSE (MMSE), Maximum A Posteriori (MAP), etc.

Linear Estimation

Frequently used in practice when there is a limitation in computational complexity – *Real-time operation*

Linear Estimators

Simpler to determine: depend on the first two moments of data

- **Statistical Approach** – Optimal Linear Filters
 - ▶ Minimum Mean Square Error
 - ▶ Require second order statistics of signals
- **Deterministic Approach** – Least Squares Estimators
 - ▶ Minimum Least Squares Error
 - ▶ Require handling of a data observation matrix

Limitations of Optimal Filters and LS Estimators

- Statistics of signals may not be available or cannot be accurately estimated
- There may not be available time for statistical estimation (real-time)
- Signals and systems may be non-stationary
- Memory required may be prohibitive
- Computational load may be prohibitive

Iterative Solutions

- Search the optimal solution starting from an initial guess
- Iterative algorithms are based on classical optimization algorithms
- Require reduced computational effort per iteration
- Need several iterations to converge to the optimal solution
- These methods form the basis for the development of adaptive algorithms
- Still require the knowledge of signal statistics

Adaptive Filters

Usually **approximate iterative** algorithms and:

- Do not require previous knowledge of the signal statistics
- Have a small computational complexity per iteration
- Converge to a neighborhood of the optimal solution

Adaptive filters are good for:

- Real-time applications, when there is no time for statistical estimation
- Applications with nonstationary signals and/or systems

Properties of Adaptive Filters

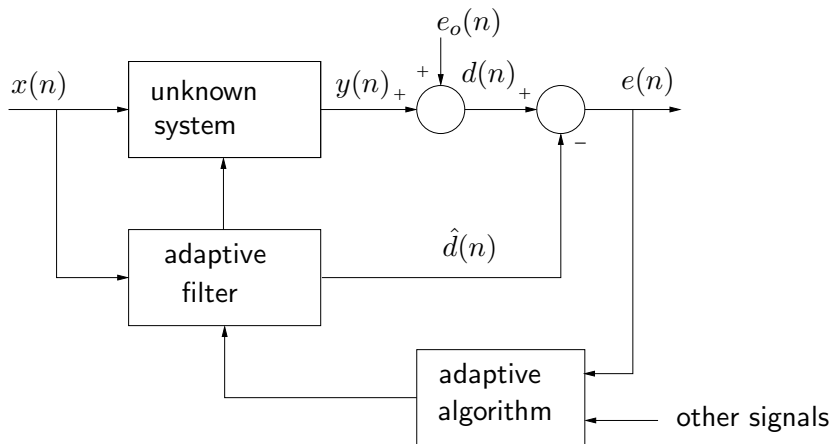
- They can operate satisfactorily in unknown and possibly time-varying environments without user intervention
- They improve their performance during operation by learning statistical characteristics from current signal observations
- They can track variations in the *signal operating environment* (SOE)

Adaptive Filtering Applications

Basic Classes of Adaptive Filtering Applications

- System Identification
- Inverse System Modeling
- Signal Prediction
- Interference Cancellation

System Identification



Applications – System Identification

Channel Estimation

- Communications systems
- Objective: model the channel to design distortion compensation
- $x(n)$: training sequence

Plant Identification

- Control systems
- Objective: model the plant to design a compensator
- $x(n)$: training sequence

Echo Cancellation

- Telephone systems and VoIP
- Echo caused by network impedance mismatches or acoustic environment
- Objective: model the echo path impulse response
- $x(n)$: transmitted signal
- $d(n)$: echo + noise

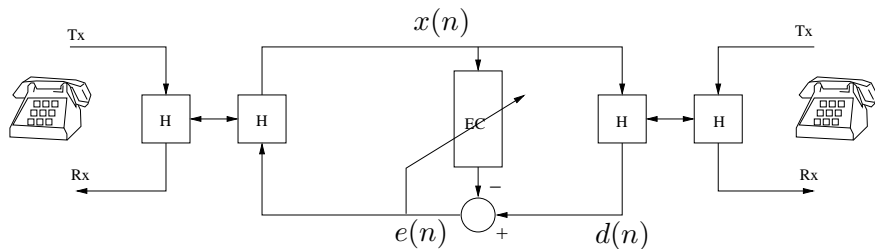
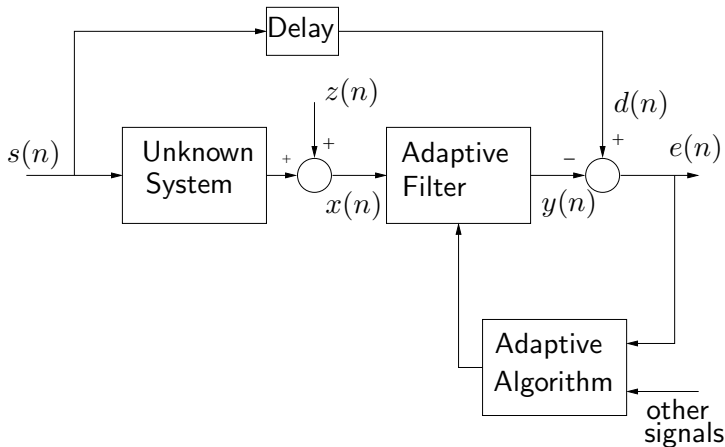


Figure: Network Echo Cancellation

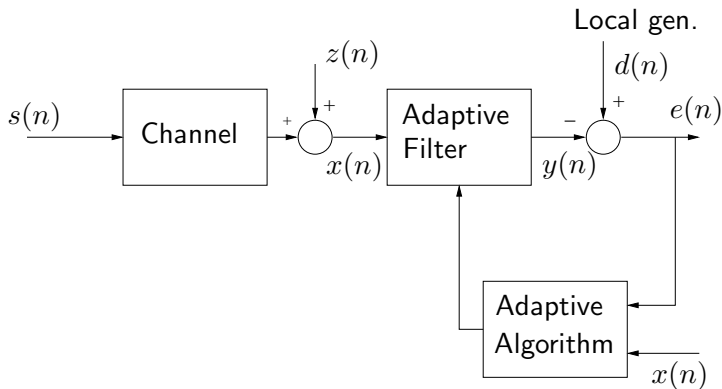
Inverse System Modeling

- Adaptive filter attempts to estimate unknown system's inverse
- Adaptive filter input usually corrupted by noise
- Desired response $d(n)$ may not be available



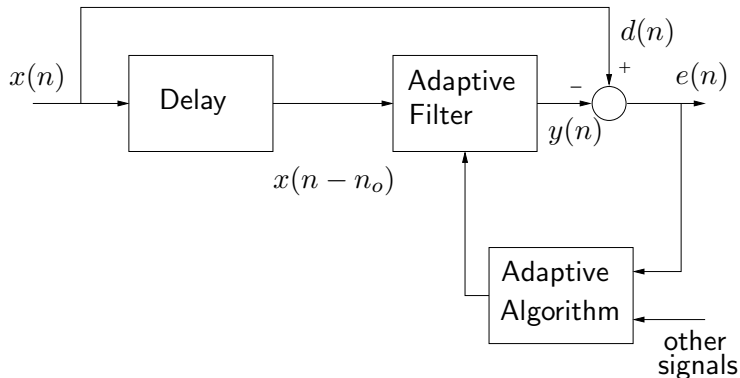
Applications – Inverse System Modeling

Channel Equalization



- Objective: reduce intersymbol interference
- Initially – training sequence in $d(n)$
- After training: $d(n)$ generated from previous decisions

Signal Prediction

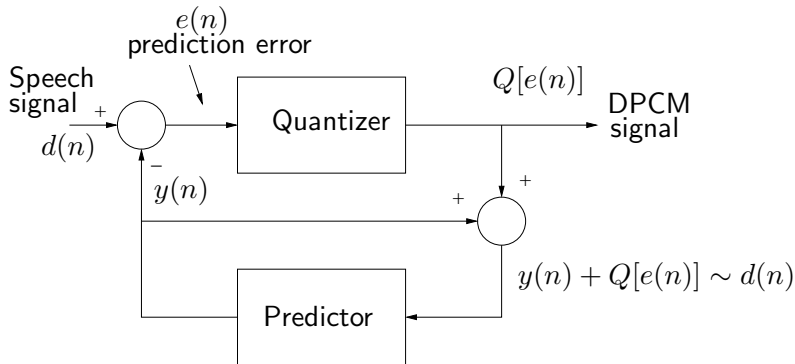


- most widely used case – forward prediction
- signal $x(n)$ to be predicted from samples $\{x(n - n_o), x(n - n_o - 1), \dots, x(n - n_o - L)\}$

Application – Signal Prediction

DPCM Speech Quantizer - Linear Predictive Coding

- Objective: Reduce speech transmission bandwidth
- Signal transmitted all the time: quantization error
- Predictor coefficients are transmitted at low rate

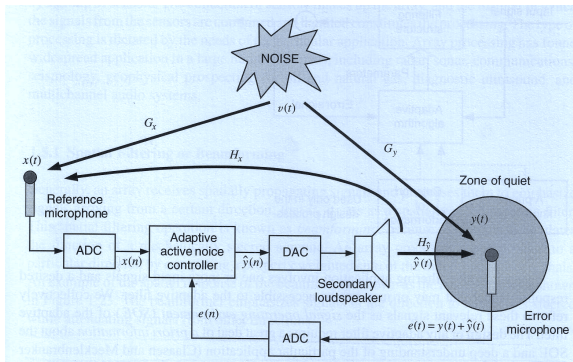


Interference Cancellation

- One or more sensor signals are used to remove interference and noise
- Reference signals correlated with the interference should also be available
- Applications:
 - ▶ array processing for radar and communications
 - ▶ biomedical sensing systems
 - ▶ active noise control systems

Application – Interference Cancellation

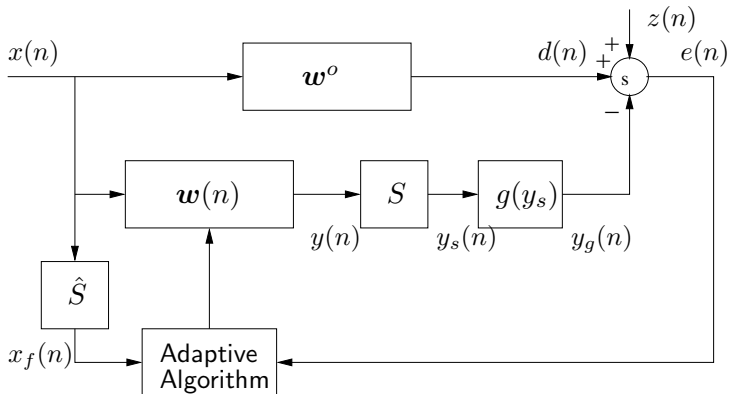
Active Noise Control



Ref: D.G. Manolakis, V.K. Ingle and S.M. Kogon, *Statistical and Adaptive Signal Processing*, 2000.

- Cancellation of acoustic noise using destructive interference
- Secondary system between the adaptive filter and the cancellation point is unavoidable
- Cancellation is performed in the acoustic environment

Active Noise Control – Block Diagram



Adaptive Filtering Principles

Adaptive Filter Features

Adaptive filters are composed of three basic modules:

- **Filtering structure**

- ▶ Determines the **output of the filter** given its input samples
- ▶ Its weights are periodically adjusted by the adaptive algorithm
- ▶ Can be linear or nonlinear, depending on the application
- ▶ Linear filters can be FIR or IIR

- **Performance criterion**

- ▶ Defined according to application and **mathematical tractability**
- ▶ Is used to derive the adaptive algorithm
- ▶ Its value at each iteration affects the adaptive weight updates

- **Adaptive algorithm**

- ▶ Uses the performance criterion value and the current signals
- ▶ Modifies the adaptive weights to improve performance
- ▶ Its form and complexity are function of the structure and of the performance criterion

Signal Operating Environment (SOE)

Comprises all informations regarding the properties of the signals and systems

- Input signals
- Desired signal
- Unknown systems

If the SOE is nonstationary

- **Aquisition or convergence mode:** from start until close to best performance
- **Tracking mode:** readjustment following SOE's time variations

Adaptation can be

- **Supervised** – desired signal is available $\Leftrightarrow e(n)$ can be evaluated
- **Unsupervised** – desired signal is unavailable

Performance Evaluation

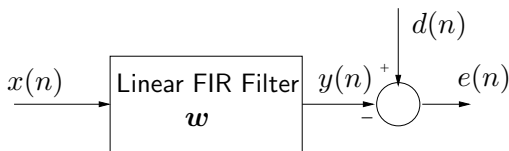
- Convergence rate
- Misadjustment
- Tracking
- Robustness (disturbances and numerical)
- Computational requirements (operations and memory)
- Structure
 - ▶ facility of implementation
 - ▶ performance surface
 - ▶ stability

Optimum versus Adaptive Filters in Linear Estimation

Conditions for this study

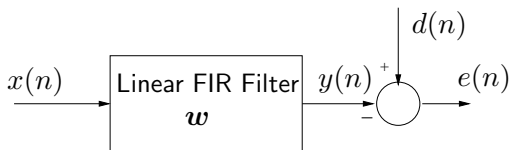
- Stationary SOE
- Filter structure is transversal FIR
- All signals are real valued
- Performance criterion: Mean-square error $E[e^2(n)]$

The Linear Estimation Problem



$$J_{ms} = E[e^2(n)]$$

The Linear Estimation Problem



$$\mathbf{x}(n) = [x(n), x(n-1), \dots, x(n-N+1)]^T$$

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

$$e(n) = d(n) - y(n) = d(n) - \mathbf{x}^T(n)\mathbf{w}$$

$$J_{ms} = E[e^2(n)] = \sigma_d^2 - 2\mathbf{p}^T\mathbf{w} + \mathbf{w}^T\mathbf{R}_{xx}\mathbf{w}$$

where

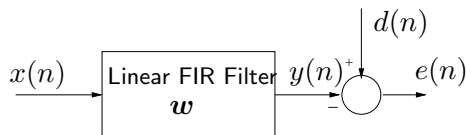
$$\mathbf{p} = E[\mathbf{x}(n)d(n)]; \quad \mathbf{R}_{xx} = E[\mathbf{x}(n)\mathbf{x}^T(n)]$$

Normal Equations

$$\mathbf{R}_{xx}\mathbf{w}^o = \mathbf{p} \quad \Leftrightarrow \quad \mathbf{w}^o = \mathbf{R}_{xx}^{-1}\mathbf{p} \quad \text{for } \mathbf{R}_{xx} > 0$$

$$J_{ms_{min}} = \sigma_d^2 - \mathbf{p}^T\mathbf{R}_{xx}^{-1}\mathbf{p}$$

What if $d(n)$ is nonstationary?



$$\mathbf{x}(n) = [x(n), x(n-1), \dots, x(n-N+1)]^T$$

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

$$e(n) = d(n) - y(n) = d(n) - \mathbf{x}^T(n) \mathbf{w}(n)$$

$$J_{ms}(n) = E[e^2(n)] = \sigma_d^2(n) - 2\mathbf{p}(n)^T \mathbf{w}(n) + \mathbf{w}^T(n) \mathbf{R}_{xx} \mathbf{w}(n)$$

where

$$\mathbf{p}(n) = E[\mathbf{x}(n)d(n)]; \quad \mathbf{R}_{xx} = E[\mathbf{x}(n)\mathbf{x}^T(n)]$$

Normal Equations

$$\mathbf{R}_{xx} \mathbf{w}^o(n) = \mathbf{p}(n) \quad \Leftrightarrow \quad \mathbf{w}^o(n) = \mathbf{R}_{xx}^{-1} \mathbf{p}(n) \quad \text{for } \mathbf{R}_{xx} > 0$$

$$J_{ms_{min}}(n) = \sigma_d^2(n) - \mathbf{p}^T(n) \mathbf{R}_{xx}^{-1} \mathbf{p}(n)$$

Optimum Filters versus Adaptive Filters

Optimum Filters

- Compute
$$\mathbf{p}(n) = E[\mathbf{x}(n)d(n)]$$
- Solve $\mathbf{R}_{xx}\mathbf{w}^o = \mathbf{p}(n)$
- Filter with $\mathbf{w}^o(n) \Leftrightarrow$
$$y(n) = \mathbf{x}^T(n)\mathbf{w}^o(n)$$

Nonstationary SOE:

Optimum filter determined
for each value of n

Adaptive Filters

- Filtering: $y(n) = \mathbf{x}^T(n)\mathbf{w}(n)$
- Evaluate error: $e(n) = d(n) - y(n)$
- Adaptive algorithm:

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \Delta\mathbf{w}[\mathbf{x}(n), e(n)]$$

$\Delta\mathbf{w}(n)$ is chosen so that $\mathbf{w}(n)$ is close to $\mathbf{w}^o(n)$ for n large

Characteristics of Adaptive Filters

- Search for the optimum solution on the performance surface
- Follow principles of optimization techniques
- Implement a recursive optimization solution
- Convergence speed may depend on initialization
- Have stability regions
- Steady-state solution fluctuates about the optimum
- Can track time varying SOEs better than optimum filters
- Performance depends on the performance surface

Iterative Solutions for the Optimum Filtering Problem

Performance (Cost) Functions

- **Mean-square error** $\Leftrightarrow E[e^2(n)]$ (Most popular)
Adaptive algorithms: Least-Mean Square (LMS), Normalized LMS (NLMS), **Affine Projection** (AP), Recursive Least Squares (RLS), etc.
- **Regularized MSE**

$$J_{rms} = E[e^2(n)] + \alpha \|\mathbf{w}(n)\|^2$$

Adaptive algorithm: **leaky** least-mean square (leaky LMS)

- **ℓ_1 norm** criterion

$$J_{\ell_1} = E[|e(n)|]$$

Adaptive algorithm: Sign-Error

Performance (Cost) Functions – continued

- Least-mean fourth (LMF) criterion

$$J_{LMF} = E[e^4(n)]$$

Adaptive algorithm: Least-Mean Fourth (LMF)

- Least-mean-mixed-norm (LMMN) criterion

$$J_{LMMN} = E[\alpha e^2(n) + \frac{1}{2}(1 - \alpha)e^4(n)]$$

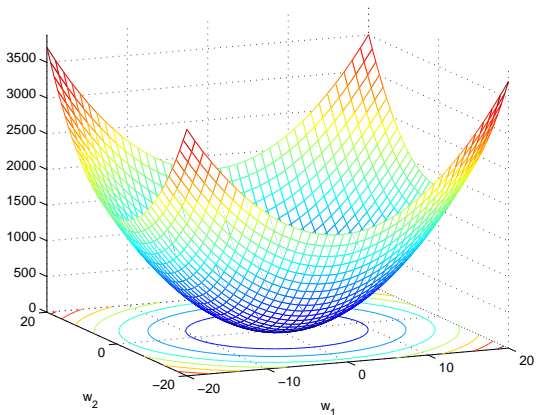
Adaptive algorithm: Least-Mean-Mixed-Norm (LMMN)

- Constant-modulus criterion

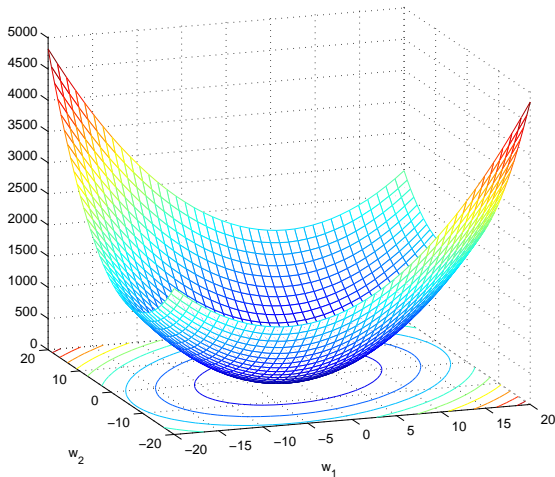
$$J_{CM} = E[(\gamma - |\mathbf{x}^T(n)\mathbf{w}(n)|^2)^2]$$

Adaptive algorithm: Constant-Modulus (CM)

MSE Performance Surface – Small Input Correlation



MSE Performance Surface – Large Input Correlation



The Steepest Descent Algorithm – Stationary SOE

Cost Function

$$J_{ms}(\mathbf{n}) = E[e^2(n)] = \sigma_d^2 - 2\mathbf{p}^T \mathbf{w}(n) + \mathbf{w}^T(n) \mathbf{R}_{xx} \mathbf{w}(n)$$

Weight Update Equation

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu \mathbf{c}(n)$$

μ : step-size

$\mathbf{c}(n)$: correction term (determines direction of $\Delta \mathbf{w}(n)$)

Steepest descent adjustment:

$$\mathbf{c}(n) = -\nabla J_{ms}(n) \quad \Leftrightarrow \quad J_{ms}(n+1) \leq J_{ms}(n)$$

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu [\mathbf{p} - \mathbf{R}_{xx} \mathbf{w}(n)]$$

Weight Update Equation About the Optimum Weights

Weight Error Update Equation

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu[\mathbf{p} - \mathbf{R}_{xx}\mathbf{w}(n)]$$

Using $\mathbf{p} = \mathbf{R}_{xx}\mathbf{w}^o$

$$\mathbf{w}(n+1) = (\mathbf{I} - \mu\mathbf{R}_{xx})\mathbf{w}(n) + \mu\mathbf{R}_{xx}\mathbf{w}^o$$

Weight error vector: $\mathbf{v}(n) = \mathbf{w}(n) - \mathbf{w}^o$

$$\boxed{\mathbf{v}(n+1) = (\mathbf{I} - \mu\mathbf{R}_{xx})\mathbf{v}(n)}$$

- Matrix $\mathbf{I} - \mu\mathbf{R}_{xx}$ must be *stable* for convergence ($|\lambda_i| < 1$)
- Assuming convergence, $\lim_{n \rightarrow \infty} \mathbf{v}(n) = \mathbf{0}$

Convergence Conditions

$$\mathbf{v}(n+1) = (\mathbf{I} - \mu \mathbf{R}_{xx})\mathbf{v}(n); \quad \mathbf{R}_{xx} \text{ positive definite}$$

Eigen-decomposition of \mathbf{R}_{xx}

$$\mathbf{R}_{xx} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$$

$$\mathbf{v}(n+1) = (\mathbf{I} - \mu \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T)\mathbf{v}(n)$$

$$\mathbf{Q}^T \mathbf{v}(n+1) = \mathbf{Q}^T \mathbf{v}(n) - \mu \mathbf{\Lambda} \mathbf{Q}^T \mathbf{v}(n)$$

Defining $\tilde{\mathbf{v}}(n+1) = \mathbf{Q}^T \mathbf{v}(n+1)$

$$\tilde{\mathbf{v}}(n+1) = (\mathbf{I} - \mu \mathbf{\Lambda})\tilde{\mathbf{v}}(n)$$

Convergence Properties

$$\tilde{\mathbf{v}}(n+1) = (\mathbf{I} - \mu\mathbf{\Lambda})\tilde{\mathbf{v}}(n)$$

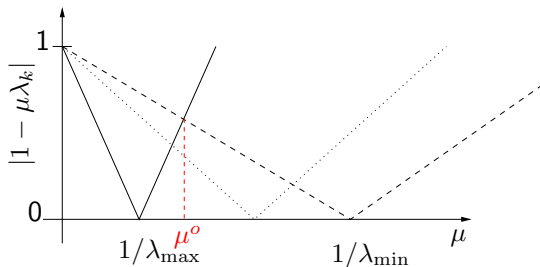
- $\tilde{v}_k(n+1) = (1 - \mu\lambda_k)\tilde{v}_k(n)$, $k = 1, \dots, N$
- $\tilde{v}_k(n) = (1 - \mu\lambda_k)^n \tilde{v}_k(0)$
- Convergence modes
 - ▶ monotonic if $0 < 1 - \mu\lambda_k < 1$
 - ▶ oscillatory if $-1 < 1 - \mu\lambda_k < 0$
- Convergence if $|1 - \mu\lambda_k| < 1 \Leftrightarrow \boxed{0 < \mu < \frac{2}{\lambda_{\max}}}$

Optimal Step-Size

$$\tilde{v}_k(n) = (1 - \mu\lambda_k)^n \tilde{v}_k(0); \text{convergence modes: } 1 - \mu\lambda_k$$

- $\max |1 - \mu\lambda_k|$: slowest mode
- $\min |1 - \mu\lambda_k|$: fastest mode

Optimal Step-Size – continued



$$\min_{\mu} \max\{|1 - \mu \lambda_k|\}$$

$$1 - \mu^o \lambda_{\min} = -(1 - \mu^o \lambda_{\max})$$

$$\boxed{\mu^o = \frac{2}{\lambda_{\max} + \lambda_{\min}}}$$

Optimal slowest modes: $\pm \frac{\rho-1}{\rho+1}$; $\rho = \frac{\lambda_{\max}}{\lambda_{\min}}$

The Learning Curve – $J_{ms}(n)$

$$J_{ms}(n) = J_{ms_{\min}} + \overbrace{\tilde{\mathbf{v}}^T(n) \mathbf{\Lambda} \tilde{\mathbf{v}}(n)}^{\text{Excess MSE}} = J_{ms_{\min}} + \sum_{k=1}^N \lambda_k \tilde{v}_k^2(n)$$

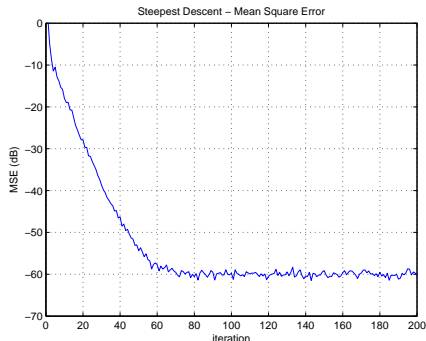
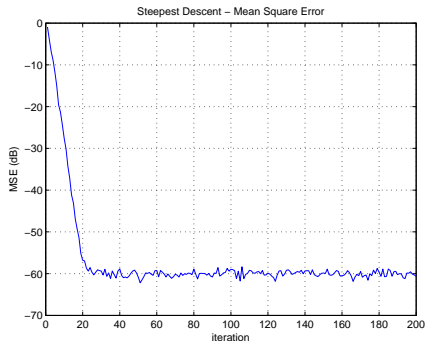
Since $\tilde{v}_k(n) = (1 - \mu\lambda_k)^n \tilde{v}_k(0)$,

$$J_{ms}(n) = J_{ms_{\min}} + \sum_{k=1}^N \lambda_k (1 - \mu\lambda_k)^{2n} \tilde{v}_k^2(0)$$

- $\lambda_k(1 - \mu\lambda_k)^2 > 0 \Leftrightarrow$ monotonic convergence
- Stability limit is again $0 < \mu < \frac{2}{\lambda_{\max}}$
- $J_{ms}(n)$ converges faster than $\mathbf{w}(n)$
- Algorithm converges faster as $\rho = \lambda_{\max}/\lambda_{\min} \rightarrow 1$

Simulation Results

$$x(n) = \alpha x(n-1) + v(n)$$



Input: White noise ($\rho = 1$)

Input: AR(1), $\alpha = 0.7$ ($\rho = 5.7$)

- Linear system identification - FIR with 20 coefficients
- Step-size $\mu = 0.3$
- Noise power $\sigma_v^2 = 10^{-6}$

The Newton Algorithm

- Steepest descent: linear approx. of J_{ms} about the operating point
- Newton's method: Quadratic approximation of J_{ms}

Expanding $J_{ms}(\mathbf{w})$ in Taylor's series about $\mathbf{w}(n)$,

$$J_{ms}(\mathbf{w}) \sim J_{ms}[\mathbf{w}(n)] + \nabla^T J_{ms}[\mathbf{w} - \mathbf{w}(n)] \\ + \frac{1}{2}[\mathbf{w} - \mathbf{w}(n)]^T \mathbf{H}(n)[\mathbf{w} - \mathbf{w}(n)]$$

Differentiating w.r.t. \mathbf{w} and equating to zero at $\mathbf{w} = \mathbf{w}(n+1)$,

$$\nabla J_{ms}[\mathbf{w}(n+1)] = \nabla J_{ms}[\mathbf{w}(n)] + \mathbf{H}[\mathbf{w}(n)][\mathbf{w}(n+1) - \mathbf{w}(n)] = \mathbf{0}$$

$$\boxed{\mathbf{w}(n+1) = \mathbf{w}(n) - \mathbf{H}^{-1}[\mathbf{w}(n)]\nabla J_{ms}[\mathbf{w}(n)]}$$

The Newton Algorithm – continued

- $\nabla J_{ms}[\mathbf{w}(n)] = -2\mathbf{p} + 2\mathbf{R}_{xx}\mathbf{w}(n)$
- $\mathbf{H}(\mathbf{w}(n)) = 2\mathbf{R}_{xx}$

Thus, adding a step-size control,

$$\mathbf{w}(n+1) = \mathbf{w}(n) - \mu \mathbf{R}_{xx}^{-1}[-\mathbf{p} + \mathbf{R}_{xx}\mathbf{w}(n)]$$

- Quadratic surface \Rightarrow conv. in *one* iteration for $\mu = 1$ 😊
- Requires the determination of \mathbf{R}_{xx}^{-1} 😞
- Can be used to derive simpler adaptive algorithms
- When $\mathbf{H}(n)$ is close to singular \Rightarrow regularization

$$\tilde{\mathbf{H}}(n) = 2\mathbf{R}_{xx} + 2\epsilon\mathbf{I}$$

Basic Adaptive Algorithms

Least Mean Squares (LMS) Algorithm

- Can be interpreted in different ways
- Each interpretation helps understanding the algorithm behavior
- Some of these interpretations are related to the steepest descent algorithm

LMS as a Stochastic Gradient Algorithm

- Suppose we use the estimate $J_{ms}(n) = E[e^2(n)] \simeq e^2(n)$
- The estimated gradient vector becomes

$$\hat{\nabla} J_{ms}(n) = \frac{\partial e^2(n)}{\partial \mathbf{w}(n)} = 2e(n) \frac{\partial e(n)}{\partial \mathbf{w}(n)}$$

- Since $e(n) = d(n) - \mathbf{x}^T(n)\mathbf{w}(n)$,

$$\hat{\nabla} J_{ms}(n) = -2e(n)\mathbf{x}(n) \quad (\text{stochastic gradient})$$

and, using the steepest descent weight update equation,

$$\boxed{\mathbf{w}(n+1) = \mathbf{w}(n) + \mu e(n)\mathbf{x}(n)} \quad (\text{LMS weight update})$$

LMS as a Stochastic Estimation Algorithm

- $\nabla J_{ms}(n) = -2\mathbf{p} + 2\mathbf{R}_{xx}\mathbf{w}(n)$
- Stochastic estimators

$$\hat{\mathbf{p}} = d(n)\mathbf{x}(n) \quad \hat{\mathbf{R}}_{xx} = \mathbf{x}(n)\mathbf{x}^T(n)$$

Then,

$$\hat{\nabla} J_{ms}(n) = -2d(n)\mathbf{x}(n) + 2\mathbf{x}(n)\mathbf{x}^T(n)\mathbf{w}(n)$$

- Using $\hat{\nabla} J_{ms}(n)$ is the steepest descent weight update,

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu e(n)\mathbf{x}(n)$$

LMS – A Solution to a Local Optimization

- Error expressions

$$e(n) = d(n) - \mathbf{x}^T(n)\mathbf{w}(n) \quad (\text{a priori error})$$

$$\epsilon(n) = d(n) - \mathbf{x}^T(n)\mathbf{w}(n+1) \quad (\text{a posteriori error})$$

- We want to maximize $|\epsilon(n) - e(n)|$ with $|\epsilon(n)| < |e(n)|$

$$\epsilon(n) - e(n) = -\mathbf{x}^T(n)\Delta\mathbf{w}(n)$$

$$\Delta\mathbf{w}(n) = \mathbf{w}(n+1) - \mathbf{w}(n)$$

- Expressing $\Delta\mathbf{w}(n)$ as $\Delta\mathbf{w}(n) = \Delta\tilde{\mathbf{w}}(n)e(n)$

$$\epsilon(n) - e(n) = -\mathbf{x}^T(n)\Delta\tilde{\mathbf{w}}(n)e(n)$$

- For $\max |\epsilon(n) - e(n)| \Leftrightarrow \Delta\tilde{\mathbf{w}}(n)$ in the direction of $\mathbf{x}(n)$

- $\Leftrightarrow \Delta \tilde{\mathbf{w}}(n) = \mu \mathbf{x}(n)$ and

$$\Delta \mathbf{w}(n) = \mu \mathbf{x}(n) e(n)$$

and

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu e(n) \mathbf{x}(n)$$

- As

$$\epsilon(n) - e(n) = -\mathbf{x}^T(n) \Delta \mathbf{w}(n) = -\mu \mathbf{x}^T(n) \mathbf{x}(n) e(n)$$

$|\epsilon(n)| < |e(n)|$ requires $|1 - \mu \mathbf{x}^T(n) \mathbf{x}(n)| < 1$, or

$$0 < \mu < \frac{2}{\|\mathbf{x}(n)\|^2} \quad (\text{stability region})$$

Observations - LMS Algorithm

- LMS is a noisy approximation of the steepest descent algorithm
- The gradient estimate is unbiased
- The errors in the gradient estimate lead to $J_{mse}(\infty) \neq 0$
- Vector $\mathbf{w}(n)$ is now random
- Steepest descent properties are no longer guaranteed
 - ⇒ *LMS analysis required*
- The **instantaneous estimates** allow tracking without redesign

Some Research Results

- J. C. M. Bermudez and N. J. Bershad, "A nonlinear analytical model for the quantized LMS algorithm - the arbitrary step size case," IEEE Transactions on Signal Processing, vol.44, No. 5, pp. 1175-1183, May 1996.
- J. C. M. Bermudez and N. J. Bershad, "Transient and tracking performance analysis of the quantized LMS algorithm for time-varying system identification," IEEE Transactions on Signal Processing, vol.44, No. 8, pp. 1990-1997, August 1996.
- N. J. Bershad and J. C. M. Bermudez, "A nonlinear analytical model for the quantized LMS algorithm - the power-of-two step size case," IEEE Transactions on Signal Processing, vol.44, No. 11, pp. 2895- 2900, November 1996.
- N. J. Bershad and J. C. M. Bermudez, "Sinusoidal interference rejection analysis of an LMS adaptive feedforward controller with a noisy periodic reference," IEEE Transactions on Signal Processing, vol.46, No. 5, pp. 1298-1313, May 1998.
- J. C. M. Bermudez and N. J. Bershad, "Non-Wiener behavior of the Filtered-X LMS algorithm," IEEE Trans. on Circuits and Systems II - Analog and Digital Signal Processing, vol.46, No. 8, pp. 1110-1114, Aug 1999.

Some Research Results – continued

- O. J. Tobias, J. C. M. Bermudez and N. J. Bershad, “Mean weight behavior of the Filtered-X LMS algorithm,” IEEE Transactions on Signal Processing, vol. 48, No. 4, pp. 1061-1075, April 2000.
- M. H. Costa, J. C. M. Bermudez and N. J. Bershad, “Stochastic analysis of the LMS algorithm with a saturation nonlinearity following the adaptive filter output,” IEEE Transactions on Signal Processing, vol. 49, No. 7, pp. 1370-1387, July 2001.
- M. H. Costa, J. C. M. Bermudez and N. J. Bershad, “Stochastic analysis of the Filtered-X LMS algorithm in systems with nonlinear secondary paths,” IEEE Transactions on Signal Processing, vol. 50, No. 6, pp. 1327-1342, June 2002.
- M. H. Costa, J. C. M. Bermudez and N. J. Bershad, “The performance surface in filtered nonlinear mean square estimation,” IEEE Transactions on Circuits and Systems I, vol. 50, No. 3, p. 445-447, March 2003.
- G. Barrault, J. C. M. Bermudez and A. Lenzi, “New Analytical Model for the Filtered-x Least Mean Squares Algorithm Verified Through Active Noise Control Experiment,” Mechanical Systems and Signal Processing, v. 21, p. 1839-1852, 2007.

Some Research Results – continued

- N. J. Bershad, J. C. M. Bermudez and J. Y. Tournet, “Stochastic Analysis of the LMS Algorithm for System Identification with Subspace Inputs,” IEEE Trans. on Signal Process., v. 56, p. 1018-1027, 2008.
- J. C. M. Bermudez, N. J. Bershad and J. Y. Tournet, “An Affine Combination of Two LMS Adaptive Filters Transient Mean-Square Analysis,” IEEE Trans. on Signal Process., v. 56, p. 1853-1864, 2008.
- M. H. Costa, L. R. Ximenes and J. C. M. Bermudez, “Statistical Analysis of the LMS Adaptive Algorithm Subjected to a Symmetric Dead-Zone Nonlinearity at the Adaptive Filter Output,” Signal Processing, v. 88, p. 1485-1495, 2008.
- M. H. Costa and J. C. M. Bermudez, “A Noise Resilient Variable Step-Size LMS Algorithm,” Signal Processing, v. 88, no. 3, p. 733-748, March 2008.
- P. Honeine, C. Richard, J. C. M. Bermudez, J. Chen and H. Snoussi, “A Decentralized Approach for Nonlinear Prediction of Time Series Data in Sensor Networks,” EURASIP Journal on Wireless Communications and Networking, v. 2010, p. 1-13, 2010. (KNLMS)

The Normalized LMS Algorithm – NLMS

- Most Employed adaptive algorithm in real-time applications
- Like LMS has different interpretations (even more)
- Alleviates a drawback of the LMS algorithm

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu e(n) \mathbf{x}(n)$$

- ▶ If amplitude of $\mathbf{x}(n)$ is large \Leftrightarrow Gradient noise amplification
- ▶ Sub-optimal performance when σ_x^2 varies with time (for instance, speech)

NLMS – A Solution to a Local Optimization Problem

- Error expressions

$$e(n) = d(n) - \mathbf{x}^T(n)\mathbf{w}(n) \quad (\text{a priori error})$$

$$\epsilon(n) = d(n) - \mathbf{x}^T(n)\mathbf{w}(n+1) \quad (\text{a posteriori error})$$

- We want to maximize $|\epsilon(n) - e(n)|$ with $|\epsilon(n)| < |e(n)|$

$$\epsilon(n) - e(n) = -\mathbf{x}^T(n)\Delta\mathbf{w}(n) \quad (\text{A})$$

$$\Delta\mathbf{w}(n) = \mathbf{w}(n+1) - \mathbf{w}(n)$$

- For $|\epsilon(n)| < |e(n)|$ we impose the restriction

$$\epsilon(n) = (1 - \mu)e(n), \quad |1 - \mu| < 1$$

$$\Leftrightarrow \epsilon(n) - e(n) = -\mu e(n) \quad (\text{B})$$

- For $\max |\epsilon(n) - e(n)| \Leftrightarrow \Delta\mathbf{w}(n)$ in the direction of $\mathbf{x}(n)$

$$\Leftrightarrow \Delta\mathbf{w}(n) = k\mathbf{x}(n) \quad (\text{C})$$

- Using (A), (B) and (C),

$$k = \mu \frac{e(n)}{\mathbf{x}^T(n)\mathbf{x}(n)}$$

and

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu \frac{e(n)\mathbf{x}(n)}{\mathbf{x}^T(n)\mathbf{x}(n)}$$

(NLMS weight update)

NLMS – Solution to a Constrained Optimization Problem

- Error sequence

$$y(n) = \mathbf{x}^T(n) \mathbf{w}(n) \quad (\text{estimate of } d(n))$$

$$e(n) = d(n) - y(n) = d(n) - \mathbf{x}^T(n) \mathbf{w}(n) \quad (\text{estimation error})$$

- Optimization (principle of minimal disturbance)

$$\begin{aligned} &\text{Minimize} \quad \|\Delta \mathbf{w}(n)\|^2 = \|\mathbf{w}(n+1) - \mathbf{w}(n)\|^2 \\ &\text{subject to:} \quad \mathbf{x}^T(n) \mathbf{w}(n+1) = d(n) \end{aligned}$$

This problem can be solved using the method of Lagrange multipliers

- Using Lagrange multipliers, we minimize

$$f[\mathbf{w}(n+1)] = \|\mathbf{w}(n+1) - \mathbf{w}(n)\|^2 + \lambda[d(n) - \mathbf{x}^T(n)\mathbf{w}(n+1)]$$

- Differentiating w.r.t. $\mathbf{w}(n+1)$ and equating the result to zero,

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \frac{1}{2}\lambda\mathbf{x}(n) \quad (*)$$

Using this result in $\mathbf{x}^T(n)\mathbf{w}(n+1) = d(n)$ yields

$$\lambda = \frac{2e(n)}{\mathbf{x}^T(n)\mathbf{x}(n)}$$

- Using this result in (*) yields

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu \frac{e(n)\mathbf{x}(n)}{\mathbf{x}^T(n)\mathbf{x}(n)}$$

NLMS as an Orthogonalization Process

- Conditions for the analysis

- ▶ $e(n) = d(n) - \mathbf{x}^T(n)\mathbf{w}(n)$
- ▶ \mathbf{w}^o is the optimal solution (Wiener solution)
- ▶ $d(n) = \mathbf{x}^T(n)\mathbf{w}^o$ (no noise)
- ▶ $\mathbf{v}(n) = \mathbf{w}(n) - \mathbf{w}^o$ (weight error vector)

- Error signal

$$\begin{aligned}e(n) &= d(n) - \mathbf{x}^T(n)[\mathbf{v}(n) + \mathbf{w}^o] \\&= \mathbf{x}^T(n)\mathbf{w}^o - \mathbf{x}^T(n)[\mathbf{v}(n) + \mathbf{w}^o] \\&= -\mathbf{x}^T(n)\mathbf{v}(n)\end{aligned}$$

$$e(n) = -\mathbf{x}^T(n)\mathbf{v}(n)$$

- *Interpretation:* To minimize the error $\mathbf{v}(n)$ should be orthogonal to all input vectors
Restriction: We have only one vector $\mathbf{x}(n)$
- *Iterative solution:*
 We can subtract from $\mathbf{v}(n)$ its component in the direction of $\mathbf{x}(n)$ at each iteration
- If there are N adaptive coefficients, $\mathbf{v}(n)$ could be reduced to zero after N orthogonal input vectors
- Iterative projection extraction
 ⇨ Gram-Schmidt orthogonalization

- Recursive orthogonalization

$$n = 0 : \quad \mathbf{v}(1) = \mathbf{v}(0) - \text{proj. of } \mathbf{v}(0) \text{ onto } \mathbf{x}(0)$$

$$n = 1 : \quad \mathbf{v}(2) = \mathbf{v}(1) - \text{proj. of } \mathbf{v}(1) \text{ onto } \mathbf{x}(1)$$

$$\vdots \quad \vdots$$

$$n + 1 : \quad \mathbf{v}(n + 1) = \mathbf{v}(n) - \text{proj. of } \mathbf{v}(n) \text{ onto } \mathbf{x}(n)$$

- Projection of $\mathbf{v}(n)$ onto $\mathbf{x}(n)$

$$\mathcal{P}_{\mathbf{x}(n)}[\mathbf{v}(n)] = \{ \mathbf{x}(n)[\mathbf{x}^T(n)\mathbf{x}(n)]^{-1}\mathbf{x}^T(n) \} \mathbf{v}(n)$$

- Weight update equation

$$\begin{aligned} \mathbf{v}(n + 1) &= \mathbf{v}(n) - \underbrace{\mu}_{\text{scalar}} \underbrace{[\mathbf{x}^T(n)\mathbf{x}(n)]^{-1}\mathbf{x}^T(n)\mathbf{v}(n)}_{-e(n)} \mathbf{x}(n) \\ &= \mathbf{v}(n) + \mu \frac{e(n)\mathbf{x}(n)}{\mathbf{x}^T(n)\mathbf{x}(n)} \end{aligned}$$

NLMS has its own problems!

- NLMS solves the LMS gradient error amplification problem, but ...
- What happens if $\|\mathbf{x}(n)\|^2$ gets too small?
- One needs to add some regularization

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu \frac{e(n)\mathbf{x}(n)}{\mathbf{x}^T(n)\mathbf{x}(n) + \epsilon} \quad (\epsilon\text{-NLMS})$$

ε -NLMS – Stochastic Approximation of Regularized Newton

- Regularized Newton Algorithm

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu[\varepsilon\mathbf{I} + \mathbf{R}_{xx}]^{-1}[\mathbf{p} - \mathbf{R}_{xx}\mathbf{w}(n)]$$

- Instantaneous estimates

$$\begin{aligned}\hat{\mathbf{p}} &= \mathbf{x}(n)d(n) \\ \hat{\mathbf{R}}_{xx} &= \mathbf{x}(n)\mathbf{x}^T(n)\end{aligned}$$

- Using these estimates

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu[\varepsilon\mathbf{I} + \mathbf{x}(n)\mathbf{x}^T(n)]^{-1}\mathbf{x}(n) \overbrace{[d(n) - \mathbf{x}^T(n)\mathbf{w}(n)]}^{e(n)}$$

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu[\varepsilon\mathbf{I} + \mathbf{x}(n)\mathbf{x}^T(n)]^{-1}\mathbf{x}(n)e(n)$$

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu[\varepsilon \mathbf{I} + \mathbf{x}(n)\mathbf{x}^T(n)]^{-1} \mathbf{x}(n)e(n)$$

- Inversion of $\varepsilon \mathbf{I} + \mathbf{x}(n)\mathbf{x}^T(n)$?
- Matrix Inversion Formula

$$[\mathbf{A} + \mathbf{BCD}]^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}[\mathbf{C}^{-1} + \mathbf{DA}^{-1}\mathbf{B}]^{-1}\mathbf{DA}^{-1}$$

- Thus,

$$[\varepsilon \mathbf{I} + \mathbf{x}(n)\mathbf{x}^T(n)]^{-1} = \varepsilon^{-1} \mathbf{I} - \frac{\varepsilon^{-2}}{1 + \varepsilon^{-1} \mathbf{x}^T(n)\mathbf{x}(n)} \mathbf{x}(n)\mathbf{x}^T(n)$$

- Post-multiplying both sides by $\mathbf{x}(n)$ and rearranging

$$[\varepsilon \mathbf{I} + \mathbf{x}(n)\mathbf{x}^T(n)]^{-1} \mathbf{x}(n) = \frac{\mathbf{x}(n)}{\varepsilon + \mathbf{x}^T(n)\mathbf{x}(n)}$$

and

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu \frac{e(n)\mathbf{x}(n)}{\mathbf{x}^T(n)\mathbf{x}(n) + \varepsilon}$$

Some Research Results

- M. H. Costa and J. C. M. Bermudez, "An improved model for the normalized LMS algorithm with Gaussian inputs and large number of coefficients," in Proc. of the 2002 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP-2002), Orlando, Florida, pp. II- 1385-1388, May 13-17, 2002.
- J. C. M. Bermudez and M. H. Costa, "A Statistical Analysis of the ε -NLMS and NLMS Algorithms for Correlated Gaussian Signals," Journal of the Brazilian Telecommunications Society, v. 20, n. 2, p. 7-13, 2005.
- G. Barrault, M. H. Costa, J. C. M. Bermudez and A. Lenzi, "A new analytical model for the NLMS algorithm," in Proc. 2005 IEEE International Conference on Acoustics Speech and Signal Processing (ICASSP 2005) Pennsylvania, USA, vol. IV, p. 41-44, 2005.
- J. C. M. Bermudez, N. J. Bershad and J.-Y. Tournet, "An Affine Combination of Two NLMS Adaptive Filters - Transient Mean-Square Analysis," In Proc. Forty-Second Asilomar Conference on Asilomar Conference on Signals, Systems & Computers, Pacific Grove, CA, USA, 2008.
- P. Honeine, C. Richard, J. C. M. Bermudez and H. Snoussi, "Distributed prediction of time series data with kernels and adaptive filtering techniques in sensor networks," In Proc. Forty-Second Asilomar Conference on Asilomar Conference on Signals, Systems & Computers, Pacific Grove, CA, USA, 2008.

The Affine Projection Algorithm

- Consider the NLMS algorithm but using $\{\mathbf{x}(n), \mathbf{x}(n-1), \dots, \mathbf{x}(n-P)\}$

$$\begin{array}{ll} \text{Minimize} & \|\Delta \mathbf{w}(n)\|^2 = \|\mathbf{w}(n+1) - \mathbf{w}(n)\|^2 \\ \text{subject to:} & \left\{ \begin{array}{l} \mathbf{x}^T(n) \mathbf{w}(n+1) = d(n) \\ \mathbf{x}^T(n-1) \mathbf{w}(n+1) = d(n-1) \\ \vdots \\ \mathbf{x}^T(n-P) \mathbf{w}(n+1) = d(n-P) \end{array} \right. \end{array}$$

- Observation matrix

$$\mathbf{X}(n) = [\mathbf{x}(n), \mathbf{x}(n-1), \dots, \mathbf{x}(n-P)]$$

- Desired signal vector

$$\mathbf{d}(n) = [d(n), d(n-1), \dots, d(n-P)]$$

- Error vector

$$\mathbf{e}(n) = [e(n), e(n-1), \dots, e(n-P)] = \mathbf{d}(n) - \mathbf{X}^T(n)\mathbf{w}(n)$$

- Vector of the constraint errors

$$\mathbf{e}_c(n) = \mathbf{d}(n) - \mathbf{X}^T(n)\mathbf{w}(n+1)$$

- Using Lagrange multipliers, we minimize

$$f[\mathbf{w}(n+1)] = \|\mathbf{w}(n+1) - \mathbf{w}(n)\|^2 + \boldsymbol{\lambda}^T[\mathbf{d}(n) - \mathbf{X}^T(n)\mathbf{w}(n+1)]$$

- Differentiating w.r.t. $\mathbf{w}(n+1)$ and equating the result to zero,

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \frac{1}{2}\mathbf{X}(n)\boldsymbol{\lambda} \quad (*)$$

Using this result in $\mathbf{X}^T(n)\mathbf{w}(n+1) = \mathbf{d}(n)$ yields

$$\boldsymbol{\lambda} = 2[\mathbf{X}^T(n)\mathbf{X}(n)]^{-1}\mathbf{e}(n)$$

- Using this result in (*) yields

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu\mathbf{X}(n)[\mathbf{X}^T(n)\mathbf{X}(n)]^{-1}\mathbf{e}(n)$$

Affine Projection

Solution to the Underdetermined Least-Squares Problem

- We want minimize

$$\|e_c(n)\|^2 = \|\mathbf{d}(n) - \mathbf{X}^T(n)\mathbf{w}(n+1)\|^2$$

where $\mathbf{X}^T(n)$ is $(P+1) \times N$ with $(P+1) < N$

- Thus, we look for the least-squares solution of the underdetermined system

$$\mathbf{X}^T(n)\mathbf{w}(n+1) = \mathbf{d}(n)$$

- The solution is

$$\mathbf{w}(n+1) = [\mathbf{X}^T(n)]^+ \mathbf{d}(n) = \mathbf{X}(n)[\mathbf{X}^T(n)\mathbf{X}(n)]^{-1}\mathbf{d}(n)$$

- Using $\mathbf{d}(n) = \mathbf{e}(n) + \mathbf{X}^T(n)\mathbf{w}(n)$ yields

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu \mathbf{X}(n)[\mathbf{X}^T(n)\mathbf{X}(n)]^{-1}\mathbf{e}(n)$$

Observations:

- Order of the AP algorithm: $P + 1$ ($P = 0 \rightarrow$ NLMS)
- Convergence speed increases with P (but not linearly)
- Computational complexity increases with P (not linearly)
- If $\mathbf{X}^T(n)\mathbf{X}(n)$ is close to singular, use $\mathbf{X}^T(n)\mathbf{X}(n) + \varepsilon\mathbf{I}$
- The scalar error of NLMS becomes a vector error in AP (except for $\mu = 1$)

Affine Projection – Stochastic Approximation of Regularized Newton

- Regularized Newton Algorithm

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu[\tilde{\varepsilon}\mathbf{I} + \mathbf{R}_{xx}]^{-1}[\mathbf{p} - \mathbf{R}_{xx}\mathbf{w}(n)]$$

- Estimates using time window statistics

$$\hat{\mathbf{p}} = \frac{1}{P+1} \sum_{k=n-P}^n \mathbf{X}(k)d(k) = \frac{1}{P+1} \mathbf{X}(n)\mathbf{d}(n)$$

$$\hat{\mathbf{R}}_{xx} = \frac{1}{P+1} \sum_{k=n-P}^n \mathbf{X}(k)\mathbf{X}^T(k) = \frac{1}{P+1} \mathbf{X}(n)\mathbf{X}^T(n)$$

- Using these estimates with $\tilde{\varepsilon} = \varepsilon/(P+1)$,

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu[\varepsilon\mathbf{I} + \mathbf{X}(n)\mathbf{X}^T(n)]^{-1} \mathbf{X}(n) \overbrace{[\mathbf{d}(n) - \mathbf{X}^T(n)\mathbf{w}(n)]}^{e(n)}$$

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu[\varepsilon\mathbf{I} + \mathbf{X}(n)\mathbf{X}^T(n)]^{-1} \mathbf{X}(n)\mathbf{e}(n)$$

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu[\varepsilon \mathbf{I} + \mathbf{X}(n)\mathbf{X}^T(n)]^{-1} \mathbf{X}(n)\mathbf{e}(n)$$

- Using the Matrix Inversion Formula

$$[\varepsilon \mathbf{I} + \mathbf{X}(n)\mathbf{X}^T(n)]^{-1} = \mathbf{X}(n)[\varepsilon \mathbf{I} + \mathbf{X}^T(n)\mathbf{X}(n)]^{-1}$$

and

$$\boxed{\mathbf{w}(n+1) = \mathbf{w}(n) + \mu \mathbf{X}(n)[\varepsilon \mathbf{I} + \mathbf{X}^T(n)\mathbf{X}(n)]^{-1} \mathbf{e}(n)}$$

Affine Projection Algorithm as a Projection onto an Affine Subspace

- Conditions for the analysis
 - ▶ $\mathbf{e}(n) = \mathbf{d}(n) - \mathbf{X}^T(n)\mathbf{w}(n)$
 - ▶ $\mathbf{d}(n) = \mathbf{X}^T(n)\mathbf{w}^o$ defines the optimal solution in the least-squares sense
 - ▶ $\mathbf{v}(n) = \mathbf{w}(n) - \mathbf{w}^o$ (weight error vector)
- Error vector

$$\begin{aligned}\mathbf{e}(n) &= \mathbf{d}(n) - \mathbf{X}^T(n)[\mathbf{v}(n) + \mathbf{w}(n+1)] \\ &= \mathbf{X}^T(n)\mathbf{w}(n+1) - \mathbf{X}^T(n)[\mathbf{v}(n) + \mathbf{w}(n+1)] \\ &= -\mathbf{X}^T(n)\mathbf{v}(n)\end{aligned}$$

$$\mathbf{e}(n) = -\mathbf{X}^T(n)\mathbf{v}(n)$$

- *Interpretation:* To minimize the error $\mathbf{v}(n)$ should be orthogonal to all input vectors

Restriction: We are going to use only $\{\mathbf{x}(n), \mathbf{x}(n-1), \dots, \mathbf{x}(n-P)\}$

- *Iterative solution:*

We can subtract from $\mathbf{v}(n)$ its projection onto the range of $\mathbf{X}(n)$ at each iteration

$$\mathbf{v}(n+1) = \mathbf{v}(n) - \mathbf{P}_{\mathbf{X}(n)}\mathbf{v}(n) \quad (\text{Proj. onto an affine subspace})$$

- Using $\mathbf{X}^T(n)\mathbf{v}(n) = -\mathbf{e}(n)$

$$\begin{aligned} \mathbf{v}(n+1) &= \mathbf{v}(n) - \{\mathbf{X}(n)[\mathbf{X}^T(n)\mathbf{X}(n)]^{-1}\mathbf{X}^T(n)\} \mathbf{v}(n) \\ &= \mathbf{v}(n) + \mu \mathbf{X}(n)[\mathbf{X}^T(n)\mathbf{X}(n)]^{-1}\mathbf{e}(n) \quad (\text{AP algorithm}) \end{aligned}$$

Pseudo Affine Projection Algorithm

- Major problem with the AP algorithm

For $\mu \neq 1$, $e(n) \rightarrow e(n) \Rightarrow$ large computational complexity

- The Pseudo-AP algorithm replaces the input $x(n)$ with its P -th order autoregressive prediction

$$\hat{\mathbf{x}}(n) = \sum_{k=1}^P a_k \mathbf{x}(n-k)$$

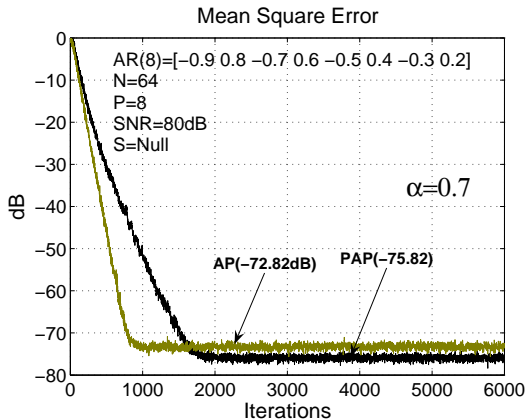
- Using the least-squares solution for $\mathbf{a} = [a_1, a_2, \dots, a_P]^T$,

$$\mathbf{a} = [\mathbf{X}_p^T(n) \mathbf{X}_p(n)]^{-1} \mathbf{X}_p^T(n) \mathbf{x}(n), \quad \mathbf{X}_p(n) = [\mathbf{x}(n-1), \dots, \mathbf{x}(n-P)]$$

- Now, we subtract from $\mathbf{x}(n)$ its projections onto the last P input vectors

$$\begin{aligned} \phi(n) &= \mathbf{x}(n) - \mathbf{X}_p(n) \mathbf{a}(n) \\ &= \mathbf{x}(n) - \{ \mathbf{X}_p(n) [\mathbf{X}_p^T(n) \mathbf{X}_p(n)]^{-1} \mathbf{X}_p^T(n) \} \mathbf{x}(n) \\ &= (\mathbf{I} - \mathbf{P}_P) \mathbf{x}(n); \quad (\mathbf{P}_P: \text{projection matrix}) \end{aligned}$$

Example – Pseudo-AP *versus* AP



- Using ϕ as the new input vector for the NLMS algorithm,

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu \frac{\phi(n)}{\phi^T(n)\phi(n)} e(n)$$

- It can be shown that Pseudo-AP is identical to AP for an AR input and $\mu = 1$
- Otherwise, Pseudo-AP is different from AP
- Simpler to implement than AP for $\mu \neq 1$
- Can lead even to better steady-state results than AP for AR inputs
- For AR inputs: NLMS with input “orthogonalization”

Some Research Results

- S. J. M. de Almeida, J. C. M. Bermudez, N. J. Bershad and M. H. Costa, "A statistical analysis of the affine projection algorithm for unity step size and autoregressive inputs," IEEE Transactions on Circuits and Systems - I, vol. 52, pp. 1394-1405, July 2005.
- S. J. M. Almeida, J. C. M. Bermudez and N. J. Bershad, "A Stochastic Model for a Pseudo Affine Projection Algorithm," IEEE Transactions on Signal Processing, v. 57, p. 107-118, 2009.
- S. J. M. Almeida, M. H. Costa and J. C. M. Bermudez, "A Stochastic Model for the Deficient Length Pseudo Affine Projection Adaptive Algorithm," In Proc. 17th European Signal Processing Conference (EUSIPCO), Aug. 24-28, Glasgow, Scotland, 2009.
- S. J. M. Almeida, M. H. Costa and J. C. M. Bermudez, "A stochastic model for the deficient order affine projection algorithm," in Proc. ISSPA 2010, Kuala Lumpur, Malaysia, May 2010.
- C. Richard, J. C. M. Bermudez and P. Honeine, "Online Prediction of Time Series Data With Kernels," IEEE Transactions on Signal Processing, v. 57, p. 1058-1067, 2009.
- P. Honeine, C. Richard, J. C. M. Bermudez, H. Snoussi, M. Essoloh and F. Vincent, "Functional estimation in Hilbert space for distributed learning in wireless sensor networks," In Proc. IEEE International Conference on Acoustics Speech and Signal Processing (ICASSP), Taipei, Taiwan, p. 2861-2864, 2009.

Deterministic Algorithms

Recursive Least Squares Algorithm (RLS)

- Based on a deterministic philosophy
- Designed for the present realization of the input signals
- Least squares method adapted for real time processing of temporal series
- Convergence speed is not strongly dependent on the input statistics

RLS – Problem Definition

- Define

$$\mathbf{x}_o(n) = [x(n), x(n-1), \dots, x(0)]^T$$

$$\mathbf{x}(n) = [x(n), x(n-1), \dots, x(n-N+1)]^T \quad \text{input to the } N\text{-tap filter}$$

- Desired signal $d(n)$
- Estimator of $d(n)$: $y(n) = \mathbf{x}^T(n)\mathbf{w}(n)$
- Cost function – Squared Error

$$\begin{aligned} J_{ls}(n) &= \sum_{k=0}^n e^2(k) = \sum_{k=0}^n [d(\textcolor{red}{k}) - \mathbf{x}^T(\textcolor{red}{k})\mathbf{w}(\textcolor{red}{n})]^2 \\ &= \mathbf{e}^T(n)\mathbf{e}(n), \quad \mathbf{e}(n) = [e(n), e(n-1), \dots, e(0)]^T \end{aligned}$$

In vector form

$$\mathbf{d}(n) = [d(n), d(n-1), \dots, d(0)]^T$$

$$\mathbf{y}(n) = [y(n), y(n-1), \dots, y(0)]^T, \quad y(k) = \mathbf{x}^T(k) \mathbf{w}(n)$$

$$\mathbf{y}(n) = \sum_{k=0}^{N-1} w_k(n) \mathbf{x}_o(n-k) = \mathbf{\Theta}(n) \mathbf{w}(n)$$

$$\begin{aligned} \mathbf{\Theta}(n) &= \begin{bmatrix} x(n) & x(n-1) & \cdots & x(n-N+1) \\ x(n-1) & x(n-2) & \cdots & x(n-N) \\ \vdots & \vdots & \vdots & \vdots \\ x(0) & x(-1) & \cdots & x(-N+1) \end{bmatrix} \\ &= [\mathbf{x}_o(n), \mathbf{x}_o(n-1), \dots, \mathbf{x}_o(n-N+1)] \end{aligned}$$

$$\boxed{\mathbf{e}(n) = \mathbf{d}(n) - \mathbf{\Theta}(n) \mathbf{w}(n)}$$

$$\mathbf{e}(n) = \mathbf{d}(n) - \mathbf{\Theta}(n)\mathbf{w}(n)$$

and we minimize

$$J_{ls}(n) = \|\mathbf{e}(n)\|^2$$

- Normal Equations

$$\mathbf{\Theta}^T(n)\mathbf{\Theta}(n)\mathbf{w}(n) = \mathbf{\Theta}^T(n)\mathbf{d}(n)$$

$$\tilde{\mathbf{R}}(n)\mathbf{w}(n) = \tilde{\mathbf{p}}(n)$$

- Alternative representation for $\tilde{\mathbf{R}}$ and $\tilde{\mathbf{p}}$

$$\tilde{\mathbf{R}}(n) = \mathbf{\Theta}^T(n)\mathbf{\Theta}(n) = \sum_{k=0}^n \mathbf{x}(k)\mathbf{x}^T(k)$$

$$\tilde{\mathbf{p}}(n) = \mathbf{\Theta}^T(n)\mathbf{d}(n) = \sum_{k=0}^n \mathbf{x}(k)d(k)$$

RLS – Observations

- $\mathbf{w}(n)$ remains fixed for $0 \leq k \leq n$ to determine $J_{ls}(n)$
- Optimum vector $\mathbf{w}(n)$ is $\mathbf{w}^o(n) = \tilde{\mathbf{R}}^{-1}(n)\tilde{\mathbf{p}}$
- The algorithm has growing memory
- In an adaptive implementation,

$$\mathbf{w}^o(n) = \tilde{\mathbf{R}}^{-1}(n)\tilde{\mathbf{p}}$$

must be solved for each iteration n

- Problems for an adaptive implementation
 - ▶ Difficulties with nonstationary signals (infinite data window)
 - ▶ Computational complexity

RLS - Forgetting the Past

- Modified cost function

$$\begin{aligned} J_{ls}(n) &= \sum_{k=0}^n \lambda^{n-k} e^2(k) = \sum_{k=0}^n \lambda^{n-k} [d(\textcolor{red}{k}) - \mathbf{x}^T(\textcolor{red}{k}) \mathbf{w}(\textcolor{red}{n})]^2 \\ &= \mathbf{e}^T(n) \mathbf{\Lambda} \mathbf{e}(n), \quad \mathbf{\Lambda} = \text{diag}[1, \lambda, \lambda^2, \dots, \lambda^n] \end{aligned}$$

- Modified Normal Equations

$$\mathbf{\Theta}^T(n) \mathbf{\Lambda} \mathbf{\Theta}(n) \mathbf{w}(n) = \mathbf{\Theta}^T(n) \mathbf{\Lambda} \mathbf{d}(n)$$

$$\hat{\mathbf{R}}(n) \mathbf{w}(n) = \hat{\mathbf{p}}(n)$$

where

$$\hat{\mathbf{R}}(n) = \mathbf{\Theta}^T(n) \mathbf{\Lambda} \mathbf{\Theta}(n) = \sum_{k=0}^n \lambda^{n-k} \mathbf{x}(k) \mathbf{x}^T(k)$$

$$\hat{\mathbf{p}}(n) = \mathbf{\Theta}^T(n) \mathbf{\Lambda} \mathbf{d}(n) = \sum_{k=0}^n \lambda^{n-k} \mathbf{x}(k) d(k)$$

RLS – Recursive Updating

- Correlation Matrix

$$\begin{aligned}\hat{\mathbf{R}}(n) &= \sum_{k=0}^n \lambda^{n-k} \mathbf{x}(k) \mathbf{x}^T(k) \\ &= \sum_{k=0}^{n-1} \lambda^{n-k} \mathbf{x}(k) \mathbf{x}^T(k) + \mathbf{x}(n) \mathbf{x}^T(n) \\ &= \lambda \hat{\mathbf{R}}(n-1) + \mathbf{x}(n) \mathbf{x}^T(n)\end{aligned}$$

- Cross-correlation vector

$$\begin{aligned}\hat{\mathbf{p}}(n) &= \sum_{k=0}^n \lambda^{n-k} \mathbf{x}(k) d(k) \\ &= \sum_{k=0}^{n-1} \lambda^{n-k} \mathbf{x}(k) d(k) + \mathbf{x}(n) d(n) \\ &= \lambda \hat{\mathbf{p}}(n-1) + \mathbf{x}(n) d(n)\end{aligned}$$

- We have recursive updating expressions for $\hat{\mathbf{R}}(n)$ and $\hat{\mathbf{p}}(n)$
- However, we need a recursive updating for $\hat{\mathbf{R}}^{-1}(n)$, as

$$\mathbf{w}^o(n) = \hat{\mathbf{R}}^{-1}(n)\hat{\mathbf{p}}(n)$$

- Applying the matrix inversion lemma to

$$\hat{\mathbf{R}}(n) = \lambda \hat{\mathbf{R}}(n-1) + \mathbf{x}(n)\mathbf{x}^T(n)$$

and defining $\mathbf{P}(n) = \hat{\mathbf{R}}^{-1}(n)$ yields

$$\mathbf{P}(n) = \lambda^{-1}\mathbf{P}(n-1) - \frac{\lambda^{-2}\mathbf{P}(n-1)\mathbf{x}(n)\mathbf{x}^T(n)\mathbf{P}(n-1)}{1 + \lambda^{-1}\mathbf{x}^T(n)\mathbf{P}(n-1)\mathbf{x}(n)}$$

The Gain Vector

- Definition

$$\mathbf{k}(n) = \frac{\lambda^{-1} \mathbf{P}(n-1) \mathbf{x}(n)}{1 + \lambda^{-1} \mathbf{x}^T(n) \mathbf{P}(n-1) \mathbf{x}(n)}$$

- Using this definition

$$\mathbf{P}(n) = \lambda^{-1} \mathbf{P}(n-1) - \lambda^{-1} \mathbf{k}(n) \mathbf{x}^T(n) \mathbf{P}(n-1)$$

- $\mathbf{k}(n)$ can be written as

$$\begin{aligned} \mathbf{k}(n) &= \overbrace{[\lambda^{-1} \mathbf{P}(n-1) - \lambda^{-1} \mathbf{k}(n) \mathbf{x}^T(n) \mathbf{P}(n-1)]}^{\mathbf{P}(n)} \mathbf{x}(n) \\ &= \hat{\mathbf{R}}^{-1}(n) \mathbf{x}(n) \end{aligned}$$

$\Leftrightarrow \mathbf{k}(n)$ is the repres. of $\mathbf{x}(n)$ on the column space of $\hat{\mathbf{R}}(n)$

RLS – Recursive Weight Update

- $\mathbf{w}(n+1) = \mathbf{w}^o(n)$
- Using

$$\mathbf{w}(n+1) = \hat{\mathbf{R}}^{-1}(n)\hat{\mathbf{p}}(n) = \mathbf{P}(n)\hat{\mathbf{p}}(n)$$

$$\hat{\mathbf{p}}(n) = \lambda\hat{\mathbf{p}}(n-1) + \mathbf{x}(n)d(n)$$

$$\mathbf{P}(n) = \lambda^{-1}\mathbf{P}(n-1) - \lambda^{-1}\mathbf{k}(n)\mathbf{x}^T(n)\mathbf{P}(n-1)$$

$$\mathbf{k}(n) = \mathbf{P}(n)\mathbf{x}(n)$$

$$e(n) = d(n) - \mathbf{x}^T(n)\mathbf{w}(n)$$

yields

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mathbf{k}(n)e(n)$$

RLS weight update

RLS – Observations

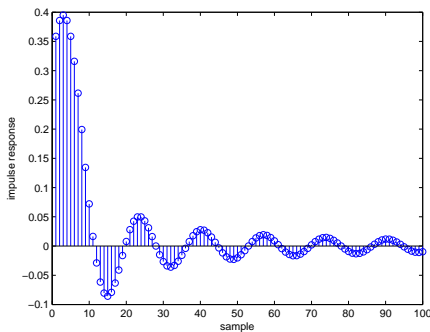
- The RLS convergence speed is not affected by the eigenvalues of $\hat{\mathbf{R}}(n)$
- Initialization: $\hat{\mathbf{p}}(0) = \mathbf{0}$, $\hat{\mathbf{R}}(0) = \delta \mathbf{I}$, $\delta \simeq (1 - \lambda)\sigma_x^2$
- This results in $\hat{\mathbf{R}}(n)$ changed to $\lambda^n \delta \mathbf{I} + \hat{\mathbf{R}}(n)$
 - ⇒ Biased estimate of $\mathbf{w}^o(n)$ for small n
(No problem for $\lambda < 1$ and large n)
- Numerical problems in finite precision
 - ▶ Unstable for $\lambda = 1$
 - ▶ Loss of symmetry in $\mathbf{P}(n)$
 - ★ Evaluate only upper (lower) part and diagonal of $\mathbf{P}(n)$
 - ★ Replace $\mathbf{P}(n)$ with $[\mathbf{P}(n) + \mathbf{P}^T(n)]/2$ after updating from $\mathbf{P}(n-1)$
- Numerical problems when $\mathbf{x}(n) \rightarrow \mathbf{0}$ and $\lambda < 1$

Some Research Results

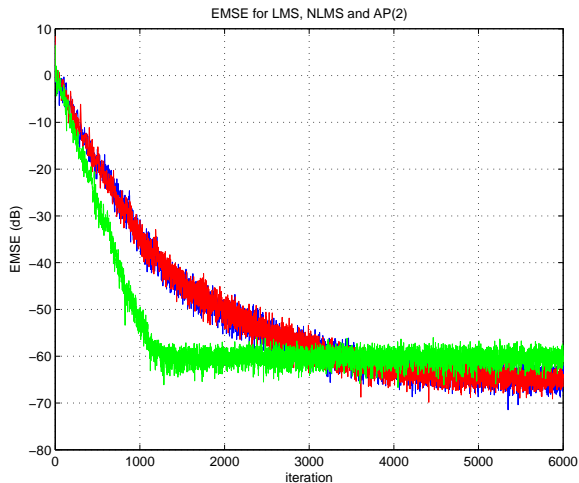
- C. Ludovico and J. C. M. Bermudez, "A recursive least squares algorithm robust to low- power excitation," in. Proc. 2004 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP-2004), Montreal, Canada, vol. II, p. 673-677, 2004.
- C. Ludovico and J. C. M. Bermudez, "An improved recursive least squares algorithm robust to input power variation," in Proc. IEEE Statistical Signal Processing Workshop, Bordeaux, France, 2005.

Performance Comparison

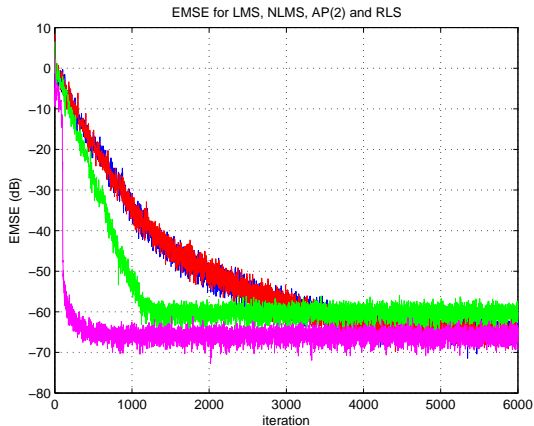
- System identification, $N=100$
- Input – AR(1) $x(n) = 0.9x(n-1) + v(n)$
- AP algorithm with $P = 2$
- Step sizes designed for equal LMS and NLMS performances
- Impulse response to be estimated



Excess Mean Square Error - LMS, NLMS and AP



Excess Mean Square Error - LMS, NLMS, AP and RLS



Performance Comparisons

Computational Complexity

- Adaptive filter with N real coefficients and real signals
- For the AP algorithm, $K = P + 1$

Algorithm	\times	$+$	$/$
LMS	$2N + 1$	$2N$	
NLMS	$3N + 1$	$3N$	1
AP	$(K^2 + 2K)N + K^3 + K$	$(K^2 + 2K)N + K^3 + K^2$	
RLS	$N^2 + 5N + 1$	$N^2 + 3N$	1

For $N = 100$, $P = 2$

Algorithm	\times	$+$	$/$	\simeq factor
LMS	201	200		1
NLMS	301	300	1	1.5
AP	1,530	1,536		7.5
RLS	10,501	10,300	1	52.5

Typical values for acoustic echo cancellation ($N = 1024$, $P = 2$)

Algorithm	\times	$+$	$/$	\simeq factor
LMS	2,049	2,048		1
NLMS	3,073	3,072	1	1.5
AP	15,390	15,396		7.5
RLS	1,053,697	1,051,648	1	514

How to Deal with Computational Complexity?

- Not an easy task!!!
 - There are "fast" versions for some algorithms (especially RLS)
 - What is usually not said is that ... speed can bring
 - ▶ Instability
 - ▶ Increased need for memory
- ⇒ *Most applications rely on simple solutions*

Understanding the Adaptive Filter Behavior

What are Adaptive Filters?

Adaptive filters are, by design, systems that are

- Time-variant ($\mathbf{w}(n)$ is time-variant)
- Nonlinear ($y(n)$ is a nonlinear function of $\mathbf{x}(n)$)
- Stochastic ($\mathbf{w}(n)$ is random)

LMS:

$$\begin{aligned}\mathbf{w}(n+1) &= \mathbf{w}(n) + \mu e(n) \mathbf{x}(n) \\ &= [\mathbf{I} - \mu \mathbf{x}(n) \mathbf{x}^T(n)] \mathbf{x}(n) + \mu d(n) \mathbf{x}(n) \\ y(n) &= \mathbf{x}^T(n) \mathbf{w}(n)\end{aligned}$$

- They are difficult to understand
- Different applications and signals require different analyses
- Simplifying assumptions are necessary

⇒ *Good design requires good analytical models.*

LMS Analysis

Basic equations:

$$d(n) = \mathbf{x}^T(n) \mathbf{w}^o + z(n)$$

$$e(n) = d(n) - \mathbf{x}^T(n) \mathbf{w}(n)$$

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu e(n) \mathbf{x}(n)$$

$$\mathbf{v}(n) = \mathbf{w}(n) - \mathbf{w}^o \quad (\text{study about the optimum weight vector})$$

Mean Weight Behavior:

Neglecting dependence between $\mathbf{v}(n)$ and $\mathbf{x}(n)\mathbf{x}^T(n)$,

$$E[\mathbf{v}(n+1)] = (\mathbf{I} - \mu \mathbf{R}_{xx}) E[\mathbf{v}(n)]$$

- Follows the steepest descent trajectory
- Convergence condition: $0 < \mu < 2/\lambda_{\max}$
- $\lim_{n \rightarrow \infty} E[\mathbf{w}(n)] = \mathbf{w}^o$

However: $\mathbf{w}(n)$ is random

⇒ We need to study its fluctuations about $E[\mathbf{w}(n)] \rightarrow \text{MSE}$

The Mean-Square Estimation Error

$$e(n) = e_o(n) - \mathbf{x}^T(n)\mathbf{v}(n), \quad e_o(n) = d(n) - \mathbf{x}^T(n)\mathbf{w}^o$$

- Square the error equation
- Take the expected value
- Neglect the statistical dep. between $\mathbf{v}(n)$ and $\mathbf{x}(n)\mathbf{x}^T(n)$

$$J_{ms}(n) = E[e_o^2(n)] + Tr\{\mathbf{R}_{xx}\mathbf{K}(n)\}, \quad \mathbf{K}(n) = E[\mathbf{v}(n)\mathbf{v}^T(n)]$$

- This expression is independent of the adaptive algorithm
- Effect of the alg. on $J_{ms}(n)$ is determined by $\mathbf{K}(n)$
- $J_{ms_{\min}} = E[e_o^2(n)]$
- $J_{ms_{\text{ex}}} = Tr\{\mathbf{R}_{xx}\mathbf{K}(n)\}$

The Behavior of $\mathbf{K}(n)$ for LMS

- Using the basic equations

$$\mathbf{v}(n+1) = [\mathbf{I} + \mu \mathbf{x}(n)\mathbf{x}^T(n)] \mathbf{v}(n) + \mu e_o(n)\mathbf{x}(n)$$

- Post-multiply by $\mathbf{v}^T(n+1)$
- Take the expected value assuming $\mathbf{v}(n)$ and $\mathbf{x}(n)\mathbf{x}^T(n)$ independent
- Evaluate the expectations
 - ▶ Higher order moments of $x(n)$ require input pdf
- Assuming Gaussian inputs

$$\begin{aligned}\mathbf{K}(n+1) = & \mathbf{K}(n) - \mu [\mathbf{R}_{xx}\mathbf{K}(n) + \mathbf{K}(n)\mathbf{R}_{xx}] \\ & + \mu^2 \{ \mathbf{R}_{xx} \text{Tr}[\mathbf{R}_{xx}\mathbf{K}(n)] + 2\mathbf{R}_{xx}\mathbf{K}(n)\mathbf{R}_{xx} \} \\ & + \mu^2 \mathbf{R}_{xx} J_{ms_{\min}}\end{aligned}$$

Design Guidelines Extrated from the Model

- Stability

$$0 < \mu < \frac{m}{\text{Tr}[\mathbf{R}_{xx}]}; \quad m = 2 \text{ or } m = 2/3$$

- $J_{ms}(n)$ converges as a function of

$$[(1 - \mu\lambda_i)^2 + 2\mu^2\lambda_i^2]^n, \quad i = 1, \dots, N$$

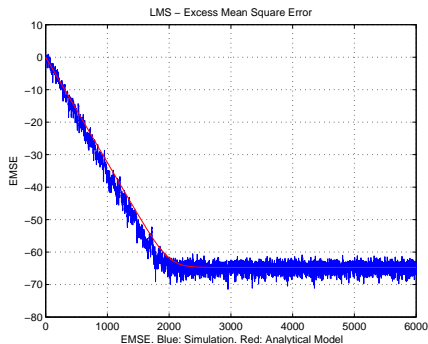
- Steady-State

$$J_{ms_{\text{ex}}}(\infty) \simeq \frac{\mu}{2} \text{Tr}[\mathbf{R}_{xx}] J_{ms_{\text{min}}}$$

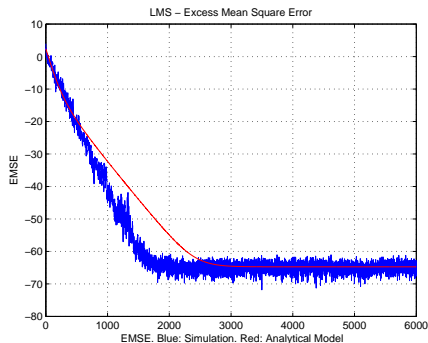
$$\mathcal{M}(\infty) = \frac{J_{ms_{\text{ex}}}(\infty)}{J_{ms_{\text{min}}}} \quad (\text{MSE Misadjustment})$$

Model Evaluation

$$x(n) = \alpha x(n-1) + v(n)$$



Input: White noise



Input: AR(1), $\alpha = 0.8$

- Linear system identification - FIR with 100 coefficients
- Step-size $\mu = 0.05$
- Noise power $\sigma_v^2 = 10^{-6}$

Merci pour votre attention!!!