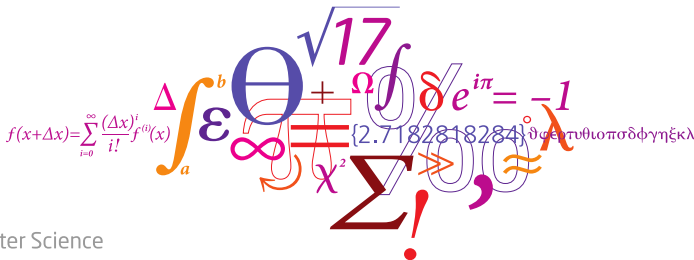


02471 Machine Learning for Signal Processing

Adaptive Linear Filtering with LMS

Tommy Sonne Alstrøm

Cognitive Systems section



Outline

- Course admin
- Last week review
- LMS and NLMS
- Justification of LMS
- Affine Projection Algorithm and normalized LMS
- Steady-state convergence analysis of LMS
- Next Week

Material: ML 2.6, 5.1–5.5.1, (skip 5.3.1), 5.6, 5.9.

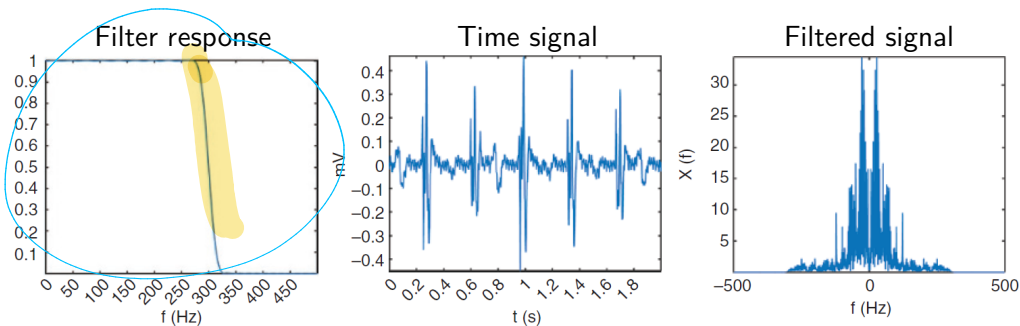
Administrative notes

- Problem set 1 second round, Sunday 6/10, pass/no pass.
- Problem set 2 first deadline, Sunday 27/10, pass/no pass.
 - Problem 2.1, 2.2, 2.4, and 2.5 should be solvable now.
- Feedback from last week
 - If lecturing lasts for 2 hours, it is better to have 2 times 10 min breaks.
 - Create a list of symbols used throughout the course (will create curated overleaf).
- Feedback from you is a critical component for improving both the course and my teaching.
- Type of feedback
 - Mention one thing that worked?
 - Mention one thing should be improved (both in current lecture and last weeks exercise)?
 - Mention one thing you would change if you gave the lecture.

Last week review

Two approaches to filter design

Design filter with specific frequency response:



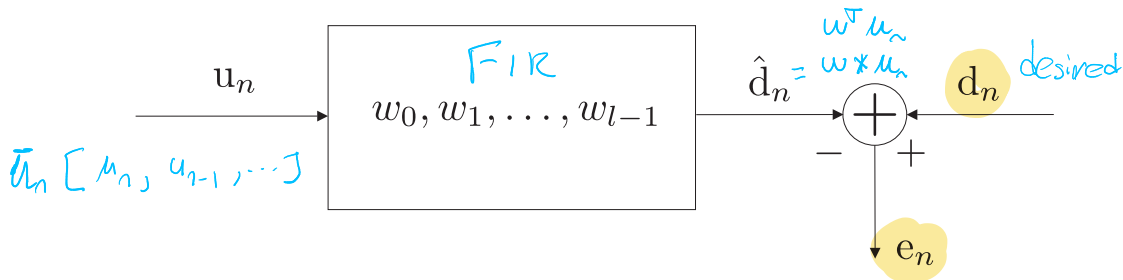
Idea:

Create a filter as a learning problem, ie. one that **minimized a cost function.**

Last week review

Linear filtering as learning problem

Filtering as a learning problem



The filter weights θ (denoted as $\mathbf{w} = [w_0, w_1, \dots, w_{l-1}]^T$ to indicate a linear filter) are learned such that the error is minimized

$$e_n = d_n - \hat{d}_n$$

If we choose to minimize the mean squared error, the filter is called a **Wiener filter**.

To use such filter, we need to specify a **desired** signal, d_n .

The normal equations

Minimizing the mean squared error $J(\boldsymbol{\theta}) = \mathbb{E}[(d_n - \hat{d}_n)^2]$ results in the normal equations

Normal Equations

$$\Sigma_u \boldsymbol{\theta}_* = \mathbf{r}_{du} \quad (\mathbf{r}_{du} \text{ is denoted } \mathbf{p} \text{ in the book})$$

Cross-cor $\mathbf{p} = [r_{du}(0) \ r_{du}(1) \ \cdots \ r_{du}(l-1)]^T$

Covar./
Auto-corr $\Sigma_u =$

$$\begin{bmatrix} r_u(0) & r_u(1) & \cdots & r_u(l-1) \\ r_u(1) & r_u(0) & \cdots & r_u(l-2) \\ \vdots & \vdots & \ddots & \vdots \\ r_u(l-1) & r_u(l-2) & \cdots & r_u(0) \end{bmatrix}$$

Wiener filter solution

$$\boldsymbol{\theta}_* = \Sigma_u^{-1} \mathbf{r}_{du}$$

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

Stochastic processes

Auto-correlation and cross-correlation expressions for WSS processes (unnormalized)

$$\mu_n = \mu \quad (\text{constant mean})$$

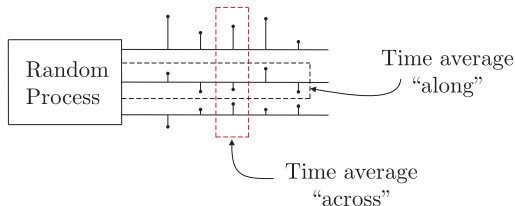
$$r_u(k) = \mathbb{E}[u_n u_{n-k}] \rightarrow \text{func. of } k$$

$$r_{uv}(k) = \mathbb{E}[u_n v_{n-k}]$$



These are unnormalized, see <https://en.wikipedia.org/wiki/Autocorrelation>

Additionally, for mean-ergodic and covariance-ergodic processes



$$\mu_u = \frac{1}{N} \sum_{n=1}^N u_n$$

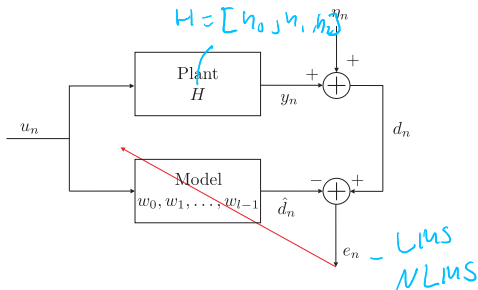
$$r_{uv}(k) = \frac{1}{N} \sum_{n=0}^{N-1} u(n)v(n-k), \quad k = 0, 1, \dots, N-1$$

$$r'_{uv}(k) = \frac{1}{N-k} \sum_{n=k}^{N-1} u(n)v(n-k), \quad k = 0, 1, \dots, N-1$$

Typical applications

System identification

Goal: Model the impulse response of H



We have access to:

- The input signal u_n , and the noisy output d_n

Common for all: usually the signals are not wide-sense stationary but may be locally wide-sense stationary.

Other applications

- Interference cancellation
- Echo cancellation
- Channel Equalization
- EEG signal analysis
- ...

Last week summary

- Instead of designing filters by hand, we can train the filters from data. If we use a LTI filter and minimize the MSE the filter is called Wiener filter (or Linear filtering).
- We apply theory from stochastic processes to analyze and design wiener filters.
- We will in general be working with wide-sense stationary and ergodic processes. That means we can estimate the mean and correlation functions using only “one” realization (example). This makes the system applicable to real data.
- Many of the applications require adaptive filtering / online learning to be useful in practice.

LMS and NLMS

Adaptive filtering problem statement

- The signal statistics are slowly changing, so we need to change our filter weights as well.
- We need to have access to second-order statistics.
- In practice we rarely have access to second-order statistics, so we will focus on algorithm that can learn the statistics iteratively using data.

The Least-Mean-Squares Algorithm

The LMS algorithm – algorithm 5.1 in the book

• Initialize

- $\theta_{-1} = \mathbf{0} \in \mathbb{R}^l$ filter length
- Select the value of μ - step size

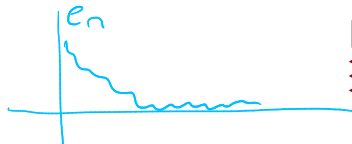
• For $n = 0, 1, \dots$, Do

- $e_n = y_n - \theta_{n-1}^T \mathbf{x}_n$ target - output, $\bar{\mathbf{x}}_n = [x_n, x_{n-l}, \dots, x_{n-l+1}]$
- $\theta_n = \theta_{n-1} + \mu e_n \mathbf{x}_n$ - vector ??
scalar

• End For

Parameters:

 μ is the step size. l is a filter length.



NLMS

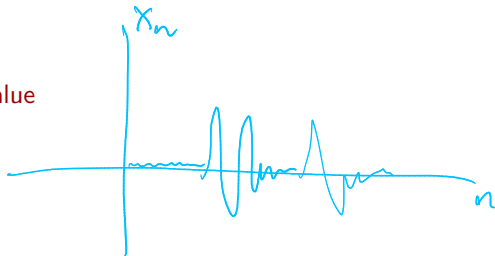
• Initialize

- $\theta_{-1} = \mathbf{0} \in \mathbb{R}^l$
- Select the value of $0 < \mu < 2$, and a small δ value

• For $n = 0, 1, \dots$, Do

- $e_n = y_n - \theta_{n-1}^T \mathbf{x}_n$
- $\theta_n = \theta_{n-1} + \frac{\mu}{\mathbf{x}_n^T \mathbf{x}_n + \delta} e_n \mathbf{x}_n$

• End For

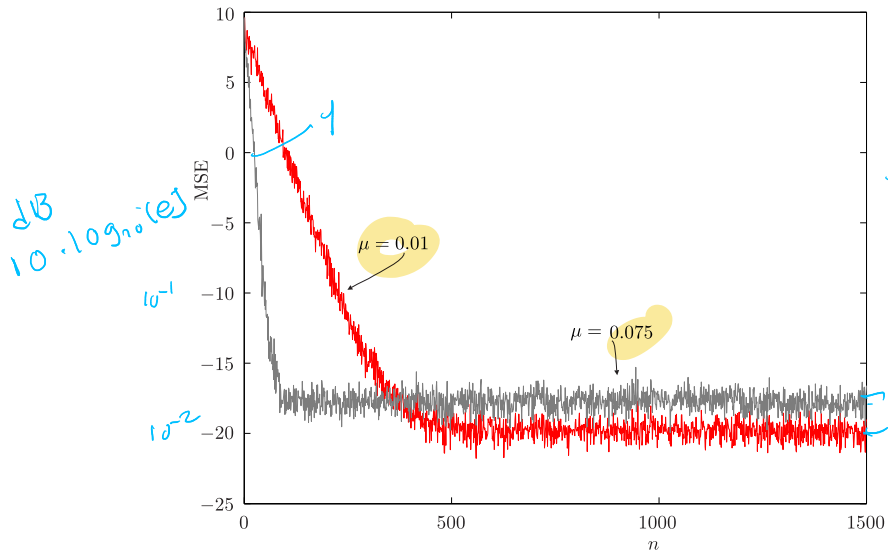


Parameters:

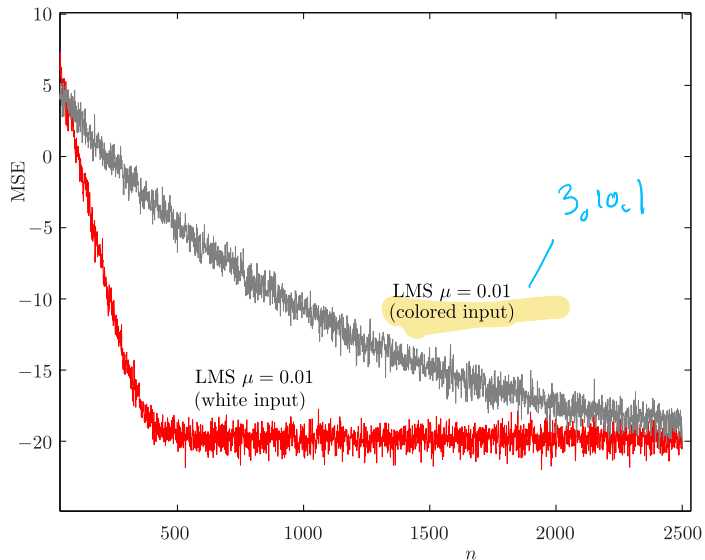
 μ is the step size. l is a filter length.

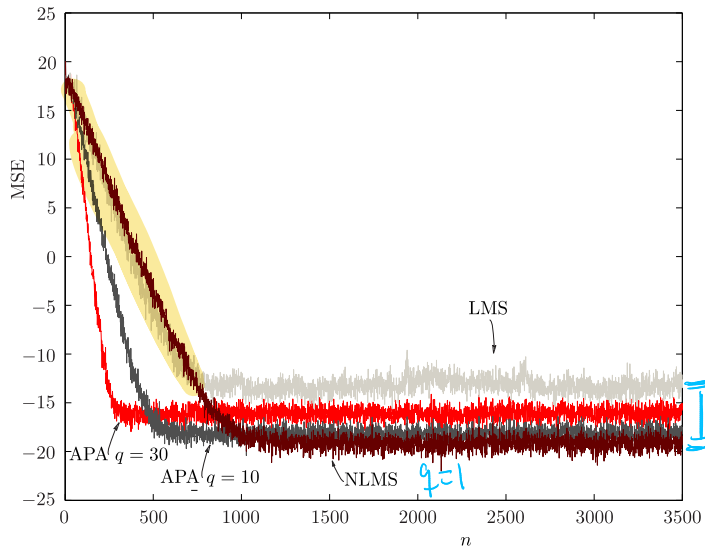
Q What is the role
power of the signal

Convergence rate LMS

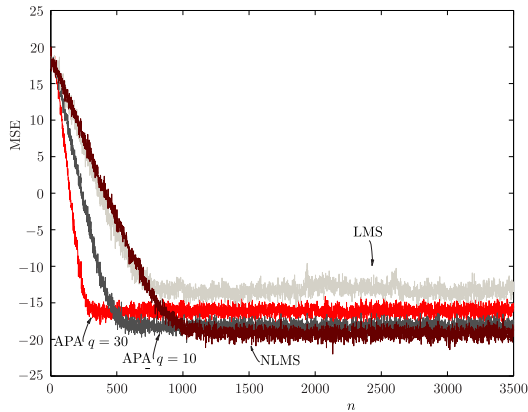
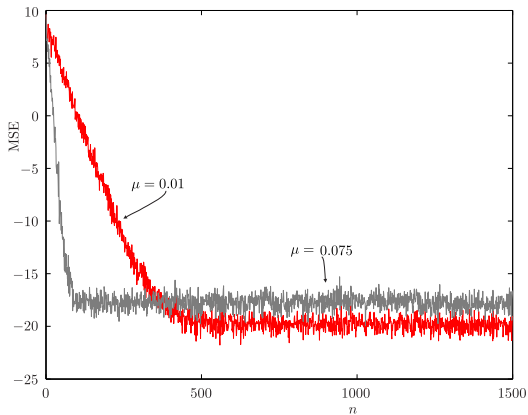


Convergence rate LMS



Convergence rate **NLMS**

LMS and NLMS Convergences



Discuss pros and cons of learning curves - can you think of applications appropriate of the difference curves?

- Many applications require adaptive algorithms (also called sequential learning, or online learning in machine learning).
- The value of μ affects the performance greatly (need theory to guide us).
- NLMS is generally preferred over LMS.

Justification of LMS

Steps involved in developing LMS

- ✓ ① Develop iterative scheme for when statistics are known ($r_x(k)$ and $r_{dx}(k)$ are known).
- ✓ ② Convert iterative scheme to data-driven approach ($r_x(k)$ and $r_{dx}(k)$ are estimated).
- ✓ ③ Ensure algorithm has agility (can adapt to changing environment).
- ④ Perform convergence analysis for steady state.

How did we arrive at LMS?

We want to develop an iterative scheme *adjustment*

$$\boldsymbol{\theta}^{(i)} = \boldsymbol{\theta}^{(i-1)} + \boxed{\mu_i \Delta \boldsymbol{\theta}^{(i)}}, \quad \mu_i > 0, \quad i > 0$$

such that

$$J(\boldsymbol{\theta}^{(i)}) < J(\boldsymbol{\theta}^{(i-1)})$$

where $J(\cdot)$ is the differentiable cost function.

Justification of LMS

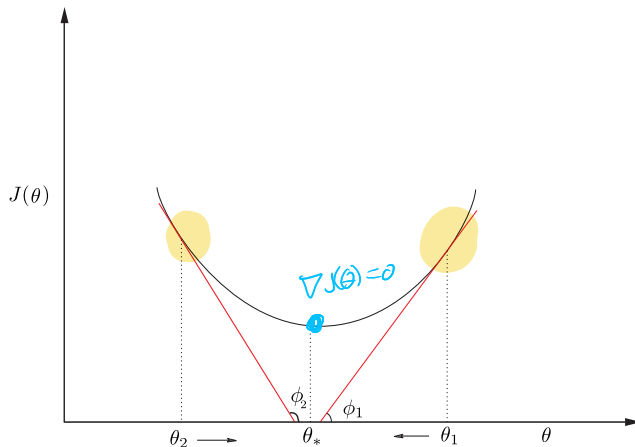
1-dimensional example

Iterative scheme:

$$\theta^{(i)} = \theta^{(i-1)} + \mu_i \Delta \theta^{(i)}$$

$\Delta \theta^{(i)} = -\nabla J(\theta)$

If $\Delta \theta^{(i)}$ is the negative gradient direction we will end up at the minima.



From gradient descent to LMS

Gradient Descent

$$\theta^{(i)} = \theta^{(i-1)} - \mu_i \nabla J(\theta^{(i-1)})$$

From last week, we found the gradient, $\nabla J(\theta) = \Sigma_x \theta - p$. week 3

At the optimum, $\nabla J(\theta) = \Sigma_x \theta - p = 0$.

$$\begin{array}{c} r_x(k) \\ \downarrow \\ E[x_n x_{n-q}] \end{array} \quad \begin{array}{c} r_{dx}(k) - E[x_n d_{n-k}] \end{array}$$

Issue: this only works if we have access to second order statistics.

Use Robbins–Monro scheme – stochastic approximation

If we have a cost function defined w.r.t. an expectation, we can apply the Robbins-Monro algorithm

Robbins-Monro algorithm



$$J(\theta) = \mathbb{E}[\mathcal{L}(\theta, \mathbf{y}, \mathbf{x})]$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}[\nabla \mathcal{L}(\theta, \mathbf{y}, \mathbf{x})]$$

$$\theta_n = \theta_{n-1} - \mu_n \nabla \mathcal{L}(\theta_{n-1}, \mathbf{y}_n, \mathbf{x}_n)$$

$$\sum_n \mu_n^2 < \infty, \quad \sum_n \mu_n \rightarrow \infty : \quad \text{convergence conditions.}$$

$$\mu_{n+1} < \mu_n$$

$$\mu_n = \frac{1}{n}$$

Derivation of the update

$$\begin{aligned}
 J(\boldsymbol{\theta}) &= \mathbb{E}[e^2] \quad \text{MSE} \quad e = y - \boldsymbol{\theta}^T \mathbf{x} \quad | \quad d - \hat{d} = \theta^T x \\
 \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \mathbb{E}[\nabla_{\boldsymbol{\theta}} e^2] \\
 &\stackrel{\text{Chain rule}}{=} \mathbb{E}[2e \nabla_{\boldsymbol{\theta}} e] \\
 &= -2 \mathbb{E}[(y - \boldsymbol{\theta}^T \mathbf{x}) \mathbf{x}] = -2 \mathbb{E}[e \cdot \mathbf{x}]
 \end{aligned}$$

The Robbins–Monro scheme was

$$\boldsymbol{\theta}_n = \boldsymbol{\theta}_{n-1} + \mu_n \nabla \mathcal{L}(\boldsymbol{\theta}_{n-1}, \mathbf{y}_n, \mathbf{x}_n)$$

So we get

$$\boldsymbol{\theta}_n = \boldsymbol{\theta}_{n-1} + \mu_n (y_n - \boldsymbol{\theta}_{n-1}^T \mathbf{x}_n) \mathbf{x}_n$$

$-2 \Rightarrow$ absorbed in μ

The Robbins-Monro update

$$\theta^{(i)} = \theta^{(i-1)} + \mu_n (y_n - \mathbf{x}_n^T \theta^{(i-1)}) \mathbf{x}_n \quad \text{last slide}$$

To gain agility to track non-stationary environments, we use a constant step size.

LMS update

$$\theta^{(i)} = \theta^{(i-1)} + \boxed{\mu} (y_n - \mathbf{x}_n^T \theta^{(i-1)}) \mathbf{x}_n = \theta_{n-1}^{(i-1)} + \mu \cdot e_n \cdot \mathbf{x}_n$$

fixed step size θ_{n-1}

✓ LMS

SGD

- The LMS algorithm is a stochastic gradient algorithm that estimates second-order statistics from data.
- To obtain agility, we use a fixed step-size. This means convergence guaranties provided by Robbins-Monro are invalid. $\frac{1}{2}$
 - Care must be taken for choosing the step-size (sec 5.3).
 - We can use the eigenvalues of Σ_x to guide our choice of step-size.
 - Large spread in eigenvalues means slow convergence (sec 5.3).
- LMS is a flexible algorithm that requires tuning of the step-size.

} LA

Affine Projection Algorithm and normalized LMS

The Affine Projection Algorithm (APA)

The APA algorithm minimizes the following problem

$$\theta_n = \arg \min_{\theta} \|\theta - \theta_{n-1}\|^2$$

$$\text{s.t. } \mathbf{x}_{n-i}^T \theta = y_{n-i}, \quad i = 0, 1, \dots, q-1$$

LA : $\mathbf{x}_n^T \theta = y_n$ $q=1$ assume FIR $l \geq q$

DESIGN

Lagrange multiplier

$$\begin{aligned} &\text{minimize} && J(\theta) \\ &\text{subject to} && f_i(\theta) = 0, \quad i = 1, 2, \dots, m \end{aligned}$$

This problem can be solved by optimizing a modified function, usually called the Lagrangian

$$L(\theta, \lambda) = J(\theta) - \sum_{i=1}^m \lambda_i f_i(\theta)$$

Lag. Mult

The saddle point of this function is then found by solving $\nabla L(\theta, \lambda) = 0$

Deriving updates for APA I

$$\boldsymbol{\theta}_n = \arg \min \|\boldsymbol{\theta} - \boldsymbol{\theta}_{n-1}\|^2$$

$$\text{s.t. } \mathbf{x}_{n-i}^T \boldsymbol{\theta} = y_{n-i}, \quad i = 0, 1, \dots, q-1$$

$$L(\boldsymbol{\theta}, \boldsymbol{\lambda}) = \underbrace{(\boldsymbol{\theta} - \boldsymbol{\theta}_{n-1})^T (\boldsymbol{\theta} - \boldsymbol{\theta}_{n-1})}_{J(\boldsymbol{\theta})} + \sum_{i=0}^{q-1} \lambda_i \left(y_{n-i} - \underbrace{\boldsymbol{\theta}^T \mathbf{x}_{n-i}}_{\mathbf{x}_{n-i}^T \boldsymbol{\theta}} \right)$$

$$\frac{d}{d\boldsymbol{\theta}} L(\boldsymbol{\theta}, \boldsymbol{\lambda}) = 2\boldsymbol{\theta} - 2\boldsymbol{\theta}_{n-1} - \sum_{i=0}^{q-1} \lambda_i \mathbf{x}_{n-i}$$

$$\frac{d}{d\boldsymbol{\theta}} L(\boldsymbol{\theta}, \boldsymbol{\lambda}) = \mathbf{0} \Rightarrow \boldsymbol{\theta} = \boldsymbol{\theta}_{n-1} + \frac{1}{2} \sum_{i=0}^{q-1} \lambda_i \mathbf{x}_{n-i}$$

$$\mathbf{X}_n := [\mathbf{x}_n, \dots, \mathbf{x}_{n-q+1}]^T, \quad \boldsymbol{\lambda}^T := [\lambda_0, \dots, \lambda_{q-1}], \quad \mathbf{y}_n^T := [y_n, \dots, y_{n-q+1}]$$

Deriving updates for APA II

$$\boldsymbol{\theta} = \boldsymbol{\theta}_{n-1} + \frac{1}{2} X_n^T \boldsymbol{\lambda}$$

$$X_n \boldsymbol{\theta} = \mathbf{y}_n$$

Solve the two equations for the two unknowns:

$$X_n \left(\boldsymbol{\theta}_{n-1} + \frac{1}{2} X_n^T \boldsymbol{\lambda} \right) = \mathbf{y}_n \quad \Leftrightarrow$$

$$\frac{1}{2} X_n X_n^T \boldsymbol{\lambda} = \mathbf{y}_n - X_n \boldsymbol{\theta}_{n-1} \quad \Leftrightarrow$$

$$\frac{1}{2} \boldsymbol{\lambda} = \left(X_n X_n^T \right)^{-1} (\mathbf{y}_n - X_n \boldsymbol{\theta}_{n-1}) \Rightarrow$$

$$\boldsymbol{\theta} = \boldsymbol{\theta}_{n-1} + X_n^T \left(X_n X_n^T \right)^{-1} (\mathbf{y}_n - X_n \boldsymbol{\theta}_{n-1})$$

Affine Projection Algorithm and normalized LMS

Deriving updates for APA III

$$\theta = \theta_{n-1} + X_n^T (X_n X_n^T)^{-1} (y_n - X_n \theta_{n-1})$$

$$e_n := y_n - X_n \theta_{n-1}$$

$$\theta = \theta_{n-1} + X_n^T (X_n X_n^T)^{-1} e_n$$

APA sol. (no step size)

The update algorithm then becomes

$$e_n = y_n - X_n \theta_{n-1}$$

$$\theta = \theta_{n-1} + \mu X_n^T (\delta I + X_n X_n^T)^{-1} e_n$$

For $q=1$

NLMS

$$e_n = y_n - x_n \theta_{n-1}$$

$$\theta_n = \theta_{n-1} + \frac{\mu}{x_n^T x_n + \delta} e_n x_n$$

we hit x_n
iff $\mu=1$

The Affine Projection Algorithm (APA) and NLMS

APA and NLMS

The APA algorithm minimizes the following problem

$$\begin{aligned} \theta_n &= \arg \min_{\theta} \|\theta - \theta_{n-1}\|^2 \\ \text{s.t.} \quad & \mathbf{x}_{n-i}^T \theta = y_{n-i}, \quad i = 0, 1, \dots, q-1 \end{aligned}$$

Define the following matrix

$$X_n^T := [\mathbf{x}_n, \dots, \mathbf{x}_{n-q+1}]$$

Then, by using Lagrange multipliers, we get the following update, $(0 \leq \mu \leq 2)$

$$\begin{aligned} e_n &= \mathbf{y}_n - X_n \theta_{n-1} \\ \theta &= \theta_{n-1} + \mu X_n^T (\delta I + X_n X_n^T)^{-1} e_n \end{aligned}$$

For $q = 1$, we call the algorithm nlms.

Steady-state convergence analysis of LMS

Steady-state convergence analysis of LMS

Convergence analysis of coefficients

5.5.1

Define the coefficient error vector as

$$\mathbf{c}_n := \boldsymbol{\theta}_n - \boldsymbol{\theta}_*$$

Analyzing convergence in the mean,

$$\mathbb{E}[\mathbf{c}_n] \rightarrow 0$$

as $n \rightarrow \infty$ we get the following results

$$0 < \mu < \frac{2}{\lambda_{\max}}$$

eigenval(Σ_x) Σ_x eq 1

$y: 1\% + 1\%$
 $N: 2\%$
 $??: 97\%$

However analyzing the covariance matrix, $n \rightarrow \infty$

$$\Sigma_{c,n} := \mathbb{E}[\mathbf{c}_n \mathbf{c}_n^T]$$

$$0 < \mu < \frac{2}{\text{trace}\{\Sigma_x\}} = \frac{2}{l \cdot r_x(0)}$$

eq. 2

$y: 5\%$
 $N: 1\%$
 $??: 94\%$

Steady-state convergence analysis of LMS

Convergence analysis of error

ML 5.5.1

Define the misalignment as

Misalignment

$$\mathcal{M} := \frac{J_{\text{exc}}}{J_{\text{min}}}$$

excess sol MSE
Best sol MSE

Analyzing the mean squared error behavior, we get the error at step n :

$$J_n := \mathbb{E} [e_n^2] = \dots (\text{ML eq 5.46} - 5.47) \dots = J_{\text{min}} + \text{trace} \{ \Sigma_x \Sigma_{c,n-1} \}$$

Excess MSE at time instant n

$$J_{\text{exc},n} = \text{trace} \{ \Sigma_x \Sigma_{c,n-1} \} : \quad \text{excess MSE at time instant } n$$

Expected misalignment, for small values of μ

$$\mathcal{M} \simeq \frac{1}{2} \mu \text{trace} \{ \Sigma_x \} : \quad \text{misadjustment}$$

is $\text{Tr}(\Sigma_x)$ constant
y: sol
N:

Steady-state convergence analysis of LMS

Summary



- The convergence analysis involve a fair amount of assumptions and is lengthy to derive.
- You are expected to be familiar with the material but can skip any derivation that the book defers to “Problem X derives this”.
- The convergence analysis provides theoretical foundation for choosing the step-size μ . We need agility in the algorithm, but also low misalignment. This is a trade-off.

Steady-state convergence analysis of LMS

Lecture summary



- Adaptive algorithms, such as LMS and NLMS can be used to estimate second-order statistics from data.
- The algorithms need to be carefully tuned per application basis, as there is trade-offs between agility and total error.
- In general, NLMS is considered more stable and better than LMS.
- These algorithms are deployed in many real-world problems to enable filtering in a changing environment.

ML

2022! prob 2.4

Material: ML 5.12, 6.1–6.3 (until "The LS estimator is BLUE"), 6.6–6.8, 6.12.

- Adaptive filtering continued
- Tracking performance in non-stationary environment
- Recursive Least-Squares (RLS) adaptive algorithm