

Machine learning for signal processing *[5LSL0]*

Ruud van Sloun, Rik Vullings

April 2019

TU/e

Technische Universiteit
Eindhoven
University of Technology

Where innovation starts

Main purpose of this course:

*Describe machine learning from a signal processing perspective +
hands-on experience*

- Main content:

- Main content:

- ▶ Optimum linear filters
- ▶ Adaptive Signal Processing
- ▶ Activation functions for classification
- ▶ Optimization and regularization
- ▶ Deep (convolutional) neural networks
- ▶ Variational neural networks

- Organization:

- ▶ Lab sessions: work in groups of 2 students and submit report

- Organization:

- ▶ Lab sessions: work in groups of 2 students and submit report
- ▶ Project: work in groups of 5 students and give presentation

- Organization:

- ▶ Lab sessions: work in groups of 2 students and submit report
- ▶ Project: work in groups of 5 students and give presentation

- Assessment:

- Organization:

- ▶ Lab sessions: work in groups of 2 students and submit report
- ▶ Project: work in groups of 5 students and give presentation

- Assessment:

- ▶ Lab sessions: 20% of grade

- Organization:

- ▶ Lab sessions: work in groups of 2 students and submit report
- ▶ Project: work in groups of 5 students and give presentation

- Assessment:

- ▶ Lab sessions: 20% of grade
- ▶ Oral exam:

- Organization:

- ▶ Lab sessions: work in groups of 2 students and submit report
- ▶ Project: work in groups of 5 students and give presentation

- Assessment:

- ▶ Lab sessions: 20% of grade
- ▶ Oral exam:
 - With group in week 26 and 27

- Organization:

- ▶ Lab sessions: work in groups of 2 students and submit report
- ▶ Project: work in groups of 5 students and give presentation

- Assessment:

- ▶ Lab sessions: 20% of grade
- ▶ Oral exam:
 - With group in week 26 and 27
 - Counts for 80% of grade

- Organization:

- ▶ Lab sessions: work in groups of 2 students and submit report
- ▶ Project: work in groups of 5 students and give presentation

- Assessment:

- ▶ Lab sessions: 20% of grade
- ▶ Oral exam:
 - With group in week 26 and 27
 - Counts for 80% of grade
 - Each student assessed individually

- Organization:

- ▶ Lab sessions: work in groups of 2 students and submit report
- ▶ Project: work in groups of 5 students and give presentation

- Assessment:

- ▶ Lab sessions: 20% of grade
- ▶ Oral exam:
 - With group in week 26 and 27
 - Counts for 80% of grade
 - Each student assessed individually
 - (anonymous) peer review

- Organization:

- ▶ Lab sessions: work in groups of 2 students and submit report
- ▶ Project: work in groups of 5 students and give presentation

- Assessment:

- ▶ Lab sessions: 20% of grade
- ▶ Oral exam:
 - With group in week 26 and 27
 - Counts for 80% of grade
 - Each student assessed individually
 - (anonymous) peer review
- ▶ **To pass 5LSLO \rightarrow ORAL ≥ 5 AND LABS+ORAL ≥ 5.5**

1. Book: "Statistical and Adaptive Array Signal Processing: Spectral estimation, signal modeling, adaptive filtering and array processing";
Dimitris G. Manolakis, Vinay K. Ingle and Stephen M. Kogon;
McGraw Hill; 2003
 - Optimum linear filters (Chapter 6)
 - Adaptive filters (Chapter 10)

1. Book: "Statistical and Adaptive Array Signal Processing: Spectral estimation, signal modeling, adaptive filtering and array processing";
Dimitris G. Manolakis, Vinay K. Ingle and Stephen M. Kogon; McGraw Hill; 2003
 - Optimum linear filters (Chapter 6)
 - Adaptive filters (Chapter 10)
2. Book: "Deep learning";
Ian Goodfellow, Yoshua Bengio and Aaron Courville; The MIT Press; 2016.
 - Machine learning basics (Chapter 5)
 - Deep feedforward networks (Chapter 6)
 - Regularization (Chapter 7)
 - Optimization for training deep models (Chapter 8)
 - Convolutional networks (Chapter 9)

3. These slides

- 3. These slides
- 4. Udacity

3. These slides

4. Udacity

Books available online at

https://lr.ttu.ee/~ameister/materjale/Dig_spe_anal/Statistical%20and%20Adaptive%20Signal%20Processing.pdf

and

<https://www.deeplearningbook.org/>

Week	Monday	Thursday
<i>Monday: hour 3-4, Thursday: hour 7-8</i>		
Apr 22	<i>No lecture</i>	Optimum linear filters
Apr 29	Adaptive filters	Adaptive filters
May 6	Activation functions	Optimization
May 13	Regularization	Regularization
May 20	CNN	CNN
May 27	Variational network	Variational network
Jun 3	Intro project	Project
Jun 10	Project	Project
Jun 17	Project	Project

Code	Deadline	Credits
Ass 1	May 6, 10:30	5
Ass 2	May 16, 13:30	5
Ass 3	May 23, 13:30	5
Ass 4	June 3, 10:30	5
Oral	June 24- July 6	80
Total		100

Code	Deadline	Credits
Ass 1	May 6, 10:30	5
Ass 2	May 16, 13:30	5
Ass 3	May 23, 13:30	5
Ass 4	June 3, 10:30	5
Oral	June 24- July 6	80
Total		100

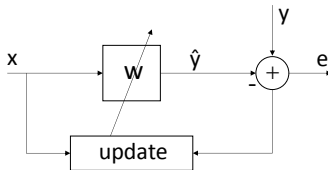
To pass 5LSLO \Rightarrow ORAL ≥ 5 and TOTAL ≥ 5.5

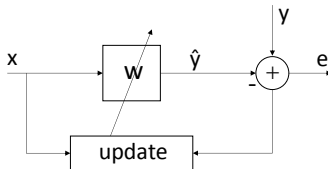
Optimum linear filters and adaptive filters

Focus on single channel adaptive algorithms using FIR structures

Focus on single channel adaptive algorithms using FIR structures

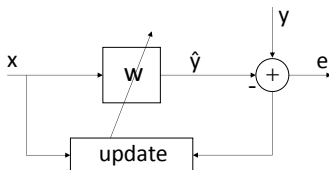
- ▶ Minimum Mean Squared Error
- ▶ Gradient Descent Algorithm
- ▶ Adaptive (N)LMS
- ▶ Newton algorithm
- ▶ Recursive Least Squares (RLS)





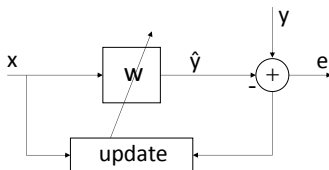
Notes:

- ▶ Input signal x and desired response y correlated



Notes:

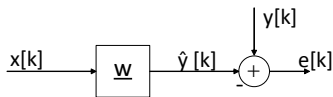
- ▶ Input signal x and desired response y correlated
- ▶ Pragmatic choices:
 - All signals have zero average
 - Filter w : FIR



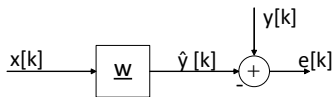
Notes:

- ▶ Input signal x and desired response y correlated
- ▶ Pragmatic choices:
 - All signals have zero average
 - Filter w : FIR
- ▶ Calculation of weight of filter w :
 - Use quadratic cost function: $J = f(e^2)$
 - **First fixed weights** (MMSE), then adaptive

General Minimum Mean Squared Error (MMSE) model:



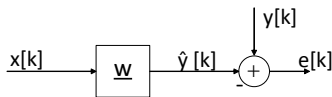
General Minimum Mean Squared Error (MMSE) model:



Goal:

Given N samples $\underline{x}[k] = (x[k], x[k-1], \dots, x[k-N+1])^t$ calculate coefficients fixed filter $\underline{w} = (w_0, w_1, \dots, w_{N-1})^t$ such that Mean Squared Error (MSE) $J = E\{e^2[k]\} = E\{(y[k] - \hat{y}[k])^2\}$ is minimized.

General Minimum Mean Squared Error (MMSE) model:



Goal:

Given N samples $\underline{x}[k] = (x[k], x[k-1], \dots, x[k-N+1])^t$ calculate coefficients fixed filter $\underline{w} = (w_0, w_1, \dots, w_{N-1})^t$ such that Mean Squared Error (MSE) $J = E \{e^2[k]\} = E\{(y[k] - \hat{y}[k])^2\}$ is minimized.

MMSE Optimization problem:

Given FIR samples $x[k-i]$ for $i = 0, 1, \dots, N-1$

$$\underline{w}_o = \arg \min_{\underline{w}} (E \{e^2[k]\})$$

$$\begin{aligned} J &= E\{(y[k] - \underline{\mathbf{w}}^t \cdot \underline{\mathbf{x}}[k]) \cdot (y[k] - \underline{\mathbf{x}}^t[k] \cdot \underline{\mathbf{w}})\} \\ &= E\{y^2[k]\} - \underline{\mathbf{w}}^t E\{\underline{\mathbf{x}}[k]y[k]\} - E\{y[k]\underline{\mathbf{x}}^t[k]\}\underline{\mathbf{w}} + \underline{\mathbf{w}}^t E\{\underline{\mathbf{x}}[k]\underline{\mathbf{x}}^t[k]\}\underline{\mathbf{w}} \end{aligned}$$

$$\begin{aligned} J &= E\{(y[k] - \underline{\mathbf{w}}^t \cdot \underline{\mathbf{x}}[k]) \cdot (y[k] - \underline{\mathbf{x}}^t[k] \cdot \underline{\mathbf{w}})\} \\ &= E\{y^2[k]\} - \underline{\mathbf{w}}^t E\{\underline{\mathbf{x}}[k]y[k]\} - E\{y[k]\underline{\mathbf{x}}^t[k]\}\underline{\mathbf{w}} + \underline{\mathbf{w}}^t E\{\underline{\mathbf{x}}[k]\underline{\mathbf{x}}^t[k]\}\underline{\mathbf{w}} \end{aligned}$$

$$\Rightarrow \boxed{J = E\{y^2[k]\} - \underline{\mathbf{w}}^t \underline{\mathbf{r}}_{yx} - \underline{\mathbf{r}}_{yx}^t \underline{\mathbf{w}} + \underline{\mathbf{w}}^t \underline{\mathbf{R}}_x \underline{\mathbf{w}}}$$

$$\begin{aligned} J &= E\{(y[k] - \underline{\mathbf{w}}^t \cdot \underline{\mathbf{x}}[k]) \cdot (y[k] - \underline{\mathbf{x}}^t[k] \cdot \underline{\mathbf{w}})\} \\ &= E\{y^2[k]\} - \underline{\mathbf{w}}^t E\{\underline{\mathbf{x}}[k]y[k]\} - E\{y[k]\underline{\mathbf{x}}^t[k]\}\underline{\mathbf{w}} + \underline{\mathbf{w}}^t E\{\underline{\mathbf{x}}[k]\underline{\mathbf{x}}^t[k]\}\underline{\mathbf{w}} \end{aligned}$$

$$\Rightarrow \boxed{J = E\{y^2[k]\} - \underline{\mathbf{w}}^t \underline{\mathbf{r}}_{yx} - \underline{\mathbf{r}}_{yx}^t \underline{\mathbf{w}} + \underline{\mathbf{w}}^t \underline{\mathbf{R}}_x \underline{\mathbf{w}}}$$

with cross correlation $\rho_{yx}[\tau] = E\{y[k]x[k - \tau]\}$:

$$\underline{\mathbf{r}}_{yx} = E\{y[k]\underline{\mathbf{x}}[k]\} = (\rho_{yx}[0], \rho_{yx}[1], \dots, \rho_{yx}[N - 1])^t$$

$$\begin{aligned} J &= E\{(y[k] - \underline{\mathbf{w}}^t \cdot \underline{\mathbf{x}}[k]) \cdot (y[k] - \underline{\mathbf{x}}^t[k] \cdot \underline{\mathbf{w}})\} \\ &= E\{y^2[k]\} - \underline{\mathbf{w}}^t E\{\underline{\mathbf{x}}[k]y[k]\} - E\{y[k]\underline{\mathbf{x}}^t[k]\}\underline{\mathbf{w}} + \underline{\mathbf{w}}^t E\{\underline{\mathbf{x}}[k]\underline{\mathbf{x}}^t[k]\}\underline{\mathbf{w}} \end{aligned}$$

$$\Rightarrow \boxed{J = E\{y^2[k]\} - \underline{\mathbf{w}}^t \underline{\mathbf{r}}_{yx} - \underline{\mathbf{r}}_{yx}^t \underline{\mathbf{w}} + \underline{\mathbf{w}}^t \mathbf{R}_x \underline{\mathbf{w}}}$$

with cross correlation $\rho_{yx}[\tau] = E\{y[k]x[k - \tau]\}$:

$$\underline{\mathbf{r}}_{yx} = E\{y[k]\underline{\mathbf{x}}[k]\} = (\rho_{yx}[0], \rho_{yx}[1], \dots, \rho_{yx}[N - 1])^t$$

and autocorrelation: $\rho_x[\tau] = E\{x[k]x[k - \tau]\} = \rho_x[-\tau]$

$$\mathbf{R}_x = E\{\underline{\mathbf{x}}[k]\underline{\mathbf{x}}^t[k]\} = \begin{pmatrix} \rho_x[0] & \rho_x[1] & \cdots & \rho_x[N - 1] \\ \rho_x[1] & \rho_x[0] & \cdots & \rho_x[N - 2] \\ \vdots & \vdots & \vdots & \vdots \\ \rho_x[N - 1] & \rho_x[N - 2] & \cdots & \rho_x[0] \end{pmatrix}$$

$$J = E\{y^2[k]\} - \underline{\mathbf{w}}^t \underline{\mathbf{r}}_{yx} - \underline{\mathbf{r}}_{yx}^t \underline{\mathbf{w}} + \underline{\mathbf{w}}^t \underline{\mathbf{R}}_x \underline{\mathbf{w}}$$

$$J = E\{y^2[k]\} - \underline{\mathbf{w}}^t \underline{\mathbf{r}}_{yx} - \underline{\mathbf{r}}_{yx}^t \underline{\mathbf{w}} + \underline{\mathbf{w}}^t \mathbf{R}_x \underline{\mathbf{w}}$$

$$\Rightarrow \text{Optimum: } \underline{\nabla} = \frac{dJ}{d\underline{\mathbf{w}}} = -2(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}) = \underline{\mathbf{0}}$$

$$J = E\{y^2[k]\} - \underline{\mathbf{w}}^t \underline{\mathbf{r}}_{yx} - \underline{\mathbf{r}}_{yx}^t \underline{\mathbf{w}} + \underline{\mathbf{w}}^t \mathbf{R}_x \underline{\mathbf{w}}$$

$$\Rightarrow \text{Optimum: } \underline{\nabla} = \frac{dJ}{d\underline{\mathbf{w}}} = -2(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}) = \underline{\mathbf{0}}$$

\Rightarrow Normal Equations

$$\boxed{\mathbf{R}_x \cdot \underline{\mathbf{w}} = \underline{\mathbf{r}}_{yx}}$$

$$J = E\{y^2[k]\} - \underline{\mathbf{w}}^t \underline{\mathbf{r}}_{yx} - \underline{\mathbf{r}}_{yx}^t \underline{\mathbf{w}} + \underline{\mathbf{w}}^t \mathbf{R}_x \underline{\mathbf{w}}$$

$$\Rightarrow \text{Optimum: } \underline{\nabla} = \frac{dJ}{d\underline{\mathbf{w}}} = -2(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}) = \underline{\mathbf{0}}$$

\Rightarrow Normal Equations

$$\mathbf{R}_x \cdot \underline{\mathbf{w}} = \underline{\mathbf{r}}_{yx}$$

\Rightarrow Wiener filter

$$\underline{\mathbf{w}}_o = \mathbf{R}_x^{-1} \cdot \underline{\mathbf{r}}_{yx}$$

$$J = E\{y^2[k]\} - \underline{\mathbf{w}}^t \underline{\mathbf{r}}_{yx} - \underline{\mathbf{r}}_{yx}^t \underline{\mathbf{w}} + \underline{\mathbf{w}}^t \mathbf{R}_x \underline{\mathbf{w}}$$

$$\Rightarrow \text{Optimum: } \underline{\nabla} = \frac{dJ}{d\underline{\mathbf{w}}} = -2(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}) = \underline{\mathbf{0}}$$

\Rightarrow Normal Equations

$$\mathbf{R}_x \cdot \underline{\mathbf{w}} = \underline{\mathbf{r}}_{yx}$$

\Rightarrow Wiener filter

$$\underline{\mathbf{w}}_o = \mathbf{R}_x^{-1} \cdot \underline{\mathbf{r}}_{yx}$$

General expression: $J = J_{min} + (\underline{\mathbf{w}} - \underline{\mathbf{w}}_o)^t \cdot \mathbf{R}_x \cdot (\underline{\mathbf{w}} - \underline{\mathbf{w}}_o)$

$$J = E\{y^2[k]\} - \underline{\mathbf{w}}^t \underline{\mathbf{r}}_{yx} - \underline{\mathbf{r}}_{yx}^t \underline{\mathbf{w}} + \underline{\mathbf{w}}^t \mathbf{R}_x \underline{\mathbf{w}}$$

$$\Rightarrow \text{Optimum: } \underline{\nabla} = \frac{dJ}{d\underline{\mathbf{w}}} = -2(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}) = \underline{\mathbf{0}}$$

\Rightarrow Normal Equations

$$\mathbf{R}_x \cdot \underline{\mathbf{w}} = \underline{\mathbf{r}}_{yx}$$

\Rightarrow Wiener filter

$$\underline{\mathbf{w}}_o = \mathbf{R}_x^{-1} \cdot \underline{\mathbf{r}}_{yx}$$

General expression:

$$J = J_{min} + (\underline{\mathbf{w}} - \underline{\mathbf{w}}_o)^t \cdot \mathbf{R}_x \cdot (\underline{\mathbf{w}} - \underline{\mathbf{w}}_o)$$

$$J_{min} = J_{\underline{\mathbf{w}}=\underline{\mathbf{w}}_o} = E\{e^2[k]\} = E\{y^2[k]\} - \underline{\mathbf{r}}_{yx}^t \mathbf{R}_x^{-1} \underline{\mathbf{r}}_{yx}$$

$$J = E\{y^2[k]\} - \underline{\mathbf{w}}^t \underline{\mathbf{r}}_{yx} - \underline{\mathbf{r}}_{yx}^t \underline{\mathbf{w}} + \underline{\mathbf{w}}^t \mathbf{R}_x \underline{\mathbf{w}}$$

$$\Rightarrow \text{Optimum: } \underline{\nabla} = \frac{dJ}{d\underline{\mathbf{w}}} = -2(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}) = \underline{\mathbf{0}}$$

\Rightarrow Normal Equations

$$\mathbf{R}_x \cdot \underline{\mathbf{w}} = \underline{\mathbf{r}}_{yx}$$

\Rightarrow Wiener filter

$$\underline{\mathbf{w}}_o = \mathbf{R}_x^{-1} \cdot \underline{\mathbf{r}}_{yx}$$

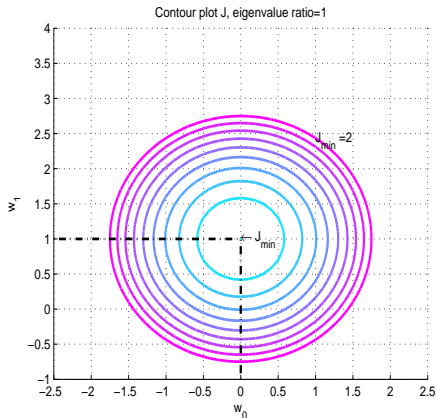
General expression:

$$J = J_{min} + (\underline{\mathbf{w}} - \underline{\mathbf{w}}_o)^t \cdot \mathbf{R}_x \cdot (\underline{\mathbf{w}} - \underline{\mathbf{w}}_o)$$

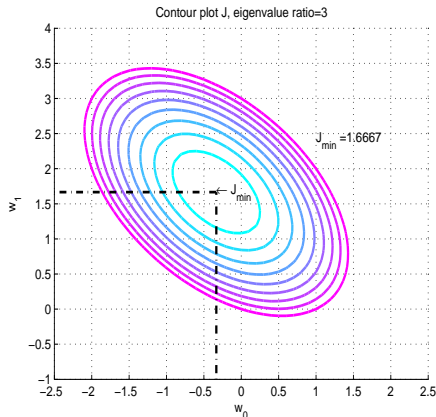
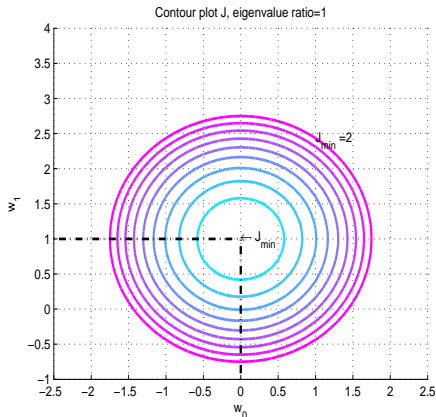
$$J_{min} = J_{\underline{\mathbf{w}}=\underline{\mathbf{w}}_o} = E\{e^2[k]\} = E\{y^2[k]\} - \underline{\mathbf{r}}_{yx}^t \mathbf{R}_x^{-1} \underline{\mathbf{r}}_{yx}$$

From general expression $\Rightarrow J$ quadratic in $\underline{\mathbf{w}}$ thus $\underline{\mathbf{w}}_o$ really minimum

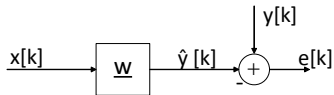
Contour plots $J = J_{min} + (\underline{\mathbf{w}} - \underline{\mathbf{w}}_o)^t \cdot \mathbf{R}_x \cdot (\underline{\mathbf{w}} - \underline{\mathbf{w}}_o)$



Contour plots $J = J_{min} + (\underline{\mathbf{w}} - \underline{\mathbf{w}}_o)^t \cdot \mathbf{R}_x \cdot (\underline{\mathbf{w}} - \underline{\mathbf{w}}_o)$



Eigenvalues: see Appendix





Different quadratic cost functions:

- Mean Square Error (MSE):

$$J_{mse} = E\{e^2[k]\} = E\{(y[k] - \underline{\mathbf{w}}^t \underline{\mathbf{x}}[k])^2\}$$

⇒ Minimum MSE (MMSE) = Wiener



Different quadratic cost functions:

- ▶ Mean Square Error (MSE):

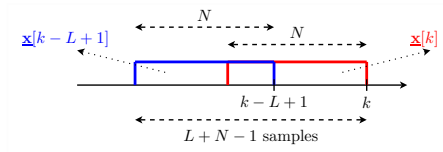
$$J_{mse} = E\{e^2[k]\} = E\{(y[k] - \underline{\mathbf{w}}^t \underline{\mathbf{x}}[k])^2\}$$

⇒ Minimum MSE (MMSE) = Wiener

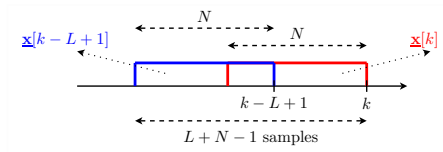
- ▶ Least Square (LS): If statistical information is not available ⇒

Use criterion based on data (thus without $E\{\cdot\}$)

Collect $L (\geq 1)$ data vectors $\underline{\mathbf{x}}[k - i]$ (each of length N)

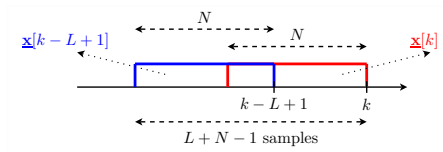


Collect $L (\geq 1)$ data vectors $\underline{\mathbf{x}}[k - i]$ (each of length N)



Available data (for $i = 0, 1, \dots, L - 1$):

Collect $L (\geq 1)$ data vectors $\underline{\mathbf{x}}[k - i]$ (each of length N)



Available data (for $i = 0, 1, \dots, L - 1$):

- *Input signal samples/ vectors* $\underline{\mathbf{x}}[k - i]$

$$\underline{\mathbf{x}}[k - i] = (x[k - i], x[k - i - 1], \dots, x[k - i - N + 1])^t$$

- *Reference signal samples:* $y[k - i]$
- *Residual signal samples:* $e[k - i] = y[k - i] - \underline{\mathbf{x}}^t[k - i] \cdot \underline{\mathbf{w}}$

Notation:

$$\mathbf{X}[k] = \begin{pmatrix} \underline{\mathbf{x}}^t[k] \\ \underline{\mathbf{x}}^t[k-1] \\ \vdots \\ \underline{\mathbf{x}}^t[k-L+1] \end{pmatrix}$$

$$\underline{\mathbf{y}}[k] = \begin{pmatrix} y[k] \\ y[k-1] \\ \vdots \\ y[k-L+1] \end{pmatrix}$$

$$\underline{\mathbf{w}} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_{N-1} \end{pmatrix}$$

$$\underline{\mathbf{e}}[k] = \begin{pmatrix} e[k] \\ e[k-1] \\ \vdots \\ e[k-L+1] \end{pmatrix}$$

Notation:

$$\mathbf{X}[k] = \begin{pmatrix} \underline{\mathbf{x}}^t[k] \\ \underline{\mathbf{x}}^t[k-1] \\ \vdots \\ \underline{\mathbf{x}}^t[k-L+1] \end{pmatrix} \quad \underline{\mathbf{w}} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_{N-1} \end{pmatrix}$$

$$\underline{\mathbf{y}}[k] = \begin{pmatrix} y[k] \\ y[k-1] \\ \vdots \\ y[k-L+1] \end{pmatrix} \quad \underline{\mathbf{e}}[k] = \begin{pmatrix} e[k] \\ e[k-1] \\ \vdots \\ e[k-L+1] \end{pmatrix}$$

Simplified notation (skip time indices):

$$\underline{\mathbf{e}} = \underline{\mathbf{y}} - \mathbf{X} \cdot \underline{\mathbf{w}}$$



LS problem formulation:

$$\underline{\mathbf{w}}_{ls,o} = \arg \min_{\underline{\mathbf{w}}} |\underline{\mathbf{y}} - \mathbf{X} \cdot \underline{\mathbf{w}}|^2$$

$$J_{ls} = \sum_{i=0}^{L-1} e^2[k-i] = \underline{\mathbf{e}}^t \cdot \underline{\mathbf{e}} = (\underline{\mathbf{y}}^t - \underline{\mathbf{w}}^t \mathbf{X}^t) \cdot (\underline{\mathbf{y}} - \mathbf{X} \underline{\mathbf{w}})$$

$$\begin{aligned} J_{ls} &= \sum_{i=0}^{L-1} e^2[k-i] = \underline{\mathbf{e}}^t \cdot \underline{\mathbf{e}} = (\underline{\mathbf{y}}^t - \underline{\mathbf{w}}^t \mathbf{X}^t) \cdot (\underline{\mathbf{y}} - \mathbf{X} \underline{\mathbf{w}}) \\ &= \underline{\mathbf{y}}^t \underline{\mathbf{y}} + \underline{\mathbf{w}}^t \mathbf{X}^t \mathbf{X} \underline{\mathbf{w}} - \underline{\mathbf{w}}^t \mathbf{X}^t \underline{\mathbf{y}} - \underline{\mathbf{y}}^t \mathbf{X} \underline{\mathbf{w}} \end{aligned}$$

$$\begin{aligned} J_{ls} &= \sum_{i=0}^{L-1} e^2[k-i] = \underline{\mathbf{e}}^t \cdot \underline{\mathbf{e}} = (\underline{\mathbf{y}}^t - \underline{\mathbf{w}}^t \mathbf{X}^t) \cdot (\underline{\mathbf{y}} - \mathbf{X} \underline{\mathbf{w}}) \\ &= \underline{\mathbf{y}}^t \underline{\mathbf{y}} + \underline{\mathbf{w}}^t \mathbf{X}^t \mathbf{X} \underline{\mathbf{w}} - \underline{\mathbf{w}}^t \mathbf{X}^t \underline{\mathbf{y}} - \underline{\mathbf{y}}^t \mathbf{X} \underline{\mathbf{w}} \end{aligned}$$

Minimum by setting gradient equal to zero:

$$\frac{dJ_{ls}}{d\underline{\mathbf{w}}} = \underline{\nabla}_{ls} = -2(\mathbf{X}^t \underline{\mathbf{y}} - \mathbf{X}^t \mathbf{X} \cdot \underline{\mathbf{w}}) = \underline{\mathbf{0}}$$

$$\begin{aligned}
 J_{ls} &= \sum_{i=0}^{L-1} e^2[k-i] = \underline{\mathbf{e}}^t \cdot \underline{\mathbf{e}} = (\underline{\mathbf{y}}^t - \underline{\mathbf{w}}^t \mathbf{X}^t) \cdot (\underline{\mathbf{y}} - \mathbf{X} \underline{\mathbf{w}}) \\
 &= \underline{\mathbf{y}}^t \underline{\mathbf{y}} + \underline{\mathbf{w}}^t \mathbf{X}^t \mathbf{X} \underline{\mathbf{w}} - \underline{\mathbf{w}}^t \mathbf{X}^t \underline{\mathbf{y}} - \underline{\mathbf{y}}^t \mathbf{X} \underline{\mathbf{w}}
 \end{aligned}$$

Minimum by setting gradient equal to zero:

$$\frac{dJ_{ls}}{d\underline{\mathbf{w}}} = \underline{\nabla}_{ls} = -2(\mathbf{X}^t \underline{\mathbf{y}} - \mathbf{X}^t \mathbf{X} \cdot \underline{\mathbf{w}}) = \underline{\mathbf{0}}$$

With $\overline{\mathbf{R}} = \mathbf{X}^t \mathbf{X}$ and $\underline{\mathbf{r}} = \mathbf{X}^t \underline{\mathbf{y}}$

$$\begin{aligned}
 J_{ls} &= \sum_{i=0}^{L-1} e^2[k-i] = \underline{\mathbf{e}}^t \cdot \underline{\mathbf{e}} = (\underline{\mathbf{y}}^t - \underline{\mathbf{w}}^t \mathbf{X}^t) \cdot (\underline{\mathbf{y}} - \mathbf{X} \underline{\mathbf{w}}) \\
 &= \underline{\mathbf{y}}^t \underline{\mathbf{y}} + \underline{\mathbf{w}}^t \mathbf{X}^t \mathbf{X} \underline{\mathbf{w}} - \underline{\mathbf{w}}^t \mathbf{X}^t \underline{\mathbf{y}} - \underline{\mathbf{y}}^t \mathbf{X} \underline{\mathbf{w}}
 \end{aligned}$$

Minimum by setting gradient equal to zero:

$$\frac{dJ_{ls}}{d\underline{\mathbf{w}}} = \underline{\nabla}_{ls} = -2(\mathbf{X}^t \underline{\mathbf{y}} - \mathbf{X}^t \mathbf{X} \cdot \underline{\mathbf{w}}) = \underline{\mathbf{0}}$$

With $\overline{\mathbf{R}} = \mathbf{X}^t \mathbf{X}$ and $\underline{\mathbf{r}} = \mathbf{X}^t \underline{\mathbf{y}}$

\Rightarrow **Normal Equations**

$$\overline{\mathbf{R}}_x \cdot \underline{\mathbf{w}} = \underline{\mathbf{r}}_{yx}$$

\Rightarrow **Wiener filter**

$$\underline{\mathbf{w}}_{ls,o} = \overline{\mathbf{R}}_x^{-1} \cdot \underline{\mathbf{r}}_{yx}$$

Use time-averaging (ergodicity):

$$\hat{\mathbf{R}}_x = \frac{1}{L} \sum_{i=0}^{L-1} \underline{\mathbf{x}}[k-i] \cdot \underline{\mathbf{x}}^t[k-i] = \frac{1}{L} \mathbf{X}^t \cdot \mathbf{X} = \frac{1}{L} \overline{\mathbf{R}}_x$$

$$\hat{\underline{\mathbf{r}}}_{yx} = \frac{1}{L} \sum_{i=0}^{L-1} \underline{\mathbf{x}}[k-i] \cdot y[k-i] = \frac{1}{L} \mathbf{X}^t \cdot \underline{\mathbf{y}} = \frac{1}{L} \overline{\underline{\mathbf{r}}}_{yx}$$

Use time-averaging (ergodicity):

$$\hat{\mathbf{R}}_x = \frac{1}{L} \sum_{i=0}^{L-1} \mathbf{x}[k-i] \cdot \mathbf{x}^t[k-i] = \frac{1}{L} \mathbf{X}^t \cdot \mathbf{X} = \frac{1}{L} \overline{\mathbf{R}}_x$$

$$\hat{\mathbf{r}}_{yx} = \frac{1}{L} \sum_{i=0}^{L-1} \mathbf{x}[k-i] \cdot y[k-i] = \frac{1}{L} \mathbf{X}^t \cdot \mathbf{y} = \frac{1}{L} \overline{\mathbf{r}}_{yx}$$

with $\hat{\mathbf{R}}_x$ estimate of \mathbf{R}_x and $\hat{\mathbf{r}}_{yx}$ estimate of \mathbf{r}_{yx}

$$\Rightarrow \hat{\mathbf{w}}_{mmse} = \left(\frac{1}{L} \overline{\mathbf{R}}_x \right)^{-1} \cdot \left(\frac{1}{L} \overline{\mathbf{r}}_{yx} \right) = \overline{\mathbf{R}}_x^{-1} \cdot \overline{\mathbf{r}}_{yx} = \mathbf{w}_{ls}$$

Use time-averaging (ergodicity):

$$\hat{\mathbf{R}}_x = \frac{1}{L} \sum_{i=0}^{L-1} \mathbf{x}[k-i] \cdot \mathbf{x}^t[k-i] = \frac{1}{L} \mathbf{X}^t \cdot \mathbf{X} = \frac{1}{L} \overline{\mathbf{R}}_x$$

$$\hat{\mathbf{r}}_{yx} = \frac{1}{L} \sum_{i=0}^{L-1} \mathbf{x}[k-i] \cdot y[k-i] = \frac{1}{L} \mathbf{X}^t \cdot \mathbf{y} = \frac{1}{L} \overline{\mathbf{r}}_{yx}$$

with $\hat{\mathbf{R}}_x$ estimate of \mathbf{R}_x and $\hat{\mathbf{r}}_{yx}$ estimate of \mathbf{r}_{yx}

$$\Rightarrow \hat{\mathbf{w}}_{mmse} = \left(\frac{1}{L} \overline{\mathbf{R}}_x \right)^{-1} \cdot \left(\frac{1}{L} \overline{\mathbf{r}}_{yx} \right) = \overline{\mathbf{R}}_x^{-1} \cdot \overline{\mathbf{r}}_{yx} = \mathbf{w}_{ls}$$

Finally note that for ergodic processes:

$$\lim_{L \rightarrow \infty} \frac{1}{L} \overline{\mathbf{R}}_x = \mathbf{R}_x ; \lim_{L \rightarrow \infty} \frac{1}{L} \overline{\mathbf{r}}_{yx} = \mathbf{r}_{yx} ; \lim_{L \rightarrow \infty} \mathbf{w}_{ls} = \mathbf{w}_{mmse}$$

Problem: Optimal Wiener involves \mathbf{R}_x^{-1}

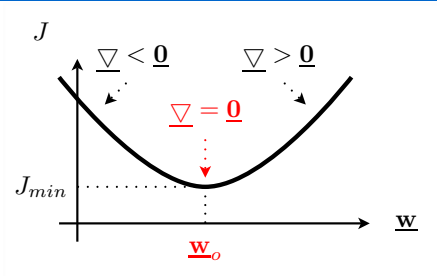
Problem: Optimal Wiener involves \mathbf{R}_x^{-1}

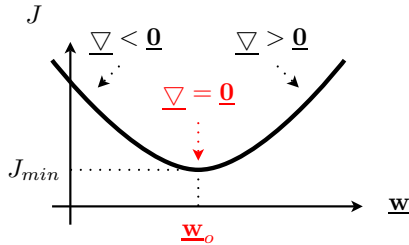
To avoid this inversion, estimate optimum *iteratively*

Problem: Optimal Wiener involves \mathbf{R}_x^{-1}

To avoid this inversion, estimate optimum *iteratively*

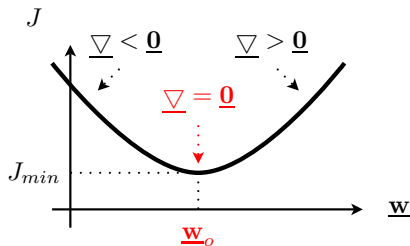
Goal: Decrease J each new iteration





GD principle: Update in negative gradient direction

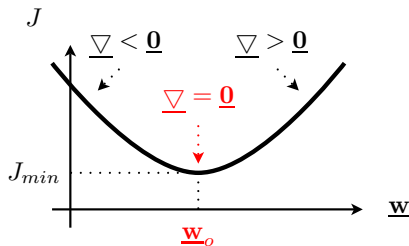
$$\Leftrightarrow \underline{\mathbf{w}} \doteq \underline{\mathbf{w}} - \alpha \underline{\nabla} \text{ with adaptation constant } \alpha \geq 0$$



GD principle: Update in negative gradient direction

$$\Leftrightarrow \underline{\mathbf{w}} \doteq \underline{\mathbf{w}} - \alpha \underline{\nabla} \text{ with adaptation constant } \alpha \geq 0$$

$$\text{With } \underline{\nabla} = -2(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$

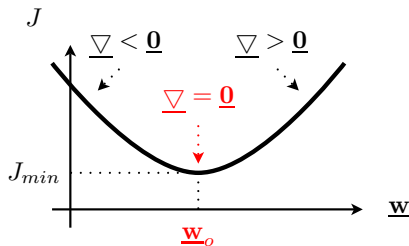


GD principle: **Update in negative gradient direction**

$$\Leftrightarrow \underline{\mathbf{w}} \doteq \underline{\mathbf{w}} - \alpha \underline{\nabla} \text{ with adaptation constant } \alpha \geq 0$$

With $\underline{\nabla} = -2(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k]) \Rightarrow$ **GD algorithm:**

$$\underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] + 2\alpha(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$



GD principle: **Update in negative gradient direction**

$$\Leftrightarrow \underline{\mathbf{w}} \doteq \underline{\mathbf{w}} - \alpha \underline{\nabla} \text{ with adaptation constant } \alpha \geq 0$$

With $\underline{\nabla} = -2(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k]) \Rightarrow$ **GD algorithm:**

$$\underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] + 2\alpha(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$

Notes: 1) No matrix inversion needed! 2) Usually $\underline{\mathbf{w}}[0] = \underline{\mathbf{0}}$

GD converges to Wiener solution:

$$\lim_{k \rightarrow \infty} \underline{\mathbf{w}}[k] \simeq \mathbf{R}_x^{-1} \cdot \underline{\mathbf{r}}_{yx}$$

GD converges to Wiener solution:

$$\lim_{k \rightarrow \infty} \underline{\mathbf{w}}[k] \simeq \mathbf{R}_x^{-1} \cdot \underline{\mathbf{r}}_{yx}$$

'Proof':

For $k \rightarrow \infty$ we have:

$$\underline{\mathbf{w}}[k+1] \simeq \underline{\mathbf{w}}[k] \simeq \underline{\mathbf{w}}[\infty]$$

GD converges to Wiener solution:

$$\lim_{k \rightarrow \infty} \underline{\mathbf{w}}[k] \simeq \mathbf{R}_x^{-1} \cdot \underline{\mathbf{r}}_{yx}$$

'Proof':

For $k \rightarrow \infty$ we have:

$$\begin{aligned} \underline{\mathbf{w}}[k+1] &\simeq \underline{\mathbf{w}}[k] \simeq \underline{\mathbf{w}}[\infty] \\ \text{GD} \Rightarrow \underline{\mathbf{w}}[\infty] &\simeq \underline{\mathbf{w}}[\infty] + 2\alpha(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[\infty]) \end{aligned}$$

GD converges to Wiener solution:

$$\lim_{k \rightarrow \infty} \underline{\mathbf{w}}[k] \simeq \mathbf{R}_x^{-1} \cdot \underline{\mathbf{r}}_{yx}$$

'Proof':

For $k \rightarrow \infty$ we have:

$$\begin{aligned} \underline{\mathbf{w}}[k+1] &\simeq \underline{\mathbf{w}}[k] \simeq \underline{\mathbf{w}}[\infty] \\ \text{GD} \Rightarrow \underline{\mathbf{w}}[\infty] &\simeq \underline{\mathbf{w}}[\infty] + 2\alpha(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[\infty]) \\ \Rightarrow \underline{\mathbf{w}}[\infty] &\simeq \mathbf{R}_x^{-1} \cdot \underline{\mathbf{r}}_{yx} \end{aligned}$$

GD converges to Wiener solution:

$$\lim_{k \rightarrow \infty} \underline{\mathbf{w}}[k] \simeq \mathbf{R}_x^{-1} \cdot \underline{\mathbf{r}}_{yx}$$

'Proof':

For $k \rightarrow \infty$ we have:

$$\begin{aligned} \underline{\mathbf{w}}[k+1] &\simeq \underline{\mathbf{w}}[k] \simeq \underline{\mathbf{w}}[\infty] \\ \text{GD} \Rightarrow \underline{\mathbf{w}}[\infty] &\simeq \underline{\mathbf{w}}[\infty] + 2\alpha(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[\infty]) \\ &\Rightarrow \underline{\mathbf{w}}[\infty] \simeq \mathbf{R}_x^{-1} \cdot \underline{\mathbf{r}}_{yx} \end{aligned}$$

For exact proof we need **stability analysis**

Define difference weight vector: $\underline{\mathbf{d}}[k] = \underline{\mathbf{w}}[k] - \underline{\mathbf{w}}_o$

Define difference weight vector: $\underline{\mathbf{d}}[k] = \underline{\mathbf{w}}[k] - \underline{\mathbf{w}}_o$

$$\underline{\mathbf{w}}[k + 1] = \underline{\mathbf{w}}[k] + 2\alpha(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$

Define difference weight vector: $\underline{\mathbf{d}}[k] = \underline{\mathbf{w}}[k] - \underline{\mathbf{w}}_o$

$$\begin{aligned}\underline{\mathbf{w}}[k+1] &= \underline{\mathbf{w}}[k] + 2\alpha(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k]) \\ \underline{\mathbf{w}}[k+1] - \underline{\mathbf{w}}_o &= (\mathbf{I} - 2\alpha \mathbf{R}_x) \cdot \underline{\mathbf{w}}[k] - \underline{\mathbf{w}}_o + 2\alpha \underline{\mathbf{r}}_{yx}\end{aligned}$$

Define difference weight vector: $\underline{\mathbf{d}}[k] = \underline{\mathbf{w}}[k] - \underline{\mathbf{w}}_o$

$$\begin{aligned}\underline{\mathbf{w}}[k+1] &= \underline{\mathbf{w}}[k] + 2\alpha(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k]) \\ \underline{\mathbf{w}}[k+1] - \underline{\mathbf{w}}_o &= (\mathbf{I} - 2\alpha \mathbf{R}_x) \cdot \underline{\mathbf{w}}[k] - \underline{\mathbf{w}}_o + 2\alpha \underline{\mathbf{r}}_{yx} \\ \Rightarrow \underline{\mathbf{d}}[k+1] &= (\mathbf{I} - 2\alpha \mathbf{R}_x) \cdot \underline{\mathbf{d}}[k]\end{aligned}$$

Define difference weight vector: $\underline{\mathbf{d}}[k] = \underline{\mathbf{w}}[k] - \underline{\mathbf{w}}_o$

$$\begin{aligned}\underline{\mathbf{w}}[k+1] &= \underline{\mathbf{w}}[k] + 2\alpha(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k]) \\ \underline{\mathbf{w}}[k+1] - \underline{\mathbf{w}}_o &= (\mathbf{I} - 2\alpha \mathbf{R}_x) \cdot \underline{\mathbf{w}}[k] - \underline{\mathbf{w}}_o + 2\alpha \underline{\mathbf{r}}_{yx} \\ \Rightarrow \underline{\mathbf{d}}[k+1] &= (\mathbf{I} - 2\alpha \mathbf{R}_x) \cdot \underline{\mathbf{d}}[k]\end{aligned}$$

Recursion:

$$\underline{\mathbf{d}}[k] = (\mathbf{I} - 2\alpha \mathbf{R}_x) \cdot \underline{\mathbf{d}}[k-1] = \dots = (\mathbf{I} - 2\alpha \mathbf{R}_x)^k \cdot \underline{\mathbf{d}}[0]$$

Define difference weight vector: $\underline{\mathbf{d}}[k] = \underline{\mathbf{w}}[k] - \underline{\mathbf{w}}_o$

$$\begin{aligned}\underline{\mathbf{w}}[k+1] &= \underline{\mathbf{w}}[k] + 2\alpha(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k]) \\ \underline{\mathbf{w}}[k+1] - \underline{\mathbf{w}}_o &= (\mathbf{I} - 2\alpha \mathbf{R}_x) \cdot \underline{\mathbf{w}}[k] - \underline{\mathbf{w}}_o + 2\alpha \underline{\mathbf{r}}_{yx} \\ \Rightarrow \underline{\mathbf{d}}[k+1] &= (\mathbf{I} - 2\alpha \mathbf{R}_x) \cdot \underline{\mathbf{d}}[k]\end{aligned}$$

Recursion:

$$\underline{\mathbf{d}}[k] = (\mathbf{I} - 2\alpha \mathbf{R}_x) \cdot \underline{\mathbf{d}}[k-1] = \dots = (\mathbf{I} - 2\alpha \mathbf{R}_x)^k \cdot \underline{\mathbf{d}}[0]$$

Converges if: $\lim_{k \rightarrow \infty} (\mathbf{I} - 2\alpha \mathbf{R}_x)^k = \mathbf{0}$

Define difference weight vector: $\underline{\mathbf{d}}[k] = \underline{\mathbf{w}}[k] - \underline{\mathbf{w}}_o$

$$\begin{aligned}\underline{\mathbf{w}}[k+1] &= \underline{\mathbf{w}}[k] + 2\alpha(\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k]) \\ \underline{\mathbf{w}}[k+1] - \underline{\mathbf{w}}_o &= (\mathbf{I} - 2\alpha \mathbf{R}_x) \cdot \underline{\mathbf{w}}[k] - \underline{\mathbf{w}}_o + 2\alpha \underline{\mathbf{r}}_{yx} \\ \Rightarrow \underline{\mathbf{d}}[k+1] &= (\mathbf{I} - 2\alpha \mathbf{R}_x) \cdot \underline{\mathbf{d}}[k]\end{aligned}$$

Recursion:

$$\underline{\mathbf{d}}[k] = (\mathbf{I} - 2\alpha \mathbf{R}_x) \cdot \underline{\mathbf{d}}[k-1] = \dots = (\mathbf{I} - 2\alpha \mathbf{R}_x)^k \cdot \underline{\mathbf{d}}[0]$$

Converges if: $\lim_{k \rightarrow \infty} (\mathbf{I} - 2\alpha \mathbf{R}_x)^k = \mathbf{0}$

Note:

When stable $\Rightarrow \underline{\mathbf{d}}[\infty] = \mathbf{0} \Rightarrow \underline{\mathbf{w}}[\infty] \simeq \mathbf{Wiener}$

How do weights converge:

How do weights converge:

Use eigenvalue decomposition (see Appendix):

How do weights converge:

Use eigenvalue decomposition (see Appendix):

With $Q^h \cdot Q = Q \cdot Q^h = I$ and $R_x = Q \Lambda Q^h$

How do weights converge:

Use eigenvalue decomposition (see Appendix):

With $\mathbf{Q}^h \cdot \mathbf{Q} = \mathbf{Q} \cdot \mathbf{Q}^h = \mathbf{I}$ and $\mathbf{R}_x = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^h$

$$\begin{aligned}\Rightarrow (\mathbf{I} - 2\alpha\mathbf{R}_x)^k &= (\mathbf{Q}\mathbf{Q}^h - 2\alpha\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^h)^k \\ &= \mathbf{Q}(\mathbf{I} - 2\alpha\mathbf{\Lambda})^k\mathbf{Q}^h\end{aligned}$$

How do weights converge:

Use eigenvalue decomposition (see Appendix):

With $\mathbf{Q}^h \cdot \mathbf{Q} = \mathbf{Q} \cdot \mathbf{Q}^h = \mathbf{I}$ and $\mathbf{R}_x = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^h$

$$\begin{aligned}\Rightarrow (\mathbf{I} - 2\alpha\mathbf{R}_x)^k &= (\mathbf{Q}\mathbf{Q}^h - 2\alpha\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^h)^k \\ &= \mathbf{Q}(\mathbf{I} - 2\alpha\mathbf{\Lambda})^k\mathbf{Q}^h\end{aligned}$$

Change of variables: $\underline{\mathbf{D}}[k] = \mathbf{Q}^h \cdot \underline{\mathbf{d}}[k]$

How do weights converge:

Use eigenvalue decomposition (see Appendix):

With $\mathbf{Q}^h \cdot \mathbf{Q} = \mathbf{Q} \cdot \mathbf{Q}^h = \mathbf{I}$ and $\mathbf{R}_x = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^h$

$$\begin{aligned}\Rightarrow (\mathbf{I} - 2\alpha\mathbf{R}_x)^k &= (\mathbf{Q}\mathbf{Q}^h - 2\alpha\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^h)^k \\ &= \mathbf{Q}(\mathbf{I} - 2\alpha\mathbf{\Lambda})^k\mathbf{Q}^h\end{aligned}$$

Change of variables: $\underline{\mathbf{D}}[k] = \mathbf{Q}^h \cdot \underline{\mathbf{d}}[k]$

$$\underline{\mathbf{d}}[k] = (\mathbf{I} - 2\alpha\mathbf{R}_x)^k \underline{\mathbf{d}}[0] \Rightarrow \underline{\mathbf{D}}[k] = (\mathbf{I} - 2\alpha\mathbf{\Lambda})^k \underline{\mathbf{D}}[0]$$

How do weights converge:

Use eigenvalue decomposition (see Appendix):

With $\mathbf{Q}^h \cdot \mathbf{Q} = \mathbf{Q} \cdot \mathbf{Q}^h = \mathbf{I}$ and $\mathbf{R}_x = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^h$

$$\begin{aligned}\Rightarrow (\mathbf{I} - 2\alpha\mathbf{R}_x)^k &= (\mathbf{Q}\mathbf{Q}^h - 2\alpha\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^h)^k \\ &= \mathbf{Q}(\mathbf{I} - 2\alpha\mathbf{\Lambda})^k\mathbf{Q}^h\end{aligned}$$

Change of variables: $\underline{\mathbf{D}}[k] = \mathbf{Q}^h \cdot \underline{\mathbf{d}}[k]$

$$\underline{\mathbf{d}}[k] = (\mathbf{I} - 2\alpha\mathbf{R}_x)^k \underline{\mathbf{d}}[0] \Rightarrow \underline{\mathbf{D}}[k] = (\mathbf{I} - 2\alpha\mathbf{\Lambda})^k \underline{\mathbf{D}}[0]$$

Recursion stable if: $\lim_{k \rightarrow \infty} (\mathbf{I} - 2\alpha\mathbf{\Lambda})^k = \mathbf{0}$

Recursion stable if: $\lim_{k \rightarrow \infty} (\mathbf{I} - 2\alpha\mathbf{\Lambda})^k = \mathbf{0}$

Recursion stable if: $\lim_{k \rightarrow \infty} (\mathbf{I} - 2\alpha\mathbf{\Lambda})^k = \mathbf{0}$

Both matrices \mathbf{I} and $\mathbf{\Lambda}$ diagonal

Recursion stable if: $\lim_{k \rightarrow \infty} (\mathbf{I} - 2\alpha\mathbf{\Lambda})^k = \mathbf{0}$

Both matrices \mathbf{I} and $\mathbf{\Lambda}$ diagonal \Rightarrow Stable if:

$$|1 - 2\alpha\lambda_i| < 1 \quad \Leftrightarrow \quad 0 < \alpha < \frac{1}{\lambda_i} \quad \text{for } i = 0, 1, \dots, N - 1$$

Recursion stable if: $\lim_{k \rightarrow \infty} (\mathbf{I} - 2\alpha\mathbf{\Lambda})^k = \mathbf{0}$

Both matrices \mathbf{I} and $\mathbf{\Lambda}$ diagonal \Rightarrow Stable if:

$$|1 - 2\alpha\lambda_i| < 1 \quad \Leftrightarrow \quad 0 < \alpha < \frac{1}{\lambda_i} \quad \text{for } i = 0, 1, \dots, N - 1$$

Thus GD algorithm stable if: $0 < \alpha < \frac{1}{\lambda_{max}}$

Recursion stable if: $\lim_{k \rightarrow \infty} (\mathbf{I} - 2\alpha\mathbf{\Lambda})^k = \mathbf{0}$

Both matrices \mathbf{I} and $\mathbf{\Lambda}$ diagonal \Rightarrow Stable if:

$$|1 - 2\alpha\lambda_i| < 1 \quad \Leftrightarrow \quad 0 < \alpha < \frac{1}{\lambda_i} \quad \text{for } i = 0, 1, \dots, N - 1$$

Thus GD algorithm stable if: $0 < \alpha < \frac{1}{\lambda_{max}}$

For adaptation constant α in this region:

$$\lim_{k \rightarrow \infty} \underline{\mathbf{w}}[k] = \underline{\mathbf{w}}_o = \mathbf{R}_x^{-1} \cdot \underline{\mathbf{r}}_{yx}$$

Recursion stable if: $\lim_{k \rightarrow \infty} (\mathbf{I} - 2\alpha\mathbf{\Lambda})^k = \mathbf{0}$

Both matrices \mathbf{I} and $\mathbf{\Lambda}$ diagonal \Rightarrow Stable if:

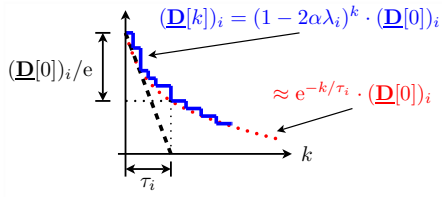
$$|1 - 2\alpha\lambda_i| < 1 \quad \Leftrightarrow \quad 0 < \alpha < \frac{1}{\lambda_i} \quad \text{for } i = 0, 1, \dots, N - 1$$

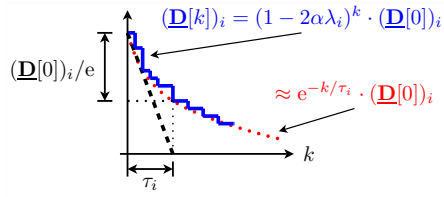
Thus GD algorithm stable if: $0 < \alpha < \frac{1}{\lambda_{max}}$

For adaptation constant α in this region:

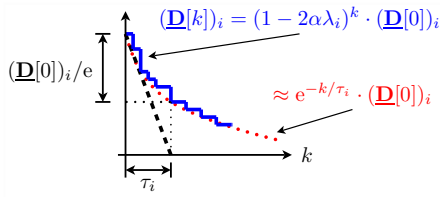
$$\lim_{k \rightarrow \infty} \underline{\mathbf{w}}[k] = \underline{\mathbf{w}}_o = \mathbf{R}_x^{-1} \cdot \underline{\mathbf{r}}_{yx}$$

$$J_{\underline{\mathbf{w}}=\underline{\mathbf{w}}_o} = E\{e^2[k]\} = J_{min} = E\{y^2\} - \underline{\mathbf{r}}_{yx}^t \mathbf{R}_x^{-1} \underline{\mathbf{r}}_{yx}$$



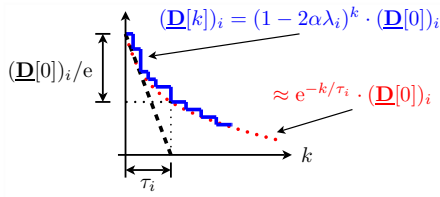


$$e^{-k/\tau_i} \cdot (\underline{D}[0])_i \approx (1 - 2\alpha\lambda_i)^k \cdot (\underline{D}[0])_i \Rightarrow$$



$$e^{-k/\tau_i} \cdot (\underline{D}[0])_i \approx (1 - 2\alpha\lambda_i)^k \cdot (\underline{D}[0])_i \Rightarrow$$

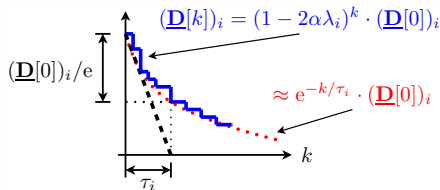
Average behavior: $\tau_{av,i} = \frac{-1}{\ln(1 - 2\alpha\lambda_i)}$



$$e^{-k/\tau_i} \cdot (\underline{\mathbf{D}}[0])_i \approx (1 - 2\alpha\lambda_i)^k \cdot (\underline{\mathbf{D}}[0])_i \Rightarrow$$

Average behavior: $\tau_{av,i} = \frac{-1}{\ln(1 - 2\alpha\lambda_i)} \Rightarrow \text{For small } \alpha$

$$\tau_{av,i} \approx \frac{1}{2\alpha\lambda_i}$$



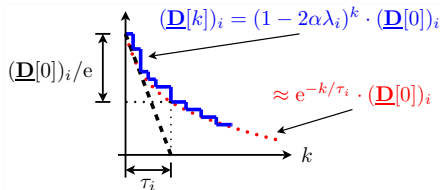
$$e^{-k/\tau_i} \cdot (\underline{D}[0])_i \approx (1 - 2\alpha\lambda_i)^k \cdot (\underline{D}[0])_i \Rightarrow$$

Average behavior: $\tau_{av,i} = \frac{-1}{\ln(1 - 2\alpha\lambda_i)} \Rightarrow \text{For small } \alpha$

$$\tau_{av,i} \approx \frac{1}{2\alpha\lambda_i}$$

Note:

Overall time constant depends on eigenvalue spread $\Gamma_x = \lambda_{max}/\lambda_{min}$.
Thus, the larger Γ_x the longer it takes for adaptation.



$$e^{-k/\tau_i} \cdot (\underline{D}[0])_i \approx (1 - 2\alpha\lambda_i)^k \cdot (\underline{D}[0])_i \Rightarrow$$

Average behavior: $\tau_{av,i} = \frac{-1}{\ln(1 - 2\alpha\lambda_i)} \Rightarrow \text{For small } \alpha$

$$\tau_{av,i} \approx \frac{1}{2\alpha\lambda_i}$$

Note:

Overall time constant depends on eigenvalue spread $\Gamma_x = \lambda_{max}/\lambda_{min}$.
Thus, the larger Γ_x the longer it takes for adaptation.

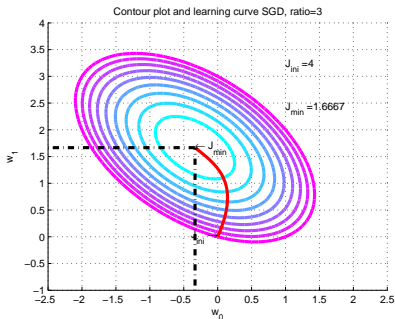
Q: What happens for white noise?

/department of electrical engineering

Example with $\Gamma_x = \lambda_{max}/\lambda_{min} = 3$

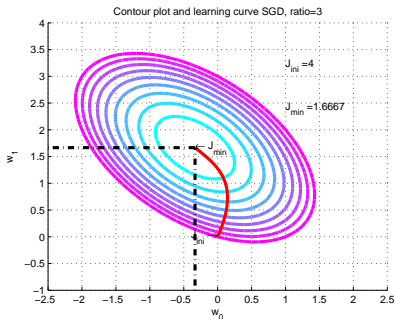
Example with $\Gamma_x = \lambda_{max}/\lambda_{min} = 3$

Learning curve in contour plot J

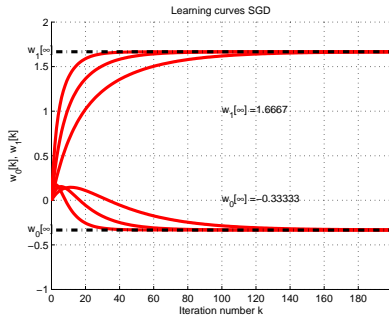


Example with $\Gamma_x = \lambda_{max}/\lambda_{min} = 3$

Learning curve in contour plot J



Learning curves for different α



Motivation: GD not practical. Gradient assumes known \mathbf{R}_x and \mathbf{r}_{yx}

Motivation: GD not practical. Gradient assumes known \mathbf{R}_x and \mathbf{r}_{yx}

LMS principle: Use instantaneous estimate of gradient:

$$\begin{aligned}\hat{\underline{\nabla}}[k] &= -2 (y[k]\underline{\mathbf{x}}[k] - \underline{\mathbf{x}}[k]\underline{\mathbf{x}}^t[k]\underline{\mathbf{w}}[k]) \\ &= -2\underline{\mathbf{x}}[k] (y[k] - \underline{\mathbf{x}}^t[k]\underline{\mathbf{w}}[k]) = -2\underline{\mathbf{x}}[k]e[k]\end{aligned}$$

Motivation: GD not practical. Gradient assumes known \mathbf{R}_x and \mathbf{r}_{yx}

LMS principle: Use instantaneous estimate of gradient:

$$\begin{aligned}\hat{\nabla}[k] &= -2 (y[k]\underline{\mathbf{x}}[k] - \underline{\mathbf{x}}[k]\underline{\mathbf{x}}^t[k]\underline{\mathbf{w}}[k]) \\ &= -2\underline{\mathbf{x}}[k] (y[k] - \underline{\mathbf{x}}^t[k]\underline{\mathbf{w}}[k]) = -2\underline{\mathbf{x}}[k]e[k]\end{aligned}$$

With $\underline{\mathbf{w}} \doteq \underline{\mathbf{w}} - \alpha \hat{\nabla} \Rightarrow$ LMS algorithm (Widrow, 1975):

$$k = 0 : \underline{\mathbf{w}}[0] = \underline{\mathbf{0}} \text{ (usually)}$$

$$k > 0 : \hat{y}[k] = \underline{\mathbf{w}}^t[k] \cdot \underline{\mathbf{x}}[k]$$

$$e[k] = y[k] - \hat{y}[k]$$

$$\underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] + 2\alpha\underline{\mathbf{x}}[k]e[k]$$

Motivation: GD not practical. Gradient assumes known \mathbf{R}_x and \mathbf{r}_{yx}

LMS principle: Use instantaneous estimate of gradient:

$$\begin{aligned}\hat{\nabla}[k] &= -2 (y[k]\underline{\mathbf{x}}[k] - \underline{\mathbf{x}}[k]\underline{\mathbf{x}}^t[k]\underline{\mathbf{w}}[k]) \\ &= -2\underline{\mathbf{x}}[k] (y[k] - \underline{\mathbf{x}}^t[k]\underline{\mathbf{w}}[k]) = -2\underline{\mathbf{x}}[k]e[k]\end{aligned}$$

With $\underline{\mathbf{w}} \doteq \underline{\mathbf{w}} - \alpha \hat{\nabla} \Rightarrow$ LMS algorithm (Widrow, 1975):

$$k = 0 : \underline{\mathbf{w}}[0] = \underline{\mathbf{0}} \text{ (usually)}$$

$$k > 0 : \hat{y}[k] = \underline{\mathbf{w}}^t[k] \cdot \underline{\mathbf{x}}[k]$$

$$e[k] = y[k] - \hat{y}[k]$$

$$\underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] + 2\alpha \underline{\mathbf{x}}[k]e[k]$$

Note: $\underline{\mathbf{w}}^t[k] \cdot \underline{\mathbf{x}}[k]$ is "convolution" and $\underline{\mathbf{x}}[k]e[k]$ "correlation"

- ▶ Convergence of LMS strongly depends on α ; "optimal" choice depends on amplitude of signals.

- ▶ Convergence of LMS strongly depends on α ; "optimal" choice depends on amplitude of signals.
- ▶ **NLMS:** LMS with normalization by $\sigma_x^2 = E\{x^2[k]\}$:

$$\underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] + \frac{2\alpha}{\sigma_x^2} \underline{\mathbf{x}}[k] e[k]$$

- ▶ Convergence of LMS strongly depends on α ; "optimal" choice depends on amplitude of signals.
- ▶ **NLMS**: LMS with normalization by $\sigma_x^2 = E\{x^2[k]\}$:

$$\underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] + \frac{2\alpha}{\sigma_x^2} \underline{\mathbf{x}}[k]e[k]$$

In practice $\hat{\sigma}_x^2[k] \Rightarrow$ time-varying step size. E.g.:

- $\hat{\sigma}_x^2[k] = \beta \hat{\sigma}_x^2[k-1] + (1 - \beta) \frac{\underline{\mathbf{x}}^t[k]\underline{\mathbf{x}}[k]}{N}$ with $0 < \beta < 1$
- $\hat{\sigma}_x^2[k] = \frac{\underline{\mathbf{x}}^t[k]\underline{\mathbf{x}}[k]}{N} + \epsilon$ with ϵ some small constant

Convergence gradient based algorithms depends on coloration input:

$$\underline{\nabla} = -2 (\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$

Convergence gradient based algorithms depends on coloration input:

$$\underline{\nabla} = -2 (\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$

Solution Newton: $\underline{\mathbf{w}}[k + 1] = \underline{\mathbf{w}}[k] - \alpha \mathbf{R}_x^{-1} \underline{\nabla} \Rightarrow$

$$\underline{\mathbf{w}}[k + 1] = \underline{\mathbf{w}}[k] + 2\alpha \mathbf{R}_x^{-1} \cdot (\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$

Convergence gradient based algorithms depends on coloration input:

$$\underline{\nabla} = -2 (\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$

Solution Newton: $\underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] - \alpha \mathbf{R}_x^{-1} \underline{\nabla} \Rightarrow$

$$\underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] + 2\alpha \mathbf{R}_x^{-1} \cdot (\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$

Convergence Newton:

$$\underline{\mathbf{d}}[k+1] = (\mathbf{I} - 2\alpha \mathbf{R}_x^{-1} \mathbf{R}_x) \underline{\mathbf{d}}[k] = (1-2\alpha) \underline{\mathbf{d}}[k] \Rightarrow \text{Convergence } 0 < \alpha < 1$$

Convergence gradient based algorithms depends on coloration input:

$$\underline{\nabla} = -2 (\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$

Solution Newton: $\underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] - \alpha \mathbf{R}_x^{-1} \underline{\nabla} \Rightarrow$

$$\underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] + 2\alpha \mathbf{R}_x^{-1} \cdot (\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$

Convergence Newton:

$$\underline{\mathbf{d}}[k+1] = (\mathbf{I} - 2\alpha \mathbf{R}_x^{-1} \mathbf{R}_x) \underline{\mathbf{d}}[k] = (1-2\alpha) \underline{\mathbf{d}}[k] \Rightarrow \text{Convergence } 0 < \alpha < 1$$

Notes:

- ▶ \mathbf{R}_x^{-1} causes whitening of input process

Convergence gradient based algorithms depends on coloration input:

$$\underline{\nabla} = -2 (\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$

Solution Newton: $\underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] - \alpha \mathbf{R}_x^{-1} \underline{\nabla} \Rightarrow$

$$\underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] + 2\alpha \mathbf{R}_x^{-1} \cdot (\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$

Convergence Newton:

$$\underline{\mathbf{d}}[k+1] = (\mathbf{I} - 2\alpha \mathbf{R}_x^{-1} \mathbf{R}_x) \underline{\mathbf{d}}[k] = (1-2\alpha) \underline{\mathbf{d}}[k] \Rightarrow \text{Convergence } 0 < \alpha < 1$$

Notes:

- ▶ \mathbf{R}_x^{-1} causes whitening of input process
- ▶ All weights have same convergence (in contrast to LMS, GD)

Convergence gradient based algorithms depends on coloration input:

$$\underline{\nabla} = -2 (\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$

Solution Newton: $\underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] - \alpha \mathbf{R}_x^{-1} \underline{\nabla} \Rightarrow$

$$\underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] + 2\alpha \mathbf{R}_x^{-1} \cdot (\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$

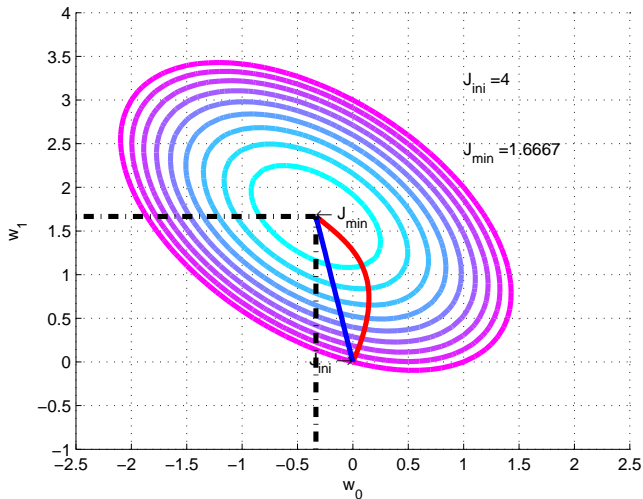
Convergence Newton:

$$\underline{\mathbf{d}}[k+1] = (\mathbf{I} - 2\alpha \mathbf{R}_x^{-1} \mathbf{R}_x) \underline{\mathbf{d}}[k] = (1-2\alpha) \underline{\mathbf{d}}[k] \Rightarrow \text{Convergence } 0 < \alpha < 1$$

Notes:

- ▶ \mathbf{R}_x^{-1} causes whitening of input process
- ▶ All weights have same convergence (in contrast to LMS, GD)
- ▶ Newton \equiv GD with white noise input!

Learning curves in contour plot: Newton vs. GD



Autocorrelation matrix \mathbf{R}_x :

Autocorrelation matrix \mathbf{R}_x :

- ▶ (In general) not known in advance
- ▶ May change during time (non-stationary process)
- ▶ Inversion is expensive (many MIPS)

Autocorrelation matrix \mathbf{R}_x :

- ▶ (In general) not known in advance
 - ▶ May change during time (non-stationary process)
 - ▶ Inversion is expensive (many MIPS)
- ⇒ Complexity Newton algorithm huge
- ⇒ Need for efficient solution with estimate of \mathbf{R}_x
- ⇒ Different algorithms, e.g. RLS.

For data block length L fixed, Least Squares problem becomes:

$$\min_{\underline{\mathbf{w}}[k]} |\underline{\mathbf{y}}[k] - \mathbf{X}[k] \cdot \underline{\mathbf{w}}[k]|^2 \Rightarrow \underline{\mathbf{w}}_{LS}[k] = (\mathbf{X}^t[k] \mathbf{X}[k])^{-1} (\mathbf{X}^t[k] \underline{\mathbf{y}}[k])$$

For data block length L fixed, Least Squares problem becomes:

$$\min_{\underline{\mathbf{w}}[k]} |\underline{\mathbf{y}}[k] - \mathbf{X}[k] \cdot \underline{\mathbf{w}}[k]|^2 \Rightarrow \underline{\mathbf{w}}_{LS}[k] = (\mathbf{X}^t[k] \mathbf{X}[k])^{-1} (\mathbf{X}^t[k] \underline{\mathbf{y}}[k])$$

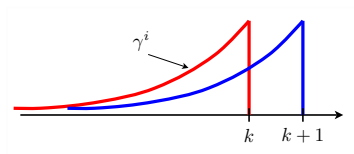
RLS: Find efficient recursive solution for LS problem from $k \rightarrow k + 1$

For data block length L fixed, Least Squares problem becomes:

$$\min_{\underline{\mathbf{w}}[k]} |\underline{\mathbf{y}}[k] - \mathbf{X}[k] \cdot \underline{\mathbf{w}}[k]|^2 \Rightarrow \underline{\mathbf{w}}_{LS}[k] = (\mathbf{X}^t[k] \mathbf{X}[k])^{-1} (\mathbf{X}^t[k] \underline{\mathbf{y}}[k])$$

RLS: Find efficient recursive solution for LS problem from $k \rightarrow k + 1$

Use exponential sliding window: Scale down data by factor γ



Forgetting factor : $0 < \gamma < 1$

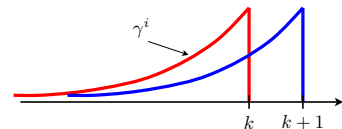
'Memory' : $\frac{1}{1 - \gamma}$

For data block length L fixed, Least Squares problem becomes:

$$\min_{\underline{\mathbf{w}}[k]} |\underline{\mathbf{y}}[k] - \mathbf{X}[k] \cdot \underline{\mathbf{w}}[k]|^2 \Rightarrow \underline{\mathbf{w}}_{LS}[k] = (\mathbf{X}^t[k] \mathbf{X}[k])^{-1} (\mathbf{X}^t[k] \underline{\mathbf{y}}[k])$$

RLS: Find efficient recursive solution for LS problem from $k \rightarrow k + 1$

Use exponential sliding window: Scale down data by factor γ



Forgetting factor : $0 < \gamma < 1$

'Memory' : $\frac{1}{1 - \gamma}$

$$\mathbf{X}[k] = \begin{pmatrix} \gamma^0 \underline{\mathbf{x}}^t[k] \\ \vdots \\ \gamma^i \underline{\mathbf{x}}^t[k - i] \\ \vdots \\ \gamma^k \underline{\mathbf{x}}^t[0] \end{pmatrix} \quad \text{and} \quad \underline{\mathbf{y}}[k] = \begin{pmatrix} \gamma^0 y[k] \\ \vdots \\ \gamma^i y[k - i] \\ \vdots \\ \gamma^k y[0] \end{pmatrix}$$

Initialization: $\underline{\bar{\mathbf{r}}}_{yx}[0] = \underline{\mathbf{0}}$; $\overline{\mathbf{R}}_x^{-1}[0] = \delta^{-1} \mathbf{I}$ with δ large

Initialization: $\bar{\mathbf{r}}_{yx}[0] = \mathbf{0}$; $\bar{\mathbf{R}}_x^{-1}[0] = \delta^{-1}\mathbf{I}$ with δ large

For $k \geq 0$:

Initialization: $\underline{\mathbf{r}}_{yx}[0] = \underline{\mathbf{0}}$; $\overline{\mathbf{R}}_x^{-1}[0] = \delta^{-1} \mathbf{I}$ with δ large

For $k \geq 0$:

$$\overline{\mathbf{R}}_x^{-1}[k+1] = \gamma^{-2} \left(\overline{\mathbf{R}}_x^{-1}[k] - \underline{\mathbf{g}}[k+1] \cdot \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \right)$$

Initialization: $\underline{\mathbf{r}}_{yx}[0] = \underline{\mathbf{0}}$; $\overline{\mathbf{R}}_x^{-1}[0] = \delta^{-1} \mathbf{I}$ with δ large

$$\begin{aligned} \text{For } k \geq 0: \underline{\mathbf{g}}[k+1] &= \frac{\overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]}{\gamma^2 + \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]} \\ \overline{\mathbf{R}}_x^{-1}[k+1] &= \gamma^{-2} \left(\overline{\mathbf{R}}_x^{-1}[k] - \underline{\mathbf{g}}[k+1] \cdot \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \right) \end{aligned}$$

Initialization: $\underline{\mathbf{r}}_{yx}[0] = \underline{\mathbf{0}}$; $\overline{\mathbf{R}}_x^{-1}[0] = \delta^{-1} \mathbf{I}$ with δ large

$$\begin{aligned} \text{For } k \geq 0: \underline{\mathbf{g}}[k+1] &= \frac{\overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]}{\gamma^2 + \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]} \\ \overline{\mathbf{R}}_x^{-1}[k+1] &= \gamma^{-2} \left(\overline{\mathbf{R}}_x^{-1}[k] - \underline{\mathbf{g}}[k+1] \cdot \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \right) \\ \underline{\mathbf{r}}_{yx}[k+1] &= \gamma^2 \underline{\mathbf{r}}_{yx}[k] + \underline{\mathbf{x}}^t[k+1] \cdot y[k+1] \end{aligned}$$

Initialization: $\underline{\mathbf{r}}_{yx}[0] = \underline{\mathbf{0}}$; $\overline{\mathbf{R}}_x^{-1}[0] = \delta^{-1} \mathbf{I}$ with δ large

$$\begin{aligned}\text{For } k \geq 0: \underline{\mathbf{g}}[k+1] &= \frac{\overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]}{\gamma^2 + \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]} \\ \overline{\mathbf{R}}_x^{-1}[k+1] &= \gamma^{-2} \left(\overline{\mathbf{R}}_x^{-1}[k] - \underline{\mathbf{g}}[k+1] \cdot \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \right) \\ \underline{\mathbf{r}}_{yx}[k+1] &= \gamma^2 \underline{\mathbf{r}}_{yx}[k] + \underline{\mathbf{x}}^t[k+1] \cdot y[k+1] \\ \underline{\mathbf{w}}[k+1] &= \overline{\mathbf{R}}_x^{-1}[k+1] \cdot \underline{\mathbf{r}}_{yx}[k+1]\end{aligned}$$

Initialization: $\underline{\mathbf{r}}_{yx}[0] = \underline{\mathbf{0}}$; $\overline{\mathbf{R}}_x^{-1}[0] = \delta^{-1} \mathbf{I}$ with δ large

$$\begin{aligned}\text{For } k \geq 0: \underline{\mathbf{g}}[k+1] &= \frac{\overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]}{\gamma^2 + \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]} \\ \overline{\mathbf{R}}_x^{-1}[k+1] &= \gamma^{-2} \left(\overline{\mathbf{R}}_x^{-1}[k] - \underline{\mathbf{g}}[k+1] \cdot \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \right) \\ \underline{\mathbf{r}}_{yx}[k+1] &= \gamma^2 \underline{\mathbf{r}}_{yx}[k] + \underline{\mathbf{x}}^t[k+1] \cdot y[k+1] \\ \underline{\mathbf{w}}[k+1] &= \overline{\mathbf{R}}_x^{-1}[k+1] \cdot \underline{\mathbf{r}}_{yx}[k+1]\end{aligned}$$

Notes:

► $\underline{\mathbf{w}}[\infty] \rightarrow \underline{\mathbf{w}}_o$

Initialization: $\underline{\mathbf{r}}_{yx}[0] = \underline{\mathbf{0}}$; $\overline{\mathbf{R}}_x^{-1}[0] = \delta^{-1} \mathbf{I}$ with δ large

$$\begin{aligned}\text{For } k \geq 0: \underline{\mathbf{g}}[k+1] &= \frac{\overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]}{\gamma^2 + \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]} \\ \overline{\mathbf{R}}_x^{-1}[k+1] &= \gamma^{-2} \left(\overline{\mathbf{R}}_x^{-1}[k] - \underline{\mathbf{g}}[k+1] \cdot \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \right) \\ \underline{\mathbf{r}}_{yx}[k+1] &= \gamma^2 \underline{\mathbf{r}}_{yx}[k] + \underline{\mathbf{x}}^t[k+1] \cdot y[k+1] \\ \underline{\mathbf{w}}[k+1] &= \overline{\mathbf{R}}_x^{-1}[k+1] \cdot \underline{\mathbf{r}}_{yx}[k+1]\end{aligned}$$

Notes:

- ▶ $\underline{\mathbf{w}}[\infty] \rightarrow \underline{\mathbf{w}}_o$
- ▶ Complexity $O(N^2)$ per time update

Initialization: $\underline{\mathbf{r}}_{yx}[0] = \underline{\mathbf{0}}$; $\overline{\mathbf{R}}_x^{-1}[0] = \delta^{-1} \mathbf{I}$ with δ large

$$\begin{aligned}\text{For } k \geq 0: \underline{\mathbf{g}}[k+1] &= \frac{\overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]}{\gamma^2 + \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]} \\ \overline{\mathbf{R}}_x^{-1}[k+1] &= \gamma^{-2} \left(\overline{\mathbf{R}}_x^{-1}[k] - \underline{\mathbf{g}}[k+1] \cdot \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \right) \\ \underline{\mathbf{r}}_{yx}[k+1] &= \gamma^2 \underline{\mathbf{r}}_{yx}[k] + \underline{\mathbf{x}}^t[k+1] \cdot y[k+1] \\ \underline{\mathbf{w}}[k+1] &= \overline{\mathbf{R}}_x^{-1}[k+1] \cdot \underline{\mathbf{r}}_{yx}[k+1]\end{aligned}$$

Notes:

- ▶ $\underline{\mathbf{w}}[\infty] \rightarrow \underline{\mathbf{w}}_o$
- ▶ Complexity $O(N^2)$ per time update
- ▶ Window length increases when time increases!

Initialization: $\underline{\mathbf{r}}_{yx}[0] = \underline{\mathbf{0}}$; $\overline{\mathbf{R}}_x^{-1}[0] = \delta^{-1} \mathbf{I}$ with δ large

$$\begin{aligned}\text{For } k \geq 0: \underline{\mathbf{g}}[k+1] &= \frac{\overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]}{\gamma^2 + \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]} \\ \overline{\mathbf{R}}_x^{-1}[k+1] &= \gamma^{-2} \left(\overline{\mathbf{R}}_x^{-1}[k] - \underline{\mathbf{g}}[k+1] \cdot \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \right) \\ \underline{\mathbf{r}}_{yx}[k+1] &= \gamma^2 \underline{\mathbf{r}}_{yx}[k] + \underline{\mathbf{x}}^t[k+1] \cdot y[k+1] \\ \underline{\mathbf{w}}[k+1] &= \overline{\mathbf{R}}_x^{-1}[k+1] \cdot \underline{\mathbf{r}}_{yx}[k+1]\end{aligned}$$

Notes:

- ▶ $\underline{\mathbf{w}}[\infty] \rightarrow \underline{\mathbf{w}}_o$
- ▶ Complexity $O(N^2)$ per time update
- ▶ Window length increases when time increases!
- ▶ Exhibits unstable roundoff error accumulation

Initialization: $\underline{\mathbf{r}}_{yx}[0] = \underline{\mathbf{0}}$; $\overline{\mathbf{R}}_x^{-1}[0] = \delta^{-1} \mathbf{I}$ with δ large

$$\begin{aligned}\text{For } k \geq 0: \underline{\mathbf{g}}[k+1] &= \frac{\overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]}{\gamma^2 + \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]} \\ \overline{\mathbf{R}}_x^{-1}[k+1] &= \gamma^{-2} \left(\overline{\mathbf{R}}_x^{-1}[k] - \underline{\mathbf{g}}[k+1] \cdot \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \right) \\ \underline{\mathbf{r}}_{yx}[k+1] &= \gamma^2 \underline{\mathbf{r}}_{yx}[k] + \underline{\mathbf{x}}^t[k+1] \cdot y[k+1] \\ \underline{\mathbf{w}}[k+1] &= \overline{\mathbf{R}}_x^{-1}[k+1] \cdot \underline{\mathbf{r}}_{yx}[k+1]\end{aligned}$$

Notes:

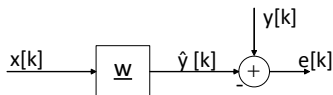
- ▶ $\underline{\mathbf{w}}[\infty] \rightarrow \underline{\mathbf{w}}_o$
- ▶ Complexity $O(N^2)$ per time update
- ▶ Window length increases when time increases!
- ▶ Exhibits unstable roundoff error accumulation
- ▶ RLS is basis for many practical algorithms

Initialization: $\underline{\mathbf{r}}_{yx}[0] = \underline{\mathbf{0}}$; $\overline{\mathbf{R}}_x^{-1}[0] = \delta^{-1} \mathbf{I}$ with δ large

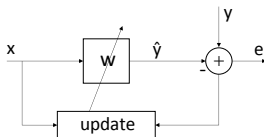
$$\begin{aligned}\text{For } k \geq 0: \underline{\mathbf{g}}[k+1] &= \frac{\overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]}{\gamma^2 + \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]} \\ \overline{\mathbf{R}}_x^{-1}[k+1] &= \gamma^{-2} \left(\overline{\mathbf{R}}_x^{-1}[k] - \underline{\mathbf{g}}[k+1] \cdot \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \right) \\ \underline{\mathbf{r}}_{yx}[k+1] &= \gamma^2 \underline{\mathbf{r}}_{yx}[k] + \underline{\mathbf{x}}^t[k+1] \cdot y[k+1] \\ \underline{\mathbf{w}}[k+1] &= \overline{\mathbf{R}}_x^{-1}[k+1] \cdot \underline{\mathbf{r}}_{yx}[k+1]\end{aligned}$$

Notes:

- ▶ $\underline{\mathbf{w}}[\infty] \rightarrow \underline{\mathbf{w}}_o$
- ▶ Complexity $O(N^2)$ per time update
- ▶ Window length increases when time increases!
- ▶ Exhibits unstable roundoff error accumulation
- ▶ RLS is basis for many practical algorithms
- ▶ Decorrelation takes place in algorithm



	MMSE	LS
Auto correlation	$\mathbf{R}_x = E\{\mathbf{x}[k] \cdot \mathbf{x}^t[k]\}$	$\overline{\mathbf{R}}_x = \mathbf{X}^t \cdot \mathbf{X}$
Cross correlation	$\mathbf{r}_{yx} = E\{y[k] \cdot \mathbf{x}[k]\}$	$\bar{\mathbf{r}}_{yx} = \mathbf{X}^t \cdot \mathbf{y}$
Error J	$E\{e^2[k]\}$	$\sum_{i=0}^{L-1} e^2[k-i]$
Criterion	$\min_{\mathbf{w}} \{E\{e^2[k]\}\}$	$\min_{\mathbf{w}} \mathbf{y} - \mathbf{X} \cdot \mathbf{w} ^2$
Opt. solution \mathbf{w}_o	$\mathbf{R}_x^{-1} \cdot \mathbf{r}_{yx}$	$\overline{\mathbf{R}}_x^{-1} \cdot \bar{\mathbf{r}}_{yx}$
Min. error J_{min}	$E\{y^2\} - \mathbf{r}_{yx}^t \mathbf{R}_x^{-1} \mathbf{r}_{yx}$	$\mathbf{y}^t \mathbf{y} - \bar{\mathbf{r}}_{yx}^t \overline{\mathbf{R}}_x^{-1} \bar{\mathbf{r}}_{yx}$



Simple adaptive algorithms (no decorrelation):

$$\text{GD} : \underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] + 2\alpha(\mathbf{r}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$

$$(\text{N})\text{LMS} : \underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] + \frac{2\alpha}{\hat{\sigma}_x^2} \mathbf{x}[k] e^*[k]$$

Algorithms with improved convergence:

$$\text{LMS/Newton} : \underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] + 2\alpha \mathbf{R}_x^{-1} \underline{\mathbf{x}}[k] r[k]$$

$$\text{Newton} : \underline{\mathbf{w}}[k+1] = \underline{\mathbf{w}}[k] + 2\alpha \mathbf{R}_x^{-1} \cdot (\underline{\mathbf{r}}_{yx} - \mathbf{R}_x \underline{\mathbf{w}}[k])$$

$$\text{RLS} : \underline{\mathbf{g}}[k+1] = \frac{\overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]}{\gamma^2 + \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \underline{\mathbf{x}}[k+1]}$$

$$\overline{\mathbf{R}}_x^{-1}[k+1] = \gamma^{-2} \left(\overline{\mathbf{R}}_x^{-1}[k] - \underline{\mathbf{g}}[k+1] \cdot \underline{\mathbf{x}}^t[k+1] \overline{\mathbf{R}}_x^{-1}[k] \right)$$

$$\underline{\mathbf{r}}_{yx}[k+1] = \gamma^2 \underline{\mathbf{r}}_{yx}[k] + \underline{\mathbf{x}}^t[k+1] \cdot y[k+1]$$

$$\underline{\mathbf{w}}[k+1] = \overline{\mathbf{R}}_x^{-1}[k+1] \cdot \underline{\mathbf{r}}_{yx}[k+1]$$

Appendix

Optimum Linear Filters & Adaptive Signal Processing

- ▶ Eigenvalue problem

Procedure: With eigenvalues λ_i and eigenvectors $\underline{\mathbf{q}}_i$:

$$\mathbf{R} \cdot \underline{\mathbf{q}}_i = \lambda_i \cdot \underline{\mathbf{q}}_i \Rightarrow (\mathbf{R} - \lambda_i \mathbf{I}) \cdot \underline{\mathbf{q}}_i = \underline{\mathbf{0}} \text{ for } i = 0, 1, \dots, N - 1$$

Procedure: With eigenvalues λ_i and eigenvectors $\underline{\mathbf{q}}_i$:

$$\mathbf{R} \cdot \underline{\mathbf{q}}_i = \lambda_i \cdot \underline{\mathbf{q}}_i \Rightarrow (\mathbf{R} - \lambda_i \mathbf{I}) \cdot \underline{\mathbf{q}}_i = \underline{\mathbf{0}} \text{ for } i = 0, 1, \dots, N-1$$

With $\mathbf{Q} = (\underline{\mathbf{q}}_0, \dots, \underline{\mathbf{q}}_{N-1})$ and $\mathbf{\Lambda} = \text{diag}\{\lambda_0, \dots, \lambda_{N-1}\}$

$$\mathbf{R} \cdot \mathbf{Q} = \mathbf{Q} \cdot \mathbf{\Lambda} \Rightarrow \mathbf{R} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1}$$

Procedure: With eigenvalues λ_i and eigenvectors \underline{q}_i :

$$\mathbf{R} \cdot \underline{q}_i = \lambda_i \cdot \underline{q}_i \Rightarrow (\mathbf{R} - \lambda_i \mathbf{I}) \cdot \underline{q}_i = \underline{0} \text{ for } i = 0, 1, \dots, N-1$$

With $\mathbf{Q} = (\underline{q}_0, \dots, \underline{q}_{N-1})$ and $\mathbf{\Lambda} = \text{diag}\{\lambda_0, \dots, \lambda_{N-1}\}$

$$\mathbf{R} \cdot \mathbf{Q} = \mathbf{Q} \cdot \mathbf{\Lambda} \Rightarrow \mathbf{R} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1}$$

Property: Eigenvectors \underline{q}_i orthogonal \Rightarrow

$$\mathbf{Q}^h \cdot \mathbf{Q} = \mathbf{Q} \cdot \mathbf{Q}^h = c \cdot \mathbf{I} \text{ with } c \text{ some constant}$$

Procedure: With eigenvalues λ_i and eigenvectors \underline{q}_i :

$$\mathbf{R} \cdot \underline{q}_i = \lambda_i \cdot \underline{q}_i \Rightarrow (\mathbf{R} - \lambda_i \mathbf{I}) \cdot \underline{q}_i = \underline{0} \text{ for } i = 0, 1, \dots, N-1$$

With $\mathbf{Q} = (\underline{q}_0, \dots, \underline{q}_{N-1})$ and $\mathbf{\Lambda} = \text{diag}\{\lambda_0, \dots, \lambda_{N-1}\}$

$$\mathbf{R} \cdot \mathbf{Q} = \mathbf{Q} \cdot \mathbf{\Lambda} \Rightarrow \mathbf{R} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1}$$

Property: Eigenvectors \underline{q}_i orthogonal \Rightarrow

$$\mathbf{Q}^h \cdot \mathbf{Q} = \mathbf{Q} \cdot \mathbf{Q}^h = c \cdot \mathbf{I} \text{ with } c \text{ some constant}$$

Main result:

Diagonalization:

$$\mathbf{Q}^h \mathbf{R} \mathbf{Q} = \mathbf{\Lambda} \Leftrightarrow \mathbf{R} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^h$$

Example MA(1):

$$x[k] = i[k] + ai[k - 1] \text{ with } E\{i[k]\} = 0 \text{ and } E\{i^2[k]\} = \sigma_i^2 \Rightarrow$$

Example MA(1):

$x[k] = i[k] + ai[k - 1]$ with $E\{i[k]\} = 0$ and $E\{i^2[k]\} = \sigma_i^2 \Rightarrow$

$\rho[0] = (1 + a^2)\sigma_i^2; \rho[1] = \rho[-1] = a\sigma_i^2; \rho[\tau] = 0$ for $|\tau| \geq 2$

Example MA(1):

$x[k] = i[k] + ai[k - 1]$ with $E\{i[k]\} = 0$ and $E\{i^2[k]\} = \sigma_i^2 \Rightarrow$

$\rho[0] = (1 + a^2)\sigma_i^2$; $\rho[1] = \rho[-1] = a\sigma_i^2$; $\rho[\tau] = 0$ for $|\tau| \geq 2$

Eigenvalues problem $\det(\mathbf{R} - \lambda\mathbf{I}) = 0$ for $N = 2$ (with $\gamma = \rho[1]/\rho[0]$):

$$\mathbf{\Lambda} = \begin{pmatrix} \lambda_0 & 0 \\ 0 & \lambda_1 \end{pmatrix} = \begin{pmatrix} 1 + \gamma & 0 \\ 0 & 1 - \gamma \end{pmatrix} ; \mathbf{Q} = (\underline{\mathbf{q}}_0, \underline{\mathbf{q}}_1) = c \cdot \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$$

Example MA(1):

$x[k] = i[k] + ai[k - 1]$ with $E\{i[k]\} = 0$ and $E\{i^2[k]\} = \sigma_i^2 \Rightarrow$

$$\rho[0] = (1 + a^2)\sigma_i^2; \rho[1] = \rho[-1] = a\sigma_i^2; \rho[\tau] = 0 \text{ for } |\tau| \geq 2$$

Eigenvalues problem $\det(\mathbf{R} - \lambda\mathbf{I}) = 0$ for $N = 2$ (with $\gamma = \rho[1]/\rho[0]$):

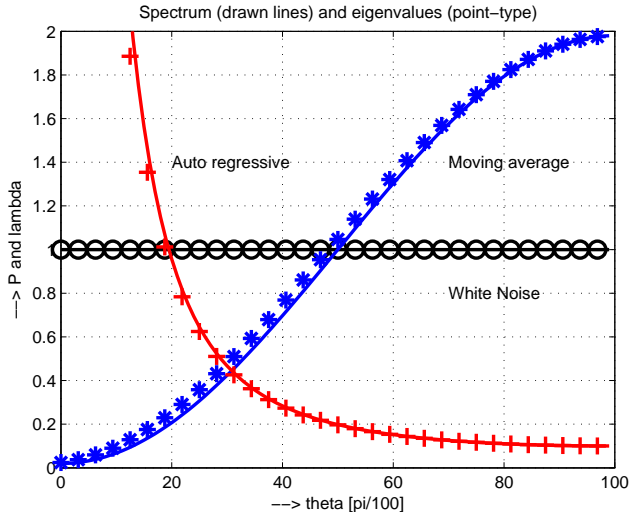
$$\mathbf{\Lambda} = \begin{pmatrix} \lambda_0 & 0 \\ 0 & \lambda_1 \end{pmatrix} = \begin{pmatrix} 1 + \gamma & 0 \\ 0 & 1 - \gamma \end{pmatrix} ; \mathbf{Q} = (\underline{\mathbf{q}}_0, \underline{\mathbf{q}}_1) = c \cdot \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$$

Notes:

- ▶ Vector $\underline{\mathbf{q}}_0$ orthogonal to $\underline{\mathbf{q}}_1$ since $\underline{\mathbf{q}}_0^t \cdot \underline{\mathbf{q}}_1 = 0$
- ▶ For white noise ($a = 0$): $\mathbf{\Lambda} = \mathbf{I}$
- ▶ For MA(1) with $N > 2$: \mathbf{R} is tri-diagonal

Example: Eigenvalues and psd for white noise, MA(1) and AR(1)

Example: Eigenvalues and psd for white noise, MA(1) and AR(1)



Example: Eigenvalues and psd for white noise, MA(1) and AR(1)

