

Machine Learning Algorithms

LMS Analysis

Markus Rupp

14.8.2020



LMS Algorithm

- Let's consider the steepest descent iteration:

$$\underline{\hat{w}}_k = \underline{\hat{w}}_{k-1} + \mu \left(r_{\underline{\mathbf{x}}\underline{\mathbf{d}}}^* - R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^* \underline{\hat{w}}_{k-1} \right); \quad k = 1, 2, \dots$$

- Obviously, the knowledge of the cross-correlation vector and the ACF matrix is required.
- If this knowledge is not present, we could operate with estimates. The most simple ones are:

$$\hat{R}_{\underline{\mathbf{x}}\underline{\mathbf{x}}} = \underline{\mathbf{x}}_k \underline{\mathbf{x}}_k^H; \quad \hat{r}_{\underline{\mathbf{x}}\underline{\mathbf{d}}} = \underline{\mathbf{x}}_k \underline{\mathbf{d}}_k^*$$

$$\underline{\mathbf{x}}_k^T = [\mathbf{x}(k), \mathbf{x}(k-1), \dots, \mathbf{x}(k-M+1)]$$



LMS Algorithm

- They are called instantaneous estimates.
- Assuming stationary processes, also other estimates are possible, for example

$$\hat{R}_{\underline{\mathbf{x}}\underline{\mathbf{x}}} = \frac{1}{N} \sum_{l=0}^{N-1} \underline{\mathbf{x}}_{k-l} \underline{\mathbf{x}}_{k-l}^H; \quad \hat{r}_{\underline{\mathbf{x}}\underline{\mathbf{d}}} = \frac{1}{N} \sum_{l=0}^{N-1} \underline{\mathbf{x}}_{k-l} \mathbf{d}_{k-l}^*$$

- We will see in the next chapter that this choice leads to a special form, the so-called RLS algorithm.
- Plugging in the steepest descent iterations the instantaneous estimates, we obtain:

$$\underline{\hat{\mathbf{w}}}_k = \underline{\hat{\mathbf{w}}}_{k-1} + \mu \underline{\mathbf{x}}_k^* \left(\mathbf{d}_k - \underline{\mathbf{x}}_k^T \underline{\hat{\mathbf{w}}}_{k-1} \right); \quad k = 1, 2, \dots$$



LMS Algorithm

- Utilizing such instantaneous values,

$$\underline{\hat{\mathbf{w}}}_k = \underline{\hat{\mathbf{w}}}_{k-1} + \mu \underline{\mathbf{x}}_k^* \underbrace{(\underline{\mathbf{d}}_k - \underline{\mathbf{x}}_k^T \underline{\hat{\mathbf{w}}}_{k-1})}_{\tilde{\mathbf{e}}_{a,k}}, k = 1, 2, \dots$$

- The update equation of the LMS algorithm is obtained.
- The error signal

$$\tilde{\mathbf{e}}_{a,k} = \underline{\mathbf{d}}_k - \underline{\mathbf{x}}_k^T \underline{\hat{\mathbf{w}}}_{k-1}$$

- is called the disturbed/distorted a-priori error signal.
- Due to the reference model, there is also an undisturbed a-priori error signal:

$$\mathbf{e}_{a,k} = \underline{\mathbf{x}}_k^T \underline{\mathbf{w}}_o - \underline{\mathbf{x}}_k^T \underline{\hat{\mathbf{w}}}_{k-1} = \underline{\mathbf{x}}_k^T \underline{\tilde{\mathbf{w}}}_{k-1}$$



LMS Algorithm

- Correspondingly, would we use the estimates at time instant k rather than $k-1$, we call the error a-posteriori error signals:

$$\tilde{\mathbf{e}}_{p,k} = \mathbf{d}_k - \underline{\mathbf{x}}_k^T \hat{\underline{\mathbf{w}}}_k$$

$$\mathbf{e}_{p,k} = \underline{\mathbf{x}}_k^T \underline{\mathbf{w}}_o - \underline{\mathbf{x}}_k^T \hat{\underline{\mathbf{w}}}_k = \underline{\mathbf{x}}_k^T \tilde{\underline{\mathbf{w}}}_k$$

- Note that due to the nature of the recursive algorithm, the estimates are random processes:

$$\hat{\underline{\mathbf{w}}}_k = \hat{\underline{\mathbf{w}}}_{k-1} + \mu \underline{\mathbf{x}}_k^* \left(\mathbf{d}_k - \underline{\mathbf{x}}_k^T \hat{\underline{\mathbf{w}}}_{k-1} \right), \quad k = 1, 2, \dots$$



LMS Algorithm

- Many different choices for the step-size are possible:

arbitrary μ_k : stochastic gradient algorithm

$\mu_k = \frac{\alpha}{\|\underline{\mathbf{x}}_k\|_2^2}$ **with $\alpha > 0$, (data-) Normalized LMS algorithm = NLMS**

$\mu_k = \frac{\alpha}{1 + \alpha \|\underline{\mathbf{x}}_k\|_2^2}$ **with $\alpha > 0$, a - posteriori form of the LMS algorithm**

$\mu_k = \frac{\alpha}{\varepsilon + \|\underline{\mathbf{x}}_k\|_2^2}$ **with $\varepsilon > 0$, the so - called ε – NLMS algorithm**



LMS Algorithm

- Its popularity is also deserved by its mathematical simplicity. Not only

$$E\left[\left|\tilde{\mathbf{e}}_{a,k}\right|^2\right]$$

- Can be minimized but also

$$E\left[\left|f\left[\tilde{\mathbf{e}}_{a,k}\right]\right|^2\right]$$

- With arbitrary nonlinear functions $f[\cdot]$. For example $f[x]=x^{K/2}$ leads to the well known Least-Mean-K algorithm (LMK):

$$\underline{\hat{\mathbf{w}}}_k = \underline{\hat{\mathbf{w}}}_{k-1} + \mu \underline{\mathbf{x}}_k^* \left|\tilde{\mathbf{e}}_{a,k}\right|^{K-2} \tilde{\mathbf{e}}_{a,k}$$



LMS Algorithm

- Other variants of the LMS algorithm are:

Least – Mean Fourth : $\underline{\hat{\mathbf{w}}}_k = \underline{\hat{\mathbf{w}}}_{k-1} + \mu \underline{\mathbf{x}}_k^* |\tilde{\mathbf{e}}_{a,k}|^2 \tilde{\mathbf{e}}_{a,k}$

Least – Mean Mixed Norm : $\underline{\hat{\mathbf{w}}}_k = \underline{\hat{\mathbf{w}}}_{k-1} + \mu \underline{\mathbf{x}}_k^* \left(\beta + (1 - \beta) |\tilde{\mathbf{e}}_{a,k}|^2 \right) \tilde{\mathbf{e}}_{a,k}$

Leaky - LMS : $\underline{\hat{\mathbf{w}}}_k = \beta \underline{\hat{\mathbf{w}}}_{k-1} + \mu \underline{\mathbf{x}}_k^* \tilde{\mathbf{e}}_{a,k}$

Sign - error LMS : $\underline{\hat{\mathbf{w}}}_k = \underline{\hat{\mathbf{w}}}_{k-1} + \mu \underline{\mathbf{x}}_k^* \text{sgn}[\tilde{\mathbf{e}}_{a,k}]$



LMS Algorithm

- Since the LMS algorithm uses instantaneous estimates of cross correlation vector and ACF matrix, we cannot expect that the algorithm behaves identical to the steepest descent algorithm.
- We therefore have to analyze the algorithm by computing its behavior in mean and mean square.



LMS Algorithm Analysis

- Independence assumptions
 - The observation of d_k originates from a reference model: $\mathbf{d}_k = \underline{\mathbf{w}}_o^T \underline{\mathbf{x}}_k + \mathbf{v}_k$, \mathbf{v}_k and $\underline{\mathbf{x}}_k$ are of zero mean.
 - The regression vectors $\underline{\mathbf{x}}_k$ are statistically independent for different time instants: $f_{\underline{\mathbf{x}}\underline{\mathbf{x}}}(\underline{\mathbf{x}}_k, \underline{\mathbf{x}}_l) = f_{\underline{\mathbf{x}}\underline{\mathbf{x}}}(\underline{\mathbf{x}}_k) f_{\underline{\mathbf{x}}\underline{\mathbf{x}}}(\underline{\mathbf{x}}_l)$ for k and l different.
 - The input process \mathbf{x}_k is spherically invariant and Gaussian distributed with zero mean.
 - The additive noise \mathbf{v}_k is statistically independent of the driving process \mathbf{x}_k .
- Note that through these conditions, $\underline{\mathbf{w}}_k$ is statistically independent of $\underline{\mathbf{x}}_l$; $l > k$.



Spherically Invariant Processes

- A complex-valued variable $\mathbf{z}=\mathbf{x}+j\mathbf{y}$ is called Gaussian when \mathbf{x} and \mathbf{y} are joint Gaussian.

- The second moment of \mathbf{z} is given by:

$$R_{\mathbf{zz}} = E[\mathbf{zz}^*] = R_{\mathbf{xx}} + R_{\mathbf{yy}} + j(R_{\mathbf{yx}} - R_{\mathbf{xy}})$$

- Note that the information of $R_{\mathbf{zz}}$ is not sufficient to conclude back to $R_{\mathbf{xx}}, R_{\mathbf{yy}}$ and $R_{\mathbf{xy}}$. For this more information is required, for example:

$$R_{\mathbf{zz}^*} = E[\mathbf{zz}] = R_{\mathbf{xx}} - R_{\mathbf{yy}} + j(R_{\mathbf{yx}} + R_{\mathbf{xy}})$$

- For spherically invariant processes, we have $R_{\mathbf{zz}^*}=0$, or equivalently $R_{\mathbf{xx}}=R_{\mathbf{yy}}$ and $R_{\mathbf{xy}}=-R_{\mathbf{yx}}$.

- The joint density function of a vector $\underline{\mathbf{z}}$ is thus given by:

$$f_{\underline{\mathbf{x}}, \underline{\mathbf{y}}}(x, y) = f_{\underline{\mathbf{z}}}(\underline{\mathbf{z}}) = \frac{1}{\pi^p} \frac{1}{\det(R_{\mathbf{zz}})} \exp\left(-\underline{\mathbf{z}}^H R_{\mathbf{zz}}^{-1} \underline{\mathbf{z}}\right)$$



LMS Algorithm Analysis

- The parameter error vector in the mean:

$$\underline{\tilde{\mathbf{w}}}_k = \underline{w}_o - \underline{\hat{\mathbf{w}}}_k$$

$$\underline{\tilde{\mathbf{w}}}_k = \left(I - \mu \underline{\mathbf{x}}_k^* \underline{\mathbf{x}}_k^T \right) \underline{\tilde{\mathbf{w}}}_{k-1} - \mu \underline{\mathbf{x}}_k^* \mathbf{v}_k$$

$$E[\underline{\tilde{\mathbf{w}}}_k] = E\left[\left(I - \mu \underline{\mathbf{x}}_k^* \underline{\mathbf{x}}_k^T \right) \underline{\tilde{\mathbf{w}}}_{k-1}\right] - \mu E[\underline{\mathbf{x}}_k^* \mathbf{v}_k]$$

- With the independence assumptions:

$$E[\underline{\tilde{\mathbf{w}}}_k] = \left(I - \mu R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^* \right) E[\underline{\tilde{\mathbf{w}}}_{k-1}]$$

- Behavior identical to steepest descent!



LMS Algorithm Analysis

- Behavior identical to steepest descent!

$$E[\tilde{\mathbf{w}}_k] = (I - \mu R_{\mathbf{xx}}^*) E[\tilde{\mathbf{w}}_{k-1}]$$

- Thus convergence guaranteed for

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



LMS Algorithm Analysis

- The parameter error vector in the mean square sense:

$$\begin{aligned} P_k &= E\left[(\underline{w}_o - \hat{\underline{w}}_k)(\underline{w}_o - \hat{\underline{w}}_k)^H\right] = E\left[\tilde{\underline{w}}_k \tilde{\underline{w}}_k^H\right] \\ &= E\left[\left(I - \mu \underline{\mathbf{x}}_k^* \underline{\mathbf{x}}_k^T\right) P_{k-1} \left(I - \mu \underline{\mathbf{x}}_k^* \underline{\mathbf{x}}_k^T\right)^H\right] + \mu^2 E\left[\underline{\mathbf{x}}_k^* \underline{\mathbf{x}}_k^T |\mathbf{v}_k|^2\right] \\ &= P_{k-1} - \mu E\left[\underline{\mathbf{x}}_k^* \underline{\mathbf{x}}_k^T P_{k-1}\right] - \mu E\left[P_{k-1} \underline{\mathbf{x}}_k^* \underline{\mathbf{x}}_k^T\right] + \mu^2 E\left[\underline{\mathbf{x}}_k^* \underline{\mathbf{x}}_k^T P_{k-1} \underline{\mathbf{x}}_k^* \underline{\mathbf{x}}_k^T\right] \\ &\quad + \mu^2 E\left[\underline{\mathbf{x}}_k^* \underline{\mathbf{x}}_k^T |\mathbf{v}_k|^2\right] \end{aligned}$$



LMS Algorithm Analysis

- Most problematic is the term

$$\mathbb{E} \left[\underline{\mathbf{x}}_k^* \underline{\mathbf{x}}_k^T P_{k-1} \underline{\mathbf{x}}_k^* \underline{\mathbf{x}}_k^T \right]$$

- It can be solved explicitly for (spherically) Gaussian processes. We distinguish real-valued processes:

$$\mathbb{E} \left[\underline{\mathbf{x}}_k \underline{\mathbf{x}}_k^T P_{k-1} \underline{\mathbf{x}}_k \underline{\mathbf{x}}_k^T \right] = 2 R_{\underline{\mathbf{x}}\underline{\mathbf{x}}} P_{k-1} R_{\underline{\mathbf{x}}\underline{\mathbf{x}}} + R_{\underline{\mathbf{x}}\underline{\mathbf{x}}} \text{trace}(P_{k-1} R_{\underline{\mathbf{x}}\underline{\mathbf{x}}})$$

- Complex-valued processes:

$$\mathbb{E} \left[\underline{\mathbf{x}}_k^* \underline{\mathbf{x}}_k^T P_{k-1} \underline{\mathbf{x}}_k^* \underline{\mathbf{x}}_k^T \right] = R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^* P_{k-1} R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^* + R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^* \text{trace}(P_{k-1} R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^*)$$



LMS Algorithm Analysis

- The procedure can be extended to general spherically invariant processes. In the general case we obtain

$$\mathbb{E} \left[\underline{\mathbf{x}}_k^* \underline{\mathbf{x}}_k^T P_{k-1} \underline{\mathbf{x}}_k^* \underline{\mathbf{x}}_k^T \right] = \gamma R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^* P_{k-1} R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^* + R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^* \text{trace} \left(P_{k-1} R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^* \right)$$
- with some positive constant $\gamma=1$ for complex and $\gamma=2$ for real valued processes.
- The update recursion for the parameter error covariance matrix reads now:

$$\begin{aligned} P_k = & P_{k-1} - \mu R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^* P_{k-1} - \mu P_{k-1} R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^* + \mu^2 \gamma R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^* P_{k-1} R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^* \\ & + \mu^2 R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^* \text{trace} \left(P_{k-1} R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^* \right) + \mu^2 R_{\underline{\mathbf{x}}\underline{\mathbf{x}}}^* \sigma_v^2 \end{aligned}$$



LMS Algorithm Analysis

- Orthogonalization for the ACF matrix leads to:

$$Q^H R_{\underline{\mathbf{x}}}^* Q = \Lambda; \quad Q^H P_k Q = C_k$$

- We thus obtain:

$$C_k = C_{k-1} - \mu \Lambda C_{k-1} - \mu C_{k-1} \Lambda + \mu^2 \gamma \Lambda C_{k-1} \Lambda \\ + \mu^2 \Lambda \text{trace}[\Lambda C_{k-1}] + \mu^2 \Lambda \sigma_v^2$$

- But note that:

$$E[\|\tilde{\mathbf{w}}_k\|_2^2] = \text{trace}[P_k] = \text{trace}[C_k]$$



LMS Algorithm Analysis

- Thus, only the diagonal terms of C_k are of importance. Order diagonal terms in vector \underline{c}_k :

$$\begin{aligned}\underline{c}_k &= B\underline{c}_{k-1} + \mu^2 \sigma_v^2 \underline{\lambda} \\ B &= I - 2\mu\Lambda + \mu^2\gamma\Lambda^2 + \mu^2 \underline{\lambda}\underline{\lambda}^T \\ &= \begin{cases} 1 - 2\mu\lambda_i + \mu^2(1+\gamma)\lambda_i^2 & \text{main diagonal} \\ \mu^2\lambda_i\lambda_j & \text{else} \end{cases}\end{aligned}$$



LMS Algorithm Analysis

- **Theorem 3.1:** Under the given conditions, the LMS algorithm is convergent in the mean-square sense, when:

$$0 < \mu < \frac{2}{\gamma \lambda_{\max} + \text{trace}(\Lambda)}$$

- **Proof:** Convergence is guaranteed if the eigenvalues of matrix B are upper bounded by one. Since B is positive definite and symmetrical, all its eigenvalues are real-valued and positive.



LMS Algorithm Analysis

- **Proof:** A sufficient condition is that the largest eigenvalue is smaller than one. This is the l_2 norm of the matrix B . An even looser condition is that the l_1 norm of B is smaller than one:

$$\lambda_{\max}^{(B)} = \|B\|_2 \leq \|B\|_1$$

- Thus, for an arbitrary row of B , we have:

$$(1 - \mu\lambda_i)^2 + \mu^2\lambda_i((\gamma - 1)\lambda_i + \text{trace}(\Lambda)) < 1$$

- From which we can follow for μ :

$$0 < \mu \leq \frac{2}{\gamma\lambda_{\max} + \text{trace}(\Lambda)} \leq \frac{2}{\gamma\lambda_i + \text{trace}(\Lambda)}$$



LMS Algorithm: Performance Measures

(parameter vector) error

system mismatch:

$$\begin{aligned} \mathbb{E}[\tilde{\underline{\mathbf{w}}}_k^H \tilde{\underline{\mathbf{w}}}_k] &= \text{trace}(\mathbb{E}[\tilde{\underline{\mathbf{w}}}_k \tilde{\underline{\mathbf{w}}}_k^H]) \\ &= \text{trace}(P_k) \\ &= \underline{\mathbf{1}}^T \underline{\mathbf{c}}_k = \sum_{l=1}^M \gamma_l \lambda_{B,l}^k \end{aligned}$$



LMS Algorithm: Performance Measures

Steady-state system mismatch:

$$\text{trace}(P_k) = \sum_{l=1}^M \gamma_l \lambda_{B,l}^k$$

$$\begin{aligned} \lim_{k \rightarrow \infty} \underline{c}_k &= \underline{c}_\infty = B \underline{c}_\infty + \mu^2 \underline{\lambda} \sigma_v^2 \\ &= [I - B]^{-1} \mu^2 \underline{\lambda} \sigma_v^2 \\ &= \left[2\Lambda - \mu \Lambda^2 \gamma - \mu \underline{\lambda} \underline{\lambda}^T \right]^{-1} \mu \underline{\lambda} \sigma_v^2 \end{aligned}$$



LMS Algorithm: Performance Measures

Matrix-Inversion-Lemma

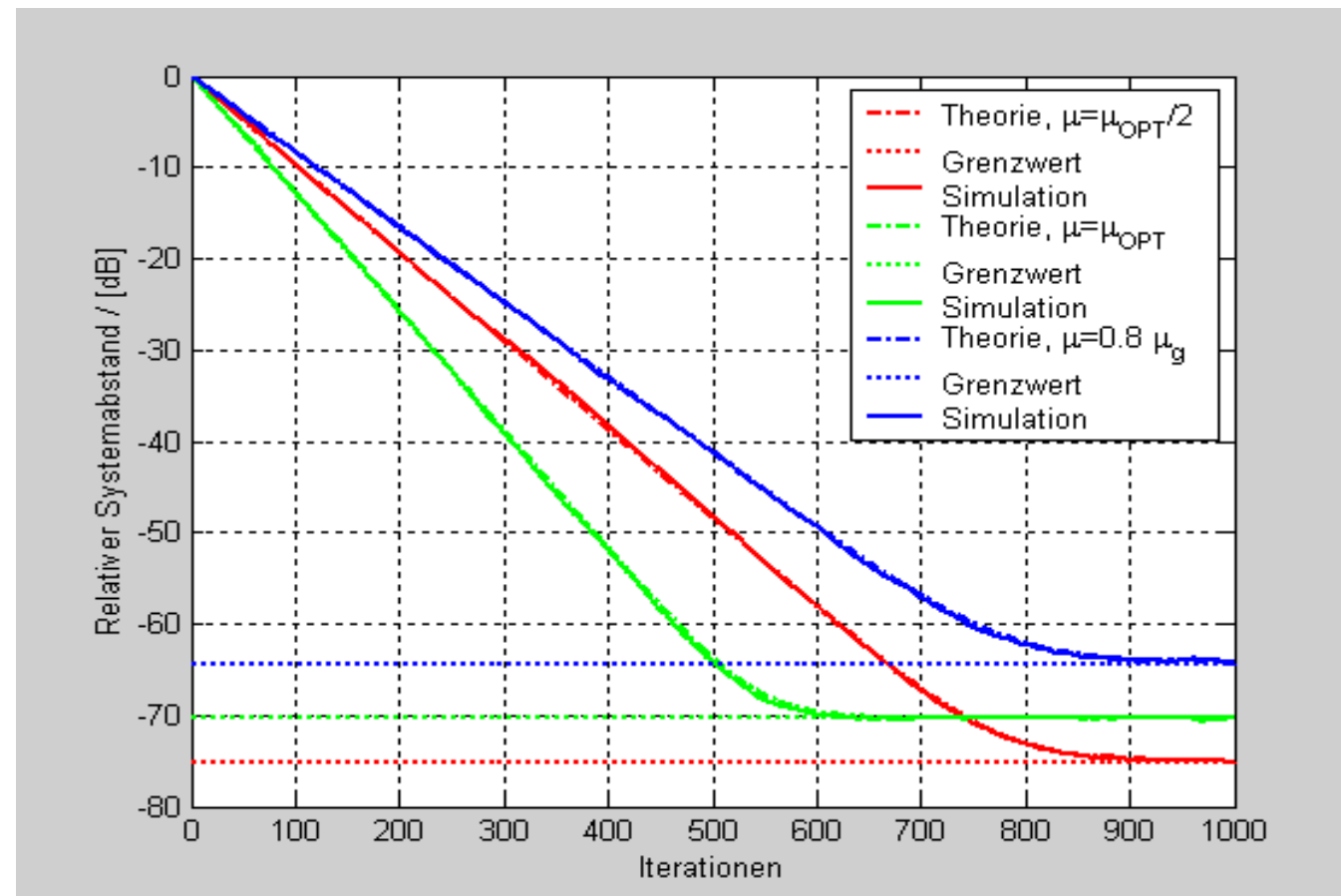
$$\underline{1}^T \underline{c}_\infty = \mu \sigma_v^2 \frac{\sum_{l=1}^M \frac{1}{2 - \mu \gamma \lambda_l}}{1 - \mu \sum_{l=1}^M \frac{\lambda_l}{2 - \mu \gamma \lambda_l}}$$

Small step-sizes μ

$$\underline{1}^T \underline{c}_\infty = \frac{\mu \sigma_v^2 M}{2}$$



LMS Algorithm: Performance Measures



LMS Algorithm: Performance Measures

Relative system mismatch

$$\frac{\underline{1}^T \underline{c}_\infty}{\|\underline{w}_o\|}$$

Distorted a-priori error:

$$\tilde{\mathbf{e}}_{a,k} = \mathbf{d}_k - \underline{\mathbf{x}}_k^T \hat{\mathbf{w}}_{k-1}$$

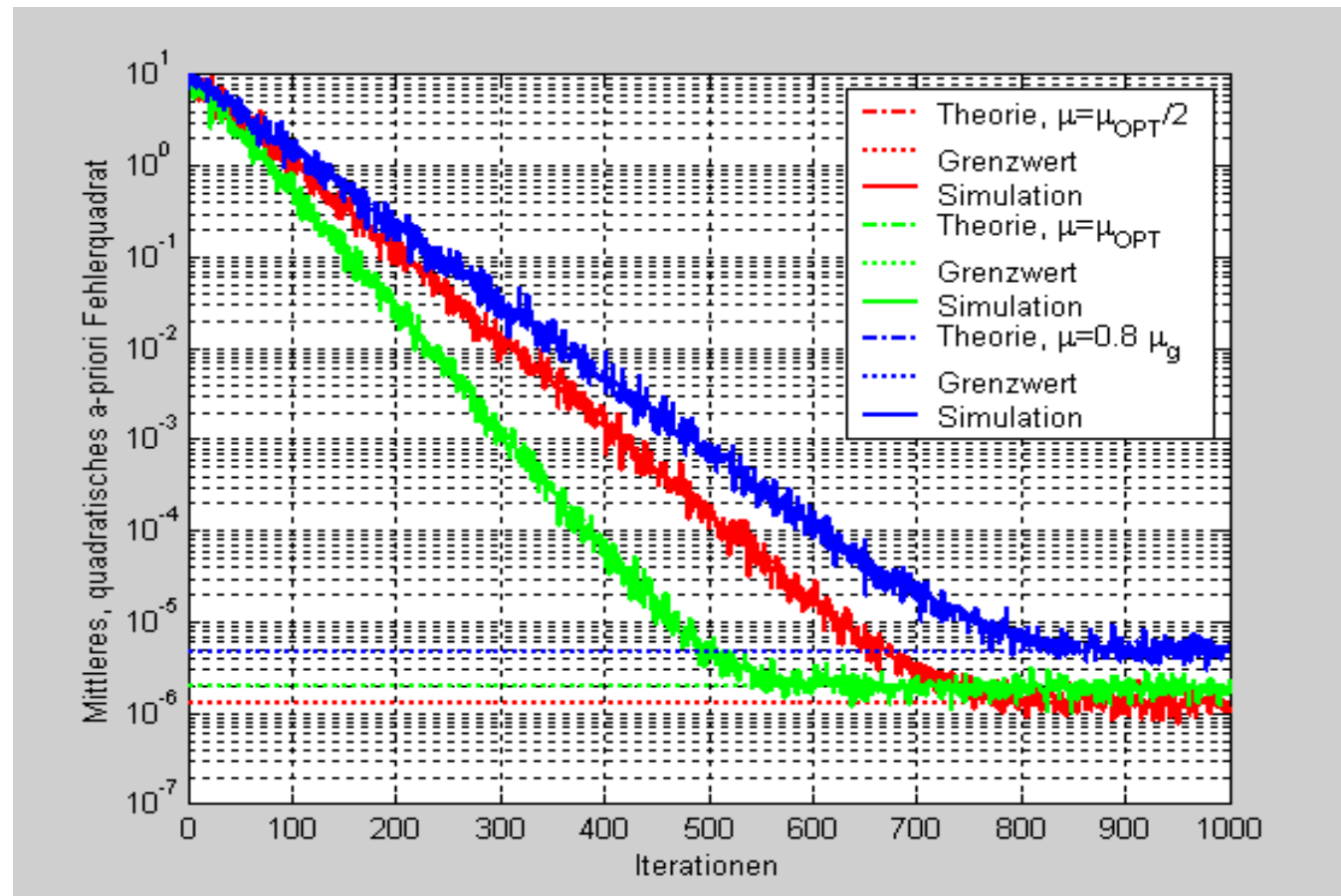


LMS Algorithm: Performance Measures

$$\begin{aligned} \mathbb{E}\left[\left|\tilde{\mathbf{e}}_{a,k}\right|^2\right] &= \mathbb{E}\left[\left|\mathbf{d}_k - \hat{\mathbf{w}}_{k-1}^T \mathbf{x}_k\right|^2\right] \\ &= \mathbb{E}\left[\left|\mathbf{v}_k - \tilde{\mathbf{w}}_{k-1}^T \mathbf{x}_k\right|^2\right] \\ &= \sigma_v^2 + \mathbb{E}\left[\left|\tilde{\mathbf{w}}_{k-1}^T \mathbf{x}_k\right|^2\right] \\ &= \sigma_v^2 + \mathbb{E}\left[\tilde{\mathbf{w}}_{k-1}^T \mathbf{x}_k \mathbf{x}_k^H \tilde{\mathbf{w}}_{k-1}^*\right] \\ &= \sigma_v^2 + \mathbb{E}\left[\tilde{\mathbf{w}}_{k-1}^T R_{\mathbf{xx}}^* \tilde{\mathbf{w}}_{k-1}^*\right] \\ &= \sigma_v^2 + \text{trace}\left(\mathbb{E}\left[R_{\mathbf{xx}}^* \tilde{\mathbf{w}}_{k-1}^* \tilde{\mathbf{w}}_{k-1}^T\right]\right) \\ &= \sigma_v^2 + \underline{\lambda}^T \underline{c}_{k-1} \end{aligned}$$



LMS Algorithm: Performance Measures



LMS Algorithm: Performance Measures

Excess Mean Square Error

$$g_{\text{ex}} = \underline{\lambda}^T \underline{c}_{\infty} = \mu \sigma_v^2 \frac{\sum_{l=1}^M \frac{\lambda_l}{2 - \mu \gamma \lambda_l}}{1 - \mu \sum_{l=1}^M \frac{\lambda_l}{2 - \mu \gamma \lambda_l}}$$

Misadjustment

$$m_{\text{LMS}} = \frac{g_{\text{ex}}}{g_o} = \mu \frac{\sum_{l=1}^M \frac{\lambda_l}{2 - \mu \gamma \lambda_l}}{1 - \mu \sum_{l=1}^M \frac{\lambda_l}{2 - \mu \gamma \lambda_l}}$$



Learning curves

- Note that there is some subtleties in averaging learning curves.
- If the outcome of an experiment is $||\underline{\mathbf{w}} - \underline{\mathbf{w}}_k||_2^2$ then averaging over N experiments is straightforward:
$$1/N \sum ||\underline{\mathbf{w}} - \underline{\mathbf{w}}_k||_2^2$$



Learning curves

- However, how do we average $||\underline{\mathbf{w}} - \underline{\mathbf{w}}_k||_2^2 / ||\underline{\mathbf{w}}||_2^2$
- The random variable $\underline{\mathbf{w}}$ now appears in numerator and denominator!
- Jensens inequality: for convex functions:
- (for concave the other way around)

$$f(E[x]) \leq E[f(x)]$$



Learning curves

- Therefore, we have to do the following:
- $Nu = 1/N \sum ||\underline{\mathbf{w}} - \underline{\mathbf{w}}_k||_2^2$
- $De = 1/N \sum ||\underline{\mathbf{w}}||_2^2$
- $Sys_{rel}(k) = Nu/De$

