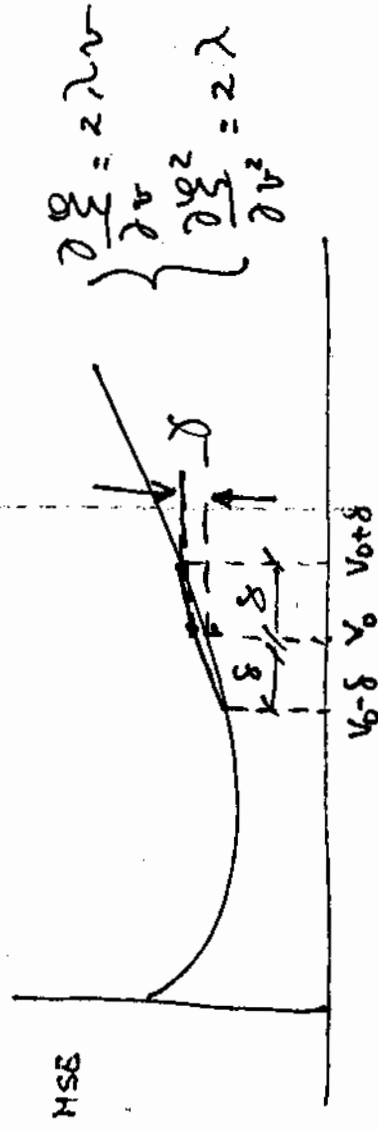


GRADIENT ESTIMATION

ESTIMATION BY DERIVATIVE

Assume we know zeta $\xi = \xi_{\min} + \lambda v^2$



If we estimate numerically by central differences

$$\frac{\partial \xi}{\partial v} \approx \frac{\xi(v_0 + \delta) - \xi(v_0 - \delta)}{2\delta}$$

When $\delta \rightarrow 0$ approximation becomes exact. (if surface is quadratic even for finite delta this is exact)

To implement the estimation we need to compute η @ $V_0 + \delta$ and $V_0 - \delta$. This quantity will be in error by γ , which is,

$$\gamma = \frac{1}{2} \left[\sum (v_0 - \delta) + \sum (v_0 + \delta) \right] - \sum (v_0)$$

or for quadratic surfaces,

$$\gamma = \frac{1}{2} \left[2 \sum_{\text{MIN}} (v_0 + \delta)^2 + \lambda (v_0 - \delta)^2 \right] - \left(\sum_{\text{MIN}} + \lambda v_0^2 \right)$$

PERFORMANCE PENALTY is defined as

$$P = \frac{\gamma}{\sum_{\text{MIN}}} = \frac{\lambda \delta^2}{\sum_{\text{MIN}}}$$

and it translates the average increase in the error normalized to η min.

If we have multiple weights

$$\begin{aligned} \mathcal{E} &= \sum_{\text{MIN}} + \mathbf{V}^T \mathbf{R} \mathbf{V} = \sum_{\text{MIN}} + [\mathbf{v}_0 \ \mathbf{v}_1] \begin{bmatrix} r_{00} & r_{01} \\ r_{10} & r_{11} \end{bmatrix} \begin{bmatrix} v_0 \\ v_1 \end{bmatrix} \\ &= r_{00} v_0^2 + r_{11} v_1^2 + 2 r_{01} v_0 v_1 \end{aligned}$$

along each axis we have

$$P_0 = \frac{r_{00} \delta^2}{\sum_{\text{MIN}}} ; P_1 = \frac{r_{11} \delta^2}{\sum_{\text{MIN}}}$$

and the average perturbation is $P = \frac{1}{2} (P_0 + P_1)$

$$P = \frac{\delta^2}{\sum_{\text{MIN}}} \left(\frac{r_{00} + r_{11}}{2} \right)$$

For L+1 weights

$$P = \frac{\delta^2}{\sum_{\text{MIN}}} \frac{\text{TRACE}[\mathbf{R}]}{L+1}$$

or (since the trace is the sum of eigenvalues)

$$P = \frac{\delta^2}{\sum_{\text{MIN}}} \sum_{n=0}^L \frac{\lambda_n}{L+1} = \frac{\delta^2}{\sum_{\text{MIN}}} \lambda_{\text{AVE.}}$$

ESTIMATION THEORY

We have assumed that we knew η . But

$$\sum_s = E\{e_k^2\}$$

$$\sum_t = \sum_{k=-\infty}^{\infty} e_k^2$$

In each we have finite data sets, i.e. we have to estimate η , the ^{second} moment of the random process. FOR THE FIRST MOMENT,

$$\mu_x = E\{x_n\} = \int_{-\infty}^{\infty} x p_x(x) dx$$

$$\mu_x = A\{x_n\} = \lim_{N \rightarrow \infty} \frac{1}{2N+1} \sum_{-N}^N x_n$$

We assume that the r.p. is ergodic, i.e. $\sum_s = \sum_t$.

It is plausible that if N is large the estimate given by the sampled mean will be good. When we use this estimate we are in fact using MAXIMUM LIKELIHOOD estimates. The r th moment is estimated by

$$\hat{\alpha}_r = \frac{1}{N} \sum_{k=1}^N (e_k)^r$$

The sampled mean for Gaussian r.v. is consistent, i.e. the bias is zero

$$\alpha_n = E[\hat{\alpha}_n]$$

$$\text{Var}[\hat{\alpha}_n] = E[(\hat{\alpha}_n - \alpha_n)^2] = \frac{\alpha_n^2 - \hat{\alpha}_n^2}{N}$$

and the variance decreases with N.

We want to estimate the gradient of the performance surface, so let us see how this theory applies.

Notice that η is defined as the second moment of $\hat{\alpha}_n$. Since the estimator is unbiased (estimated mean equals the true mean), the variance of the estimate is

$$\text{Var}[\hat{\eta}] = \frac{\alpha_1 - \alpha_2}{N}$$

To compute the power one must know the p.d.f. of $\hat{\eta}$. Let us assume that it is zero mean Gaussian, with variance σ^2

Then,

$$\text{Var}[\hat{\eta}] = \frac{3\sigma^4 - \sigma^4}{N} = 2 \frac{\sigma^2}{N}$$

In general 2 is exact or conservatively high.

Knowing the variance of $\hat{\eta}$, we are ready to compute the variance of the gradient estimate. Since

$$\frac{\partial \hat{\xi}}{\partial v} = \frac{1}{2\delta} \hat{\xi}(v+\delta) - \frac{1}{2\delta} \hat{\xi}(v-\delta)$$

and assuming independent estimates

$$\begin{aligned} \text{Var} \left[\frac{\partial \hat{\xi}}{\partial v} \right] &= \frac{1}{4\delta^2} \text{Var} \left[\hat{\xi}(v+\delta) \right] + \frac{1}{4\delta^2} \text{Var} \left[\hat{\xi}(v-\delta) \right] \\ &= \frac{1}{2N\delta^2} \left[\sum_{k=1}^N \xi^2(v+\delta) + \sum_{k=1}^N \xi^2(v-\delta) \right] \end{aligned}$$

After convergence, and for small perturbation delta,

$$\text{Var} \left[\frac{\partial \hat{\xi}}{\partial v} \right] = \frac{\sum_{k=1}^N \xi^2}{N\delta^2}$$

SO

Assuming independence of ϵ_k (therefore of the errors), and that N and δ are the same for all components of the gradient, the covariance of the estimated gradient is

$$\text{cov} \left[\hat{\nabla}_k \right] = E \left[\left(\hat{\nabla}_k - \nabla_k \right) \left(\hat{\nabla}_k - \nabla_k \right)^T \right] = \frac{\sum_{k=1}^N \text{Min}_k}{N\delta^2} \mathbf{I}$$

Notice that the covariance is inversely proportional to δ , i.e. the larger the δ the better the estimate (why?).

With this result we can now calculate the effect of the imprecision (noise) in the gradient into the weight solution. Remember that we are calculating the weights iteratively using the gradient information.

WEIGHT VECTOR COVARIANCE

NEWTON'S METHOD

Let us write the true gradient as

$$\hat{\nabla}_k = \nabla_k + N_k$$

The recursion to calculate the new weight vector is

$$W_{k+1} = W_k + \mu R^{-1} \hat{\nabla}_k$$

Substituting the noisy gradient

$$W_{k+1} = W_k + \mu R^{-1} \nabla_k - \mu R^{-1} N_k$$

System is cross coupled ($R^{-1} N$). However in the principal coordinate system

$$V_{k+1} = (1-2\mu)V_k - \mu R^{-1} N_k$$

\Rightarrow

$$V'_{k+1} = (1-2\mu)V'_k - \mu \Lambda^{-1} N'_k \quad (\theta^{-1} N = N')$$

and the equations become decoupled. By induction the solution is

$$V'_{k+1} = (1-2\mu)^k V'_0 - \mu \Lambda^{-1} \sum_{n=0}^{k-1} (1-2\mu)^n N'_{k-n-1}$$

Note that we have the previous term plus a noise component. For very large k when $\mu \approx 1/2$, the first term vanishes and

$$V'_{k+1} = -\mu \Lambda^{-1} \sum_{n=0}^{\infty} (1-2\mu)^n N'_{k-n-1}$$

This is the steady state error after adaptation.

STEEPEST DESCENT

$$W_{k+1} = W_k - \mu \nabla_k$$

Once again the recursion can be written as a function of the true gradient plus the noise component.

$$V_{k+1} = (I - z\mu R) V_k - \mu N_k$$

In the principal coordinate system,

$$V'_{k+1} = (I - z\mu \Lambda) V'_k - \mu N'_k$$

and by induction the solution can be written,

$$V'_k = (I - z\mu \Lambda)^k V'_0 - \mu \sum_{n=0}^{k-1} (I - z\mu \Lambda)^n N'_{k-n-1}$$

If μ is chosen for convergence, for large k the first term vanishes, and the steady state solution becomes

$$V'_k = -\mu \sum_{n=0}^{\infty} (I - z\mu \Lambda)^n N'_k$$

So the net effect of the noisy gradient estimates is to produce even after adaptation a change in the weight vector. The question now is how much noise power? Since we have already the covariance of the estimated gradient vector we can calculate the covariance of the weight vector solution.

NEWTON'S

The covariance is $\text{Cov}[v'_k] = E[v'_k v_k^T]$

and substituting we get

$$\text{Cov}[v'_k] = \frac{\mu(\Lambda^{-1})^2}{4(1-\mu)} \text{Cov}[N'_k]$$

In the STEEPEST DESCENT

$$\text{Cov}[v'_k] = \frac{\mu}{4} (\Lambda - \mu \Lambda^2)^{-1} \text{Cov}[N'_k]$$

Now substituting the calculated value of the gradient cov.

NEWTON'S

$$\text{Cov}[v'_k] = \frac{\mu(\Lambda^{-1})^2 \sum_{\text{MIN}} \sigma^2}{4N\delta^2(1-\mu)}$$

STEEPEST

$$\text{Cov}[v'_k] = \frac{\mu(\Lambda - \mu \Lambda^2)^{-1} \sum_{\text{MIN}} \sigma^2}{4N\delta^2}$$

$$(\text{Cov}[N'_k] = E[N'_k N_k^T]) = \sigma^2 E[N_k N_k^T] \sigma^2 = \frac{\sum_{\text{MIN}} \sigma^2}{N\delta^2} I$$

In the translated coordinate system, noting that

$$\text{cov}[V_k] = E[V_k V_k^T] = Q E[V_k R_k^{-1}]^T Q^{-1} = Q \text{cov}[V_k] Q^{-1}$$

we get

NEWTON'S

$$\text{cov}[V_k] = \frac{\mu \sigma_{\text{MIN}}^2 (R^{-1})^2}{4N \delta^2 (1-\mu)}$$

STEEPEST

$$\text{cov}[V_k] = \frac{\mu \sigma_{\text{MIN}}^2 (R - \mu R^2)^{-1}}{4N \delta^2}$$

Note that: *in the estimation of the gradient*

1. If there is no noise, $\text{cov}[V]$ is zero.
2. The larger the # of measurements, the smaller the error
3. the larger the δ , the smaller the error.
4. The smaller the μ the smaller the error.

EXCESS MEAN SQUARE ERROR

Noise in the gradient estimate implies inaccuracies in the filter coefficients, which mean that even after adaptation, the coefficients will be moving around the bottom of the "bowl".

This means that the error will not be the theoretical minimum, i.e. an excess mean square error will result from the noisy estimation of the gradient.

Notice that this does not include the effect of the perturbation required to estimate the derivative.

Let us characterize the mean of this error. Assume independence in the errors across iterations, stationarity and convergence.

In the principal component system,

$$\text{Excess MSE} = E \left[\sum_k \xi_k - \xi_{\min} \right] = E \left[V_k^T R V_k \right]$$

$$\Rightarrow E \left[V_k^T R V_k \right]$$

NEWTON'S

Define $r = 1 - 2\mu$

$$E_x \text{MSE} = \mu^2 \sum_n \sum_m r^{n+m} E \left[N_{k-n-1}^T \Lambda^{-1} N'_{k-n-1} \right]$$

Since independent,

$$E_x \text{MSE} = \mu^2 \sum_{n=0}^{\infty} r^{2n} E \left[N_{k-n-1}^T \Lambda^{-1} N'_{k-n-1} \right]$$

and stationarity,

$$E_x \text{MSE} = E \left[N_k^T \Lambda^{-1} N'_k \right] \frac{\mu^2}{1-r^2}$$

After evaluation of the matrix

$$E_x \text{MSE} = \frac{\sum_{m=0}^L \min \mu^2}{NS^2(1-r^2)} \sum_{m=0}^L \frac{1}{\lambda_m}$$

If we express the result in terms of τ , the adaptive time constant,

$$r \approx 1 - \frac{1}{2\tau} \Rightarrow \mu \approx \frac{1}{2\tau}$$

iteration

we get

$$E_{\text{excess}} \text{MSE} = \frac{\sum_{m=0}^L \min}{8NS^2\tau} \sum_{m=0}^L \frac{1}{\lambda_m}$$

Defining the average inverse eigenvalue as,

$$\sum_{m=0}^{L-1} \frac{1}{\lambda_m} = (L+1) \left(\frac{1}{\lambda} \right)_{\text{avg}}$$

and substituting

$$\text{excess MSE} \approx \frac{(L+1) \sum_{\text{min}} \lambda_{\text{av}} (1/\lambda_{\text{av}})}{8 N P \sigma^2}$$

Notice that the excess MSE is directly proportional to the number of filter coefficients, inversely proportional to the adaptive time constant and the perturbation.

STEEPEST

The method is the same. Substitute V in the definition of the excess MSE, we get

$$\text{excess MSE} = \frac{\mu \sum_{\text{min}} \lambda_{\text{av}}}{4 N P} \sum_{m=0}^{L-1} \frac{1}{1 - \mu \lambda_m}$$

Now the problem is that we do not have a single adaptive time constant.

TIME CONSTANT OF THE LEARNING CURVE, AND THE TIME CONSTANT OF ADAPTATION.

For the learning curve adaptation, since $r_{mse} = r^2$

$$\exp\left(-\frac{1}{\tau_{mse}}\right) = \exp\left(-\frac{2}{\tau}\right)$$

then

$$\tau_{mse} = \frac{\tau}{2}$$

This constant describes the learning time. The unit is the iteration number.

The units for the time constant of adaptation, T is the data sample (can be related to the sampling rate). Since we need $2N$ samples per gradient component, then $2(L+1)N$ samples per iteration.

$$T_{mse} = 2(L+1)N\tau_{mse} = N(L+1)\tau$$

$$= \frac{N(L+1)}{2\mu}$$

CAN TRADE N FOR $\mu \implies$ WHAT MATTERS IS # OF POINTS PER TIME CONSTANT

Returning to the steepest descent, the n th adaptive time constant is

$$\tau_n = 1 - 2\mu\lambda_n \quad ; \quad \mu\lambda_n = \frac{1}{2\sigma_n}$$

If we proceed to calculate the n th time constant of the learning process

$$(\sigma_{mse})_n = \frac{\sigma_n}{2}$$

and the n th t.c. for the adaptation process

$$(T_{mse})_n = N(L+1)\sigma_n \Rightarrow \frac{N(L+1)}{2\mu\lambda_n}$$

averaging over L we get

$$\mu\lambda_{ave} = \frac{N(L+1)}{2} \left(\frac{1}{T_{mse}} \right)_{ave.}$$

and assuming t_n large, we have

$$\text{excess MSE} \approx \frac{(L+1)^2 \sum_{\text{Min}}}{8P} \left(\frac{1}{T_{mse}} \right)_{ave.}$$

For the Newton's method also

NEWTON

$$\text{excess MSE} = \frac{(L+1)^2 \sum_{m=1}^M \lambda_{\text{ave}} (1/\lambda_{m4})}{8P T_{\text{mse}}}$$

STEEPEST

$$\text{excess MSE} = \frac{(L+1)^2 \sum_{m=1}^M \lambda_{\text{min}}}{8P} \left(\frac{1}{T_{\text{mse}}} \right)_{\text{ave.}}$$

These expressions show that the excess MSE:

1. decreases with larger perturbations (more accurate gradient).
2. decreases with larger time constants (more data is being used).
3. increases with the square of the weights (try always to live with the minimum # of degrees of freedom).
4. eigenvalues only affect the error if there is a large disparity between them.

Can trade time constant and N. So, what is important is the number of points used per time constant of adaptation.

MISADJUSTMENT

The excess mean square error is the average MSE less the minimum MSE.
$$M \hat{=} \left(\frac{\text{excess MSE}}{\sum \text{MIN}} \right)$$

MISADJUSTMENT is the excess MSE normalized by the minimum MSE. It is dimensionless, but now it becomes a good characterization of the cost of adaptability.

Notice that misadjustment of the Steepest and Newton is equal if all the eigenvalues of R are equal.

For the same speed of adaptation, they have different misadjustments.

COMPARISON OF NEWTON'S AND STEEPEST

We have to compare them in terms of misadjustment and speed of adaptation.

$$\text{Writing } (T_{mse})_{MAX} = \frac{(L+1) N}{2\mu \lambda_{MIN}}$$

we can express M for the steepest descent as

$$M \approx \frac{(L+1)^2 \lambda_{AVG}}{8P \lambda_{MIN} (T_{mse})_{MAX}}$$

while for the Newton's is

$$M = \frac{(L+1)^2 \lambda_{AV} (1/\lambda_{AVE})}{8P T_{mse}}$$

We see that if $T_{mse, max} = T_{mse}$ with the other things constant, Newton's method has a smaller misadjustment.

$$\frac{1}{\lambda_{MIN}} \gg \left(\frac{1}{\lambda}\right)_{AVE}$$

TOTAL MISADJUSTMENT

It may seem that we can tradeoff misadjustment by perturbation. If we add the perturbation to the misadjustment, we get the total misadjustment.

$$M_{\text{Tot}} \triangleq M + P$$

This is the measure that we should be interested in. Now the form of the total misadjustment is

$$M_{\text{Tot}} \approx P + \frac{A}{P}$$

So the best perturbation is achieved by setting the derivative = 0, which means

$$P_{\text{opt}} \approx \frac{1}{2} M_{\text{Tot}}$$

Total misadjustment is minimized when the perturbation is half of the total misadjustment.

LMS ALGORITHM

LMS is a steepest descent algorithm. Its basic features are:

1. uses a very simple estimate of the performance surface.

$$\varepsilon_k^2 = (d_k - y_k)^2 \quad \text{ESTIMATE OF } E\{\varepsilon_k^2\}$$

2. Assumes zeta quadratic

DERIVATION

$$\varepsilon_k = d_k - y_k = d_k - \mathbf{x}_{1k}^T \mathbf{w}_k$$

If zeta is quadratic then

$$\hat{\nabla}_k = \begin{bmatrix} \frac{\partial \varepsilon_k^2}{\partial w_0} \\ \vdots \\ \frac{\partial \varepsilon_k^2}{\partial w_L} \end{bmatrix} \approx \begin{bmatrix} \frac{\partial \varepsilon_k^2}{\partial w_0} \\ \vdots \\ \frac{\partial \varepsilon_k^2}{\partial w_L} \end{bmatrix} = 2\varepsilon_k \begin{bmatrix} \frac{\partial \varepsilon_k}{\partial w_0} \\ \vdots \\ \frac{\partial \varepsilon_k}{\partial w_L} \end{bmatrix} = -2\varepsilon_k \bar{\mathbf{x}}_k$$

Therefore,

$$\bar{\mathbf{w}}_{k+1} = \bar{\mathbf{w}}_k + z \underbrace{\varepsilon_k \bar{\mathbf{x}}_k}_{\hat{\nabla}_k} = \bar{\mathbf{w}}_k + \mu (\mathbf{d}_k - y_k) \bar{\mathbf{x}}_k$$

Can see that the update of coefficients involve ONLY a matrix multiplication. So the asymptotic complexity of the algorithm (the big O notation) is

$$O(n) \quad N = L + 1$$

i.e., linear in the number of weights.

This is VERY EFFICIENT!

The algorithm does not use any averaging, squaring, perturbation, etc. Each component of the gradient is obtained from a SINGLE data point. NO assumptions or knowledge about the statistics of the input/ desired signal are made.

As can be expected the estimation of the gradient is VERY NOISY. What is amazing is that the iterative process works as an effective lowpass filter.....

CONVERGENCE OF THE ALGORITHM

$E[\hat{\mathbf{w}}_k]$ is UNBIASED. In fact

$$\begin{aligned} &= -2 E[\epsilon_k \bar{X}_{1c}] = -2 E[d_k \bar{X}_{1c} - X_k X_k^T w_k] \\ &= 2(Rw - P) = \nabla \end{aligned}$$

So the LMS algorithm can be made conceptually identical to the steepest descent. Just compute the estimate, but instead of updating the weights at each step, only change them after a large number of samples. These two procedures are equivalent if the estimator is unbiased. In fact,

$$\begin{aligned} E[w_{k+1}] &= E[w_k] + z\mu E[\epsilon_k X_k] = \\ &= E[w_k] + z\mu \{P - R E[w_k]\} = \\ &= E[w_k] (I - z\mu R) + z\mu R w^* \quad (P = R w^*) \end{aligned}$$

This is equivalent to the expression for the steepest descent.

In the principal coordinate system

$$E[w_k'] = (I - z\mu \Lambda)^k V_0'$$

and we know that this converges to W^* if

$$0 < \mu < \frac{1}{\lambda_{\text{MAX}}}$$

As

$$\lambda_{\text{MAX}} < \sum \lambda_i = t_n[R]$$

and

$$0 < \mu < \frac{1}{t_n[R]}$$

$$t_n[R] = (L+1) \epsilon \{x_k^2\}$$

we can say

$$0 < \mu < \frac{1}{(L+1) (\text{SIGNAL POWER})}$$

Practically,

$$\mu \approx 0.1 \mu_{\text{THEORY}}$$

LEARNING CURVE

In the steepest descent we have defined the time constant of each mode as

$$\tau_n = \frac{1}{2\mu\lambda_n} \quad n = 0, 1, \dots, L$$

and for the time constant of the learning curve

$$\tau_{mse} = \frac{\tau}{2} \Rightarrow (\tau_{mse})_n = \frac{1}{4\mu\lambda_n}$$

Due to the fact that the LMS only uses 1 data sample per iteration,

$$(\tau_{mse})_n = (\tau_{mse})_n$$

This result is an optimistic approximation since the error may be very different from its expected value.

Notice that $\tau_{mse} = \max_i \left\{ \frac{1}{4\mu\lambda_i} \right\} = \frac{1}{4\mu\lambda_{\min}}$

But as $\mu < 1/\lambda_{\max}$, if say $\mu = \frac{\alpha}{\lambda_{\max}} \quad 0 < \alpha < 1$

then $\tau_{mse} = \frac{1}{4\alpha} \frac{\lambda_{\max}}{\lambda_{\min}}$

So, time constant of adaptation is limited by the eigenvalue spread.

NOISE IN WEIGHT SOLUTION

Need another expression because LMS does not use perturbation. Write as before

$$\hat{V}_k = \bar{V}_k + \tilde{N}_k$$

Near adaptation ($\bar{V}_k = 0$) $\hat{V}_k = \tilde{N}_k = -\epsilon \epsilon_k X_k$

$$\text{So Cov}[N_k] = 4 E[\epsilon_k^2 X_k X_k^T]$$

Assuming uncorrelated noise and signal vectors

$$= 4 E[\epsilon_k^2] \cdot E[X_k X_k^T] = 4 \sigma_{\text{MIN}}^2 R$$

In the principal coordinate system

$$\text{Cov}[N'_k] = 4 \sigma_{\text{MIN}}^2 \Lambda$$

So the covariance of the coefficient vector is

$$\text{Cov}[V'_k] = \mu \sigma_{\text{MIN}}^2 (\Lambda - \mu \Lambda^2)^{-1} \Lambda$$

If $\mu \Lambda^2 \ll \Lambda$

$$\text{cov} [V_k'] = \mu \sum_{\text{min}} I$$

The the covariance of V_k is proportional to μ .

MISADJUSTMENT

Now can write

$$\text{excess MSE} = E[V_k'^T \Lambda V_k'] = \sum_{n=0}^L \lambda_n E[V_{nk}'^2]$$

Assuming that $E[V_{nk}']$ is a component of the covariance (the analysis is done after transients died out),

$$\text{excess MSE} = \sum_{n=0}^L \lambda_n \cdot [\mu \sum_{\text{min}} I] = \mu \sum_{\text{min}} \text{tr} [R]$$

So

$$M = \mu \text{tr} [R]$$

i.e. the misadjustment in the LMS is proportional to μ .

$$\left\{ \begin{array}{l} M = \mu \tau_n [R] \\ T_{mse} = \frac{1}{4\mu} \lambda_n \end{array} \right.$$

we see that there is a trade-off between speed of adaptation and misadjustment. As

$$\tau_n [R] = \sum_0^{\infty} \lambda_n = \frac{1}{4\mu} \sum \frac{1}{(\sigma_{mse})_n} = \frac{L+1}{4\mu} \left(\frac{1}{\sigma_{mse}} \right)_{AVE}$$

SO

$$M = \frac{L+1}{4} \cdot \left(\frac{1}{\sigma_{mse}} \right)_{AVE}$$

This is a good approximation. Since after 4 time constants the adaptation transients have died out, we can say that in the LMS MISADJUSTMENT IS EQUAL TO THE NUMBER OF WEIGHTS DIVIDED BY SETTLING TIME (for equal weights).

$$M = \frac{L+1}{4 \sigma_{mse}}$$

PERFORMANCE LMS - STEEPEST DESCENT

LMS is better (less computation, faster to converge):

1. LMS misadjustment increases linearly with weights
2. For same misadjustment LMS is much faster
3. LMS shares with steepest descent the problem of adaptation time dependency on eigenvalue spread.

(i.e. Misadjustment controlled by fast modes

Time constant controlled by slowest modes)

OTHER FORMS OF LMS ALGORITHMS

COMPLEX LMS

If data is complex, then complex LMS

$$w_{k+1} = w_k + 2\mu \varepsilon_k X_k^*$$

NORMALIZED LMS

Speed of converge depends on signal power

$$\lambda_{\max} = \sum_{i=0}^L \lambda_i = E[X_k X_k^T] \cdot (L+1)$$

$$0 < \mu < \frac{1}{\lambda_{\max}}$$

Alternatively, if signal power changes, a fixed μ can lead to divergence.

To cope with this, NORMALIZE stepsize by the input power.

SIGN ERROR LMS

In the other end of things (simplicit) we can try to avoid the multiplications. (NOTE: this is not an issue anymore with DSP chips).

We saw that to update the coefficients we need $L+1$ multiplications, and $L+1$ multiplications to implement the filter. So total $\sim 2L$.

One idea is just use the sign of the error for adaptation.

$$\begin{aligned} \varepsilon_k = d_k - y_k &\implies \varepsilon'_k = \text{sgn} [d_k - y_k] \\ \text{sgn} [d_k - y_k] &= \begin{cases} 1 & \text{if } d_k > y_k \\ 0 & \text{if } d_k = y_k \\ -1 & \text{if } d_k < y_k \end{cases} \end{aligned}$$

There is a penalty paid, because the estimate of the gradient became even noisier. To obtain a similar performance to the straight LMS must use smaller μ s (which imply slower adaptation).

The method changes the scale of the update vector.

$$\bar{w}_{k+1} = \bar{w}_k + \mu \varepsilon'_k \bar{x}_k$$

With this scheme we still need $L+1$ multiplications.

There are 2 ways to cut down on this:

1. First is to approximate μ by negative powers of 2 (left RIGHT shifts). To update the coefficients just left shift the input and add. This may be a significant speed improvement.

2. Take only the sign of the input.

$$\bar{x}_k \rightarrow \text{sgn}[x_k]$$

$$\bar{\omega}_{k+1} = \bar{\omega}_k + \mu \varepsilon_k \bar{x}_k$$

This modification is more drastic because the DIRECTION of the LMS update is changed. It can lead to instability. Combining the two simplifications

$$\omega_{k+1} = \omega_k + \mu \varepsilon_k x_k$$

In the CCITT standard for ADPCM transmission this was the adopted model.

We can define

$$\mu(k) = \frac{\alpha}{X_k^T X_k} \quad 0 < \alpha < 1$$

With this μ_k

$$W_{k+1} = W_k + \mu_k \frac{\alpha}{X_k^T X_k} \epsilon_k X_k$$

Experimentally it was verified that $X_k^T X_k$ can become too small, so modify to (heuristics)

$$\bar{W}_{k+1} = \bar{W}_k + \alpha \epsilon_k \frac{\bar{X}_k}{\gamma + X_k^T X_k}$$

where gamma is a small positive constant (can be 1).

It may seem that we made the algorithm much more complex (a vector multiplication). However, this is not so.

Remember that

$$X_k^T X_k = \sum_{i=0}^{L-1} x_i^2(k-i)$$

So if we use $L+1$ storage locations, can compute X_{k+1}^T by adding $x^2(k+1)$ and subtracting $x^2(k-L)$, so just need squaring and an extra division.

Can implement this operation recursively, through an exponentially weighted time average of the input square. Notice that we can compute recursively an estimate of the power by

$$\sigma_{k+1}^2 = \alpha \sigma_k^2 + (1-\alpha) X_k^2 \quad \left(H(z) = \frac{(1-\alpha)z^{-1}}{z-\alpha} \right)$$

If you think of it, this is the output of a lowpass filter with a pole at alpha, when the input is the signal squared. For alpha close to 1 (.9 to .99) we are smoothing the input square (the envelope). Now the limit of this iteration for large k is

$$\sigma_k^2 = (1-\alpha) \sum_{j=0}^{\infty} \alpha^j x^2(k-j)$$

LMS in non stationary environments

Figures of merit of the LMS have been derived assuming stationarity conditions.

Nonstationarity may arise at least in 2 important cases:

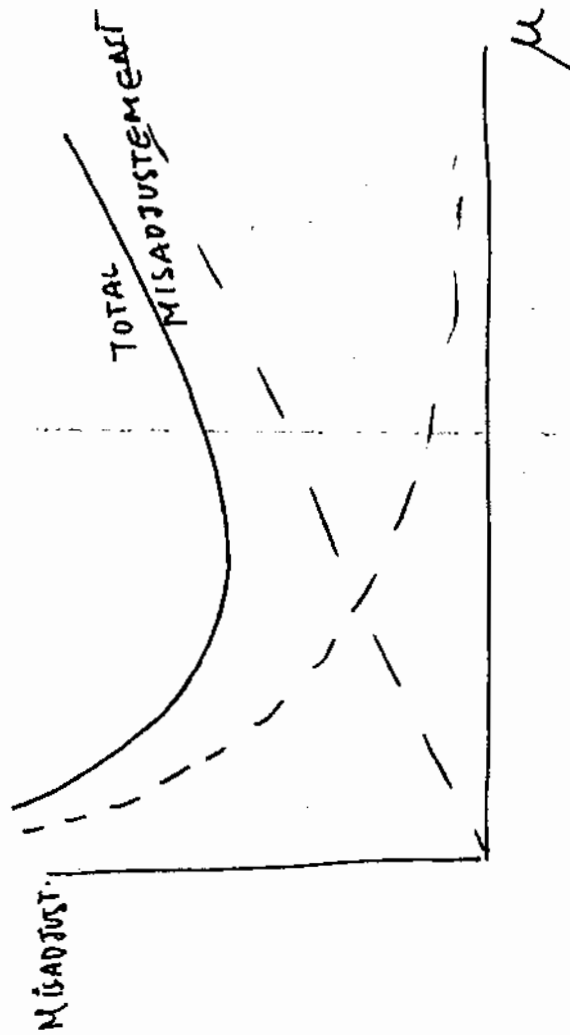
- The desired signal is time varying (system identification). Here the autocorrelation function is time invariant, but the crosscorrelation function is time varying.

- The input is non-stationary (equalization, noise cancellation). Here both the autocorrelation and crosscorrelation are non stationary.

In non-stat environments, the optimum set of weights change with time. Therefore, the LMS must not only find the minimum, but also track the change of the minimum of the performance surface.

In this case we have two error sources for the actual weights. One that is related to the errors of the estimate of the gradient. The other is related to the lag in the adaptive process.

It can be shown (Widrow and al 1976) that the first contribution is directly proportional to μ . And that the second contribution is inversely proportional to μ . We can assume them additive, so the total misadjustment has a single minimum.



IMPLEMENTATION ISSUES- ROUND OFF NOISE.

It is a known fact that the finite precision of the fixed point computer arithmetic causes errors in digital filter implementations. This error can be modelled as additive white noise, with a variance equal to weight of the least significant bit kept in the calculations.

This is difficult to model, but for fixed coefficients digital filters, we know how to estimate the roundoff noise level at the output. This noise will degrade the output dynamic range.

What I would like to address is the effect of fixed point arithmetic in the adaptation process.

In the LMS we have

Therefore in the weight update equation, we have a product of the input by the error, which must be further reduced by μ . μ is generally small, which means that we may get a truncated value of zero to update the weights..., i.e. no update.

We say that the weight stalls. This is responsible for excess mean square error.

So the theoretical result that the smaller the μ , the smaller the misadjustment is no longer valid.

This effect is dependent on the wordlength of the computation of the update equation. The longer the wordlength of the accumulator to keep the results (and the weights) the better we are.

This means practically two things:

- The precision of the coefficients for the update equation should be larger than the one used for the filtration. An example is to use double precision for the weight computations and only use the MSB to implement the filter.
- Block computation minimizes this problem.

PROBLEMS WITH LMS

1. It is slow to converge for large eigenvalue spread.
2. It uses a noisy gradient estimate (stochastic gradient).
3. It "rattles".
4. The dependency on signal power is bothersome.

How can we improve on these aspects?

1. Speed - go directly to W^* (LMS Newton)
2. Rattle - Uncouple error from the coefficients (orthogonalize error at each step, not only after adaptation)--- Lattice structures.
3. Compute the Optimum, but instead of doing it algebraically, do it recursively (Recursive Least Squares).