

HSE ESTIMATORS

Cont.(2)

RECURSIVE METHODS

Istituto di Telecomunicazioni

Prof FRANCESCO A. N. PALMIERI

Corso di "Signal Processing and Data Fusion"

RECURSIVE METHODS (GRADIENT SEARCH)

R.1

The solution to the minimum squared error could be computed iteratively without matrix inversion.

The idea is to define the cost function and use a gradient search to find the solution.

Even if the solution, also in the singular case, could be obtained easily by solving the normal equation, the gradient search is very instructive to understand the nature of the stochastic searches that will be introduced later. In fact when the parameters to be found belong to a non linear function (i.e. a neural network), we cannot write the solution in closed form, or in the linear (or affine) case, and the only way to reach a solution is via a gradient search (iterative learning).

Recall that the cost function is

$$E(\underline{c}) = E[(d[n] - \underline{c}^T x[n])^2] = E[d[n]^2] + \underline{c}^T E[x[n]x[n]^T]\underline{c} - 2E[d[n]x[n]]\underline{c}$$
$$- 2E[d[n]x[n]]\underline{c} = E_d + \underline{c}^T R_x \underline{c} - 2\underline{c}^T d_x \underline{c}$$

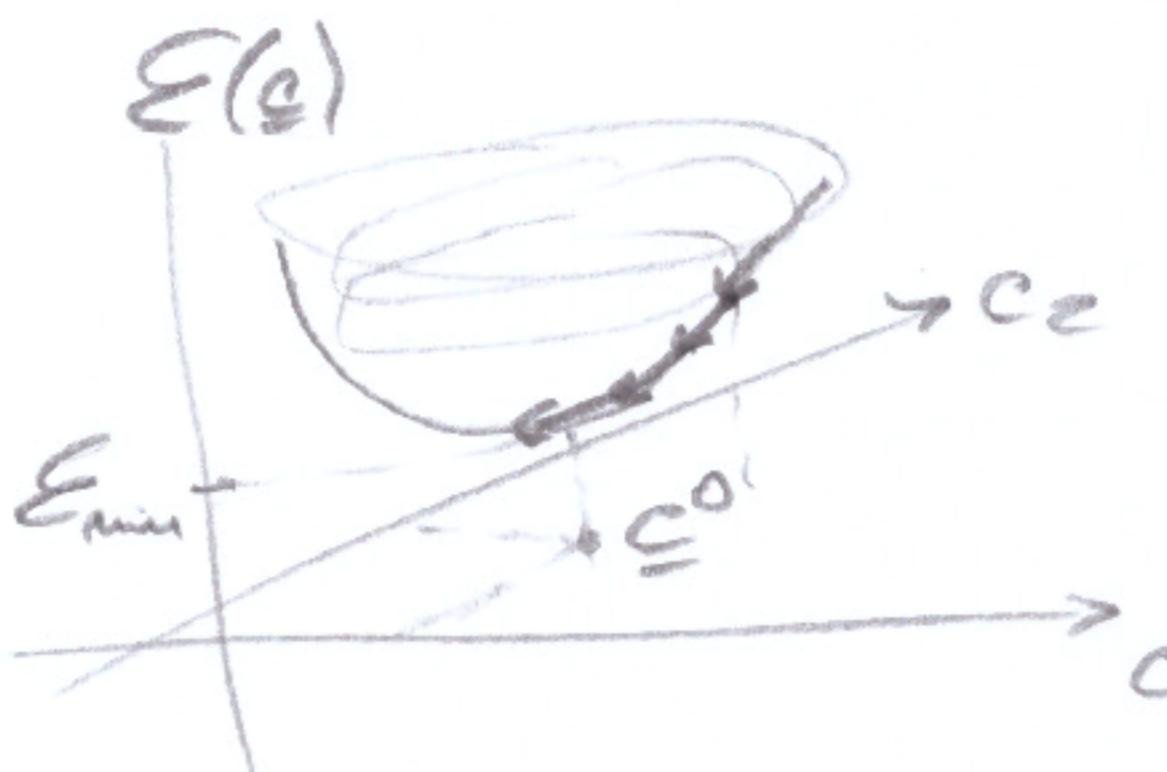
and the gradient

$$\nabla_{\underline{c}} E = 2R_x \underline{c} - 2d_x \quad (*)$$

A gradient search (Newton algorithm) starts from an initial guess $\underline{c}(0)$ and proceeds in small steps in the direction opposite to the current gradient, i.e.

$$\underline{c}(k) = \underline{c}(k-1) - \mu (\underline{R}_x \underline{c}(k-1) - \underline{\tau}_{dx}), \quad (**)$$

where μ is a stepsize parameter that controls the adaptation speed. The search is depicted in the following figure for $N=2$



The function $E(\underline{c})$ is a quadratic form and looks like a bowl. The minimum is unique because the function is convex.
To verify (we did not do it before)
let us compute the Hessian

$$H = (\nabla_{\underline{c}} \nabla_{\underline{c}}^T) E(\underline{c}) = \begin{pmatrix} \frac{\partial}{\partial c_1} & \left[\frac{\partial^2}{\partial c_1 \partial c_1} \frac{\partial^2}{\partial c_1 \partial c_2} \frac{\partial^2}{\partial c_2 \partial c_1} \right] \\ \frac{\partial}{\partial c_2} & \frac{\partial^2}{\partial c_2 \partial c_2} \end{pmatrix} E(\underline{c})$$

$$= \begin{pmatrix} \frac{\partial}{\partial c_1} \\ \vdots \\ \frac{\partial}{\partial c_N} \end{pmatrix} \left[2 \underline{c}^T \underline{R}_x - 2 \underline{\tau}_{dx}^T \right]$$

$$= 2 \underline{R}_x$$

But $\underline{R}_x \geq 0$ (def. positive),
therefore the function is convex everywhere.

The recursive equation (**) obviously has its only stationary point at the solution

$\underline{R}_x \underline{c} - \underline{r}_{dx} = 0$, but we need to choose μ to control the speed, because large μ may make the recursion wander around the solution. Small μ may cause the algorithm to be too slow.



SEARCH DYNAMICS : (COEFFICIENT EVOLUTION)

To understand the evolution of $\{\underline{c}(k)\}$, consider the difference vector

$$\underline{v}(k) = \underline{c}(k) - \underline{c}^0$$

Now we have to verify that $\underline{v}(k) \rightarrow 0$. Substituting in (**), we obtain

$$\underline{v}(k) + \cancel{\underline{c}^0} = \underline{v}(k-1) + \cancel{\underline{c}^0} - \mu (\underline{R}_x(\underline{v}(k-1) + \cancel{\underline{c}^0}) - \underline{r}_{dx})$$

$$\underline{v}(k) = \underline{v}(k-1) - \mu (\underline{R}_x \underline{v}(k-1) + \cancel{\underline{R}_x \cancel{\underline{c}^0}} - \cancel{\underline{r}_{dx}})$$

$$\underline{v}(k) = (\underline{I} - \mu \underline{R}_x) \underline{v}(k-1)$$

The convergence of this recursion is determined by the eigenvalues of $(\underline{I} - \mu \underline{R}_x)$. Use eigenvalue decomposition for \underline{R}_x

$$\underline{R}_x = \underline{Q} \underbrace{\Delta_x}_{\text{eigenvectors}} \underline{Q}^T \quad \Delta_x = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_N) \quad \begin{matrix} \lambda_x \\ \text{eigenvalues} \end{matrix}$$

To get

$$\underline{v}(k) = (\underline{I} - \mu \underline{Q} \underline{\Lambda}_x \underline{Q}^T) \underline{v}(k-1)$$

Define

$$\underline{v}(k) = \underline{Q}^T \underline{v}(k) \quad \begin{array}{l} \text{difference vector in} \\ \text{the eigenspace of } \underline{x} \end{array}$$

$$\underline{Q}^T \underline{v}(k) = \underline{Q}^T (\underline{I} - \mu \underline{Q} \underline{\Lambda}_x \underline{Q}^T) \underline{v}(k-1)$$

$$\underline{v}(k) = (\underline{Q}^T - \mu \underline{Q}^T \underline{\Lambda}_x \underline{Q}^T) \underline{v}(k-1)$$

$$\underline{v}(k) = (\underline{I} - \mu \underline{\Lambda}_x) \underline{v}(k-1)$$

The components of $\underline{v}(k)$ are decoupled

$$v_i(k) = (1 - \mu \lambda_i) v_i(k-1) \quad i=1, \dots, N$$

Convergence to zero is obtained if
 $|1 - \mu \lambda_i| < 1 \quad \forall i$

To make sure that all components converge, we need

$$|1 - \mu \lambda_{\max}| < 1$$

$$-1 < 1 - \mu \lambda_{\max} < 1$$

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

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

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

Each component converges exponentially with factor $(1 - \mu \lambda_i)$

If want to associate a time constant to each R.S mode we set

$$1 - \mu_{di} = e^{-\frac{1}{\tau_i}} \quad \forall i=1, \dots, N$$

$$\tau_i = -\frac{1}{\ln(1-\mu_{di})}$$

τ_i is the time it takes to reach $\frac{1}{e}$ of $v_i(0)$.

If μ_{di} is small $1 - \mu_{di}$ is around 1 and negative using $\ln x \approx x - 1$, $\ln(1 - \mu_{di}) \approx -\mu_{di}$

$$\boxed{\tau_i \approx \frac{1}{\mu_{di}}}$$

More detailed analysis of the evolution of $c(k)$, can be obtained observing that

$$v(k) = Q v(k) \quad (Q \text{ is unitary})$$

Therefore

$$\underline{c}(k) = \underline{c}^0 + \underline{Q} \underline{v}(k)$$

$$\underline{c}(k) = \underline{c}^0 + \sum_{i=1}^N q_i v_i(k)$$

Each element ~~can~~ be written as

$$c_j(k) = c_j^0 + \sum_{i=1}^N q_{ij} v_i(0) (1 - \mu_{di})^k \quad j=1, \dots, N$$

The overall behaviour is the exponentiation of N modes. The slowest mode is governed by the largest eigenvalue and the fastest by the smallest one.

The time to reach the solution will be within the bounds

$$\frac{-1}{\log(1-\mu \lambda_{\min})} < T_a < \frac{-1}{\log(1-\mu \lambda_{\max})}$$

SEARCH DYNAMICS : (SQUARED ERROR EVOLUTION)

In one of the previous derivations, we have showed that the error can be expressed in the following canonical form

$$E(\underline{e}(k)) = \underline{E}(e^0) + (\underline{\underline{e}}(k) - \underline{\underline{e}}^0)^T \underline{\underline{R}}_x \underline{\underline{e}}(k) - \underline{\underline{e}}^0 \quad (1)$$

Using the decomposition introduced above for the coefficient evolution, we have

$$\begin{aligned} E(k) &= E_0 + \underline{\underline{v}}(k)^T \underline{\underline{Q}} \underline{\underline{\Lambda}}_x \underline{\underline{Q}}^T \underline{\underline{v}}(k) \\ &= E_0 + \underline{\underline{v}}(k)^T \underline{\underline{\Lambda}}_x \underline{\underline{v}}(k) \\ &= E_0 + \sum_{i=1}^N (v_i(k))^2 / \lambda_i \\ &= E_0 + \sum_{i=1}^N \lambda_i (1 - \mu \lambda_i)^k v_i^2(0) \end{aligned}$$

Under the same conditions found above
 $(|1 - \mu \lambda_i| < 1)$ $E(k) \rightarrow E_0$

(1) Easy to verify expanding the quadratic form.

Now since the error is the superposition of N nodes, just as before, we have

$$\Sigma_i \varepsilon \approx -\frac{1}{2 \ln(1-\mu_i)} \approx \frac{1}{2 \mu_i}$$

R.7

The mean squared error $\bar{\varepsilon}^2(k)$ converges
is in general twice as fast as the
coefficients.

STOCHASTIC METHOD

(THE LMS ALGORITHM)

The Newton gradient descent method, presented in the previous note, requires the preliminary computation of correlation matrix R_x and ^{the} error-correlation vector \underline{c}_{dx} .

It is possible to avoid this step and get a search that is directly driven by the data. (DATA-BASED METHODS).

We start from considering the availability of a "TRAINING SET" of pairs

$$\{(x[n], d[n]), n=1, \dots\}$$

A current version of the filter $\underline{c}(k)$ correspond to the instantaneous error

$$e[n] = d[n] - \underline{c}^T(k) x[n]$$

and its squared value

$$\hat{e}^2[n] = (d[n] - \underline{c}^T(k) x[n])^2$$

The ^{instantaneous} gradient of $\hat{e}^2[n]$ with respect to $\underline{c}(k)$ is

$$\nabla_{\underline{c}(k)} \hat{e}^2[n] = 2(d[n] - \underline{c}^T(k) x[n]) x[n]$$

Now if we decide to make a small step of adaptation Δn , in the direction opposite to the current gradient, we have

$$c(u) = \underline{c}(u-1) - \mu \underline{x}^T [u] (d[u] - \underline{c}(u-1) \underline{x}^T [u])$$

this is the famous LMS algorithm.

Despite the very rough idea of making one adaptation step on every example, the algorithm converges if μ is small and if we present the examples many times.

To explain why these updates may work, consider the evolution of their means.

$$\begin{aligned} E[c[u]] &= E[c[u-1]] - \mu (E[x^T[u] d[u]] - E[x^T[u] \underline{x}^T[u] c[u-1]]) \\ &= E[c[u-1]] - \mu (\underline{r}_{\text{ax}} - E[x^T[u] y[u]]) \\ &\quad \swarrow \\ &= E[x^T[u] \underline{x}^T[u]] \underline{c}^o \\ &\quad \text{if } y[u] = y^o[u] \\ &\quad \text{in optimal conditions} \end{aligned}$$

Therefore it looks like the mean evolution of the algorithm follows the dynamics studied for the Newton method.

A complete analysis of the LMS algorithm is beyond the scope of these lectures, and can be found in classical texts on adaptive signal processing [Haykin, Adaptive Signal Processing].

In practice the converge will be very noisy and μ has to be carefully controlled to avoid excessive wandering around and excessively low speed. In fact μ can be

most time varying $\mu(k)$ and the state in the training set ~~just usually be~~ presented many times.

Each time the training set is presented, we say that one "epoch" has passed.

Various ideas have been proposed in the literature to obtain convergence and accuracy. For example, the state could be shuffled at every epoch. The step size parameter $\mu(k)$ would be adopted according to the current convergence speed
(ADAM algorithm)

The stochastic search of the LHS algorithm is the precursor of the famous ~~backpropagation~~
backpropagation algorithm that will be ~~in following lectures~~
presented for the neural network architectures where it is impossible to derive analytical solutions. In all cases the search follows a noisy gradient descent path.