

LINEAR NETWORKS

FOR DATA-DRIVEN

REGRESSION

Prof. FRANCESCO A. N. PALMIERI

UNIVERSITA' DELLA CAMPANIA
"LUIGI VANVITELLI"

CORSO DI SIGNAL PROCESSING DATA
FUSION oct 2023

INTRODUCTION

A successful regressor or classifier learned from data depends on the type of parametric function $f(x, \theta)$ chosen. Generally speaking the function has to be sufficiently "powerful" to represent the relations hidden in the data. On the other end, as pointed out before, an excessive parametrization may lead to a less overfitted solution. The issue is not so simple because we are not able to understand directly from the data the nature of the parametric function that may work, unless we try it! Generally, more free parameters the function has, the more powerful its representation capabilities are. Clearly the danger of overfitting is always there because a function with a small number of parameters may be considered "smoother", i.e. providing better generalization. However this is not necessarily the case because even if the number of parameters is large, there are many ways to "regularize" the solution. We will mention some in the following. The choice of the appropriate parametric function class is still now an "art", based on extensive experiments and trials.

(LR1)

LINEAR NETWORKS FOR REGRESSION

The simplest parametric function we could choose for our estimator, is a linear function. To begin our discussion, let us consider a regression problem where $y = f(\underline{x}, \theta) = \underline{c}^T \underline{x} + b$ where y and \underline{c} are continuous values.

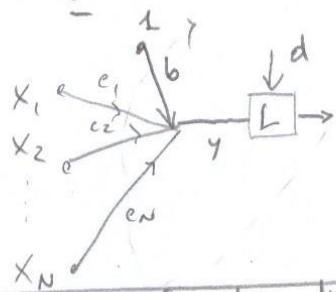
(The function is "affine" and not linear because of the bias. In the following we will ignore in the terminology this detail, referring generally to it as "linear")

In the model-based part of these lectures we have seen how this function emerges as the optimal solution when the likelihood function, or more generally but belongs to the exponential family. Here the linear function is postulated. It has a parametric form with free parameters $\theta = \begin{pmatrix} \underline{c} \\ b \end{pmatrix}$.

that have to be learned from the training set. $Z = \{(\underline{x}[u], d[u]), u=1, \dots, n_Z\}$.

We have to solve the problem

$$\underline{\theta}^0 = \underset{\theta}{\operatorname{argmin}} \sum_{u=1}^{n_Z} L(\underline{c}^T \underline{x}[u] + b, d[u])$$

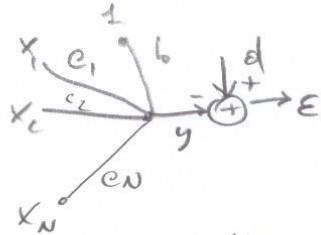


(*) A linear function has to satisfy the superposition principle, i.e. $x_1 \rightarrow y_1; x_2 \rightarrow y_2 \Rightarrow \alpha x_1 + \beta x_2 \rightarrow \alpha y_1 + \beta y_2 \quad \forall \alpha, \beta$. This is not the case for affine functions because $\underline{y} = \underline{c}^T \underline{x} + b$ $\alpha x_1 + \beta x_2 \rightarrow \alpha(\underline{c}^T \underline{x}_1 + b) + \beta(\underline{c}^T \underline{x}_2 + b) = \underline{c}^T \alpha \underline{x}_1 + \alpha b + \underline{c}^T \beta \underline{x}_2 + \beta b$

$$\underline{c}^T (\alpha \underline{x}_1 + \beta \underline{x}_2) + (\alpha + \beta)b \neq \underline{c}^T (\alpha \underline{x}_1 + \beta \underline{x}_2) + \underline{c}^T \alpha \underline{x}_1 + \underline{c}^T \beta \underline{x}_2 + \alpha b + \beta b$$

SOLUTION FOR QUADRATIC LOSS

When the loss function $L(\cdot)$ is quadratic in the difference $\epsilon = d - y$ (most typical), we have (LR2)



$$\hat{\theta}^o = \begin{bmatrix} \hat{e}^o \\ \hat{b}^o \end{bmatrix} = \underset{\begin{bmatrix} \hat{e}^o \\ \hat{b}^o \end{bmatrix}}{\text{argmin}} \sum_{u=1}^{n_2} \frac{(d[u] - (\hat{e}^T x[u] + \hat{b}))^2}{\|e\|^2}$$

The solution here can be found in closed-form as we have seen in deriving FIR filters from data (Wiener). In that case did not include the bias as we considered zero-mean sources. We re-consider the problem here for vector inputs with a bias included and with no hypotheses on the means.

The cost function is

$$\mathcal{E}(e, b) = \sum_{u=1}^{n_2} (d[u] - (e^T x[u] + b))^2$$

Setting the gradients to zero as usual, we have

$$\left\{ \begin{array}{l} \nabla_e \mathcal{E} = 2 \sum_{u=1}^{n_2} (d[u] - (e^T x[u] + b)) x[u] = 0 \\ \frac{\partial \mathcal{E}}{\partial b} = 2 \sum_{u=1}^{n_2} (d[u] - (e^T x[u] + b)) = 0 \\ \sum_{u=1}^{n_2} d[u] x[u] = \left(\sum_{u=1}^{n_2} x[u] x[u]^T \right) e^o + b \sum_{u=1}^{n_2} x[u] \\ \sum_{u=1}^{n_2} d[u] = \left(\sum_{u=1}^{n_2} x[u] \right) e^o + n_2 b^o \end{array} \right.$$

Solving for b^o

$$b^o = \frac{1}{n_2} \left(\sum_{u=1}^{n_2} d[u] - \sum_{u=1}^{n_2} x[u] e^o \right)$$

$$\sum_{u=1}^{n_2} d[u] x[u] = \left(\sum_{u=1}^{n_2} x[u] x[u]^T \right) e^o + \frac{1}{n_2} \sum_{u=1}^{n_2} d[u] \sum_{u=1}^{n_2} x[u]$$

$$- \frac{1}{n_2} \left(\sum_{u=1}^{n_2} x[u] \right) e^o = \sum_{u=1}^{n_2} x[u] - \mu_x u_2$$

Defining the required means (on the training set)

$$\mu_d \triangleq \frac{1}{n_2} \sum_{u=1}^{n_2} d[u] \quad \mu_x \triangleq \frac{1}{n_2} \sum_{u=1}^{n_2} x[u]$$

(R3)

We rewrite

$$\sum_{u=1}^{n_2} d[u] x[u] = \left(\sum_{u=1}^{n_2} x[u] x^T[u] \right) e^0 + \mu_d \mu_x^T n_2$$

$$- \mu_x^T e^0 n_2 / \mu_x$$

$$\sum_{u=1}^{n_2} d[u] x[u] - \mu_d \mu_x^T n_2 = \left(\sum_{u=1}^{n_2} x[u] x^T[u] \right) e^0 - \mu_x \mu_x^T e^0 n_2$$

$$\sum_{u=1}^{n_2} (d[u] - \mu_d)(x[u] - \mu_x) = \sum_{u=1}^{n_2} (x[u] - \mu_x)(x[u] - \mu_x)^T e^0$$

$$e^0 = \left(\sum_{u=1}^{n_2} (x[u] - \mu_x)(x[u] - \mu_x)^T \right)^{-1} \sum_{u=1}^{n_2} (d[u] - \mu_d)(x[u] - \mu_x)$$

n_2 sample covariance matrix n_2 sample errors V_{dx}
 $\sum_{u=1}^{n_2} x[u]$ covariance vector

The bias adds to the output
 The mean of $d[u]$ need subtracts
 from the output the effects of the
 mean of $x[u]$

$$b^0 = \mu_d - \mu_x^T e^0$$

The above solution generalizes the one we discussed for Wiener filters and shows the crucial role of the means and the bias to account for them.

To account for the means we can use one of two steps:

- 1) Remove the means at the beginning and use just a linear estimator e^0 . In this case the sample covariance is the sample autocorrelation matrix and the right-hand side is the error correlation vector.

2) Include a bias in the estimator that substantially takes care of the errors. (Preferred solution) LR4

When the bias is included, a ~~compact~~^{more compact} formula is obtained by considering an augmented input vector

$$\underline{x}[n] = \begin{bmatrix} 1 \\ \underline{x}[n] \end{bmatrix} \quad \begin{array}{l} \text{(we will use this)} \\ \text{"trick" in many neural network structures} \end{array}$$

and a linear mapping with an $(N+1)$ -dimensional vector \underline{w} .

$$y = \underline{w}^T \underline{x}[n] = \begin{bmatrix} b & c_1 & \dots & c_N \end{bmatrix} \begin{bmatrix} 1 \\ \underline{x}_1[n] \\ \underline{x}_2[n] \\ \vdots \\ \underline{x}_N[n] \end{bmatrix}$$

The solution is immediately derived in compact form as:

$$\underline{w}^* = \left(\sum_{n=1}^{N_a} \underline{x}_a[n] \underline{x}_a^T[n] \right)^{-1} \left(\sum_{n=1}^{N_a} \underline{x}_a[n] d[n] \right) \quad (*)$$

$$\boxed{\underline{w}^* = \underline{R}_{x_a}^{-1} \underline{\Sigma}_{x_a d}}$$

Note how the matrix $\frac{1}{N} \sum_{n=1}^{N^2} \underline{x}_a[n] \underline{x}_a^T[n] = R_{x_a}$ contains means and correlations (sample)

$$\begin{bmatrix} 1 \\ \underline{x}[n] \end{bmatrix} \begin{bmatrix} 1 & \underline{x}^T[n] \end{bmatrix} = \begin{bmatrix} 1 \\ x_1[n] \\ x_2[n] \\ \vdots \\ x_N[n] \end{bmatrix} \begin{bmatrix} 1 & x_1[n] x_2[n] \dots x_N[n] \\ x_1^2[n] & \ddots & & \\ x_2^2[n] & \ddots & \ddots & \\ \vdots & \ddots & \ddots & x_2[n] x_N[n] \\ x_N^2[n] & x_N[n] x_1[n] & x_N[n] x_2[n] & \ddots & x_N^2[n] \end{bmatrix}$$

$$= \begin{bmatrix} 1 & x_1[n] & x_2[n] & \dots & x_N[n] \\ x_1[n] & x_1^2[n] & x_1[n] x_2[n] & \dots & x_1[n] x_N[n] \\ x_2[n] & x_2[n] x_1[n] & x_2^2[n] & \dots & x_2[n] x_N[n] \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_N[n] & x_N[n] x_1[n] & x_N[n] x_2[n] & \dots & x_N^2[n] \end{bmatrix}$$

$$R_{x_a} = \frac{1}{N^2} \sum_n \underline{x}_a[n] \underline{x}_a^T[n]$$

$$= \frac{1}{N^2} \begin{bmatrix} 1 & \sum_n x_1[n] - \sum_n x_2[n] & \dots & \sum_n x_N[n] \\ \sum_n x_1[n] & \sum_n x_1^2[n] & \sum_n x_1[n] x_2[n] & \dots & \sum_n x_1[n] x_N[n] \\ \sum_n x_2[n] & \sum_n x_2[n] x_1[n] & \sum_n x_2^2[n] & \dots & \sum_n x_2[n] x_N[n] \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sum_n x_N[n] & \sum_n x_N[n] x_1[n] & \sum_n x_N[n] x_2[n] & \dots & \sum_n x_N^2[n] \end{bmatrix}$$

and the vector $\frac{1}{N^2} \sum_{n=1}^{N^2} \underline{x}_a[n] d[n] = \underline{x}_{ad}$ contains the mean of $d[n]$ and the cross-correlations between $d[n]$ and $x[n]$.

$$\underline{x}_{ad} = \frac{1}{N^2} \sum_n \underline{x}_a[n] d[n]$$

$$= \frac{1}{N^2} \begin{bmatrix} \sum_n d[n] \\ \sum_n d[n] x_1[n] \\ \sum_n d[n] x_2[n] \\ \vdots \\ \sum_n d[n] x_N[n] \end{bmatrix}$$

(All the necessary information is included)

LR6

MINIMUM MEAN SQUARED ERROR (MMSE) ON THE TRAINING SET

The cost function $E(\underline{w})$ in optimal conditions is expressed in one of the two forms:

$$E(\underline{w}^*) = E_d - \underline{x}_{ad}^T \underline{w}^* \quad (\text{I})$$

or

$$E(\underline{w}^*) = E_d - \underline{x}_{ad}^T R_x^{-1} \underline{x}_{ad} \quad (\text{II})$$

$$E_d = \frac{1}{m} \sum_{u=1}^m d[u]^2$$

Note that expression II predicts the MMSE without computing \underline{w}^* explicitly.

Proof By substitution

$$E(\underline{w}^*) = \frac{1}{m} \sum_{u=1}^m (d[u] - \underline{x}_{ad}^T \underline{w}^*)^2 = \frac{1}{m} \sum_{u=1}^m d[u]^2 + \underline{w}^T \frac{\underline{x}_{ad}^T \underline{x}_{ad}}{m} \underline{w}^* - 2 \sum_{u=1}^m d[u] \underline{x}_{ad}^T \underline{w}^*$$

$$\begin{aligned} &= \frac{1}{m} \sum_{u=1}^m d[u]^2 + \underline{w}^T \sum_{u=1}^m \underline{x}_{ad}^T d[u] \left(\frac{1}{m} \sum_{u=1}^m \underline{x}_{ad} \underline{x}_{ad}^T \right) \left(\frac{1}{m} \sum_{u=1}^m \underline{x}_{ad} \underline{x}_{ad}^T \right) \underline{w}^* \\ &\quad - 2 \sum_{u=1}^m d[u] \underline{x}_{ad}^T \underline{w}^* \\ &= \frac{1}{m} \sum_{u=1}^m d[u]^2 - \frac{1}{m} \sum_{u=1}^m \underline{x}_{ad}^T d[u] \left(\frac{1}{m} \sum_{u=1}^m \underline{x}_{ad} \underline{x}_{ad}^T \right) \underline{w}^* \end{aligned}$$

$$E(\underline{w}^*) = E_d - \underline{x}_{ad}^T \underline{w}^* \quad (\text{I})$$

$$\begin{aligned} &= \frac{1}{m} \sum_{u=1}^m d[u]^2 - \frac{1}{m} \sum_{u=1}^m \underline{x}_{ad}^T d[u] \left(\frac{1}{m} \sum_{u=1}^m \underline{x}_{ad} \underline{x}_{ad}^T \right)^{-1} \left(\frac{1}{m} \sum_{u=1}^m \underline{x}_{ad}^T d[u] \right) \\ &= E_d - \underline{x}_{ad}^T R_x^{-1} \underline{x}_{ad} \end{aligned}$$

$$E(\underline{w}^*) = E_d - \underline{x}_{ad}^T R_x^{-1} \underline{x}_{ad} \quad (\text{II})$$

LRF

MEAN SQUARED ERROR ON THE VALIDATION SET

After computing the best solution on the training set, as in any machine learning application, we have to check how good the solution is on the validation set. Note that we cannot use expressions (I) and (II). That's because we work only on the training set, but we have to plug in the solution in the general cost function expression

$$\mathcal{E}^V(\underline{w}^0) = \frac{1}{m_V} \sum_{m=1}^{m_V} (d[m] - \underline{w}^0 \underline{x}_a[m])^2$$

that we can rewrite as

$$\begin{aligned} \mathcal{E}^V(\underline{w}^0) &= \frac{1}{m_V} \sum_{m=1}^{m_V} d[m]^2 + \underline{w}^0 \sum_{m=1}^{m_V} \underline{x}_a[m] \underline{x}_a[m]^T \underline{w}^0 \\ &\quad + \frac{2}{m_V} \sum_{m=1}^{m_V} \underline{x}_a[m]^T d[m] \underline{w}^0 \end{aligned}$$

i.e. \underline{w}^0 is computed only on the training set, but $\mathcal{E}^V = R_{x_a} \mathcal{E}_{x_a}$ must be computed for the validation set.

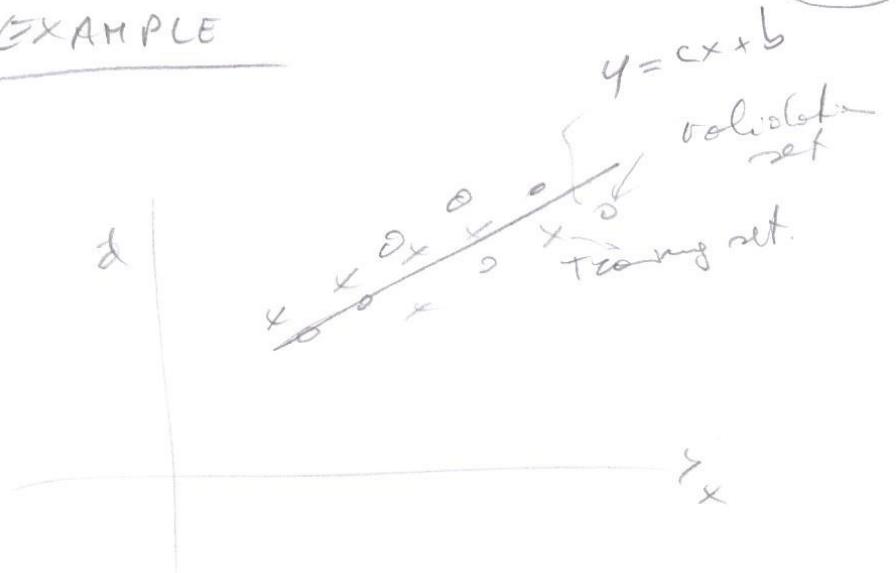
$$\mathcal{E}^V(\underline{w}^0) = \mathcal{E}_d + \underline{w}^0 R_{x_a} \underline{w}^0 - 2 \underline{x}_{x_a}^T \underline{w}^0$$

$$\left\{ \begin{array}{l} \mathcal{E}_d = \frac{1}{m_V} \sum_{m=1}^{m_V} d[m]^2 ; \quad R^V = \frac{1}{m_V} \sum_{m=1}^{m_V} \underline{x}_a[m] \underline{x}_a[m]^T \\ \underline{x}_{x_a} = \frac{1}{m_V} \sum_{m=1}^{m_V} \underline{x}_a[m] d[m] \end{array} \right.$$

In checking for generalization we compute both \mathcal{E} and \mathcal{E}^V . It is usually better compare their square root (RMS value) $\sqrt{\mathcal{E}}$ and $\sqrt{\mathcal{E}^V}$ because RMS values are more easily interpretable as they are in the same orders of magnitude of d and y .

(LR 8)

EXAMPLE

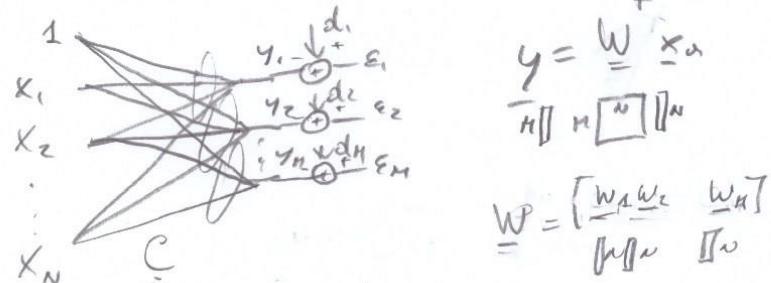


Compare E and E'

THE MULTI-OUTPUT CASE

(LRG)

Consider now the straightforward extension to the case where y and d are M -dimensional vectors



For economy of notation we consider the input to be the vector \underline{x}_a augmented from x , that includes the one in first position. Therefore the matrix \underline{W} already includes the biases and the estimator is affine. (not linear in x_a) The problem of only adding the contributions of the M outputs for the cut is

$$\underline{W}^0 = \underset{\underline{W}}{\text{argmin}} \sum_{j=1}^M \sum_{u=1}^{n_{\text{obs}}} (d_j[u] - \underline{W}_j^T \underline{x}_a[u])^2$$

which can be written in matrix form as

$$\underline{W}^0 = \underset{\underline{W}}{\text{argmin}} \sum_{u=1}^{n_{\text{obs}}} (\underline{d}[u] - \underline{W}^T \underline{x}_a[u])^T (\underline{d}[u] - \underline{W}^T \underline{x}_a[u])$$

Since each "filter" \underline{W}_j operates independently from the others, we can use the results of the previous section to write the relations

$$(*) \quad \underline{W}_j^0 = \underline{R}_{x_a}^{-1} \underline{\epsilon}_{x_a \text{adj}} \quad \text{where} \quad \left\{ \begin{array}{l} \underline{R}_{x_a} = \sum_{u=1}^{n_{\text{obs}}} \underline{x}_a[u] \underline{x}_a[u]^T \\ \underline{\epsilon}_{x_a \text{adj}} = \sum_{u=1}^{n_{\text{obs}}} \underline{x}_a[u] d_j[u] \end{array} \right.$$

(LR10)

In compact matrix form the solution can be written as

$$\underline{W}^o = \underline{R}_{x_a}^{-1} \underline{R}_{x_{adj}} \quad \text{with } \underline{R}_{x_{adj}} = \sum_{u=1}^{M_v} \underline{x}_{[u]} \underline{d}_{[u]}^T$$

$\begin{matrix} N \\ \square \\ \square \\ N \end{matrix} \quad \begin{matrix} N \\ \square \\ \square \\ N \end{matrix} \quad \begin{matrix} M \\ \square \end{matrix}$

cross correlation
matrix

The error on the training set, as in previous section, is for output j

$$\mathcal{E}_j(\underline{W}_j) = \frac{1}{N} \sum_{u=1}^{M_v} d_j^2[u] - \underline{\mathcal{X}}_j^T \underline{R}_{x_a}^{-1} \underline{\mathcal{X}}_j \quad j=1, \dots, M$$

As pointed out above, \mathcal{E}_j

There is a difference in the error computed on the training set and the one computed on the validation (test) set. $V = \{d[m], x[m]\}, m=1, \dots, M_v\}$ that is written as

$$\begin{aligned} \mathcal{E}(\underline{W}_j) &= \frac{1}{M_v} \sum_{m=1}^{M_v} (d_j[m] - \underline{W}_j^T \underline{x}[m])^2 \\ &= \frac{1}{M_v} \sum_{m=1}^{M_v} d_j^2[m] + \underline{W}_j^T \left(\sum_{m=1}^{M_v} \underline{x}[m] \underline{x}[m]^T \right) \underline{W}_j - 2 \sum_{m=1}^{M_v} d_j[m] \underline{x}[m]^T \end{aligned}$$

$\begin{matrix} V \\ \sum \\ \underline{d}_j \\ \underline{x}_j \\ \underline{x}_{adj} \end{matrix} \quad \begin{matrix} M_v \\ \sum \\ \underline{x}[m] \\ \underline{x}[m]^T \\ \underline{W}_j \end{matrix} \quad \begin{matrix} V \\ \sum \\ \underline{d}_j \\ \underline{x}_j \\ \underline{x}_{adj} \end{matrix}$

NOTATION:

When not specified $\underline{R}_{x_a} \underline{\mathcal{X}}_{adj}$ are computed on the training set $\underline{\mathcal{X}}$. $\underline{R}_{x_a}^V \underline{\mathcal{X}}_V$ are the ones computed on the validation set.

OTHER STATISTICS ON THE SOLUTION:

(LRII)

It is often useful to compute other statistics
 (other than the errors) obtained from the training set
 (on the solution obtained from the training set).

MEANS

For example the output mean

on the validation set and on the validation set
 on the training set and on the validation set

$$\underline{\mu}_y = \underline{W}^T \underline{\mu}_{x_a} \quad \underline{\mu}_y^v = \underline{W}^T \underline{\mu}_{x_a}^v \quad (\text{Note that they are not the same.})$$

The mean error on the training set is optimally

$$\underline{\mu}_E = \frac{1}{n} \sum_{n=1}^{n=2} (d[n] - \underline{W}^T \underline{\mu}_{x_a}[n]) = \underline{\mu}_d - \underline{W}^T \underline{\mu}_{x_a}$$

But since $\underline{x}_a[n]$ is augmented with "1", we
 have

$$\underline{\mu}_{x_a} = \begin{bmatrix} 1 \\ \underline{\mu}_{x_1} \\ \underline{\mu}_{x_2} \\ \vdots \\ \underline{\mu}_{x_n} \end{bmatrix} = \begin{bmatrix} 1 \\ \underline{\mu}_x \end{bmatrix} \quad \underline{W}^e = [\underline{w}_1^T \underline{w}_2^T \dots \underline{w}_n^T] = \begin{bmatrix} \underline{b}_1^T & \underline{b}_2^T & \dots & \underline{b}_n^T \\ \underline{c}_1^T & \underline{c}_2^T & \dots & \underline{c}_n^T \end{bmatrix}$$

$$\underline{\mu}_E^o = \underline{\mu}_d - \begin{bmatrix} \underline{b}_1^T \underline{c}_1^T \\ \underline{b}_2^T \underline{c}_2^T \\ \vdots \\ \underline{b}_n^T \underline{c}_n^T \end{bmatrix} \begin{bmatrix} 1 \\ \underline{\mu}_x \end{bmatrix} = \underline{\mu}_d - \begin{bmatrix} \underline{b}_1^T \underline{c}_1^T \underline{\mu}_x^T \\ \underline{b}_2^T \underline{c}_2^T \underline{\mu}_x^T \\ \vdots \\ \underline{b}_n^T \underline{c}_n^T \underline{\mu}_x^T \end{bmatrix}$$

From the previous section the optimal bias and

$$b_j^o = \mu_{dj} - \underline{c}_j^T \underline{\mu}_x$$

Therefore the errors have zero mean on
 the training set $\underline{\mu}_E^o = 0$ (but they are empirically unbiased)

These means are not necessarily null on
 the validation set

$$\underline{\mu}_E^o = \begin{bmatrix} \mu_{d_1}^o - b_1^o - c_1^T \underline{\mu}_x^o \\ \mu_{d_2}^o - b_2^o - c_2^T \underline{\mu}_x^o \\ \vdots \\ \mu_{d_M}^o - b_M^o - c_M^T \underline{\mu}_x^o \end{bmatrix} \neq 0$$

The question of the means is in practical machine learning solutions often overlooked. It is very important that the data that we present for training, for validation and therefore testing, be properly "normalized". In other words the pairs $(x_{t,i}, d_{t,i})$ in \mathcal{T} , V and T have to be in the same ranges (measured averages). Otherwise performances may be very poor. We will return on this issue later on, but for now we should observe that different means on the training set and on the validation set may produce unsatisfactory results. Therefore

- either : ~~remove~~
- Remove the means from the data using no biases and perhaps add them again at the end (μ_d). Note that this can be done on means computed on the training set and/or on the validation set inducing confusion on the results.
 - Take care that the data in \mathcal{T} , V , and T are in the same ranges and are comparable in averages and values. In this way we can let the biases take care of the means on the training set, and expect good stable working on the other sets.

COVARIANCES

In studying the solution after optimization, it may be useful to look beyond the means and analyse variation parameters, such as variances and covariances. More specifically, the error variance, discussed before may give us only a partial account of the solution: variables and put components may be better estimated than others and their values may be mutually correlated. Therefore we can consider the error covariance matrix on the leaving set.

$$\sum_{\varepsilon^0} = \frac{1}{n_2} \sum_{n=1}^{n_2} (\underline{\varepsilon}^0[n] - \underline{\mu}_{\varepsilon}) (\underline{\varepsilon}^0[n] - \underline{\mu}_{\varepsilon})^T$$

($\underline{\varepsilon}^0[n]$ is the error in optimal conditions)

$$\text{where } \underline{\varepsilon}^0[n] = \underline{d}[n] - \underline{W}^T \underline{x}_d[n]$$

But we have already shown that on the leaving set $\underline{\mu}_{\varepsilon} = \underline{0}$, therefore

$$\sum_{\varepsilon^0} = \frac{1}{n_2} \sum_{n=1}^{n_2} \underline{\varepsilon}^0[n] \underline{\varepsilon}^0[n]^T$$

Note that on the diagonal \sum_{ε^0} contains the variances of the errors on the various components, and the covariances on the other off-diagonal elements.

$$\sum_{\varepsilon^0} = \begin{pmatrix} \sigma_{\varepsilon_1^0}^2 & \sigma_{\varepsilon_1^0 \varepsilon_2^0} & \dots & \sigma_{\varepsilon_1^0 \varepsilon_n^0} \\ \sigma_{\varepsilon_2^0 \varepsilon_1^0} & \sigma_{\varepsilon_2^0}^2 & \dots & \sigma_{\varepsilon_2^0 \varepsilon_n^0} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{\varepsilon_n^0 \varepsilon_1^0} & \sigma_{\varepsilon_n^0 \varepsilon_2^0} & \dots & \sigma_{\varepsilon_n^0}^2 \end{pmatrix}$$

Inspection of the variances on the diagonal can tell us how accurate are the estimates on the various components (on the training set). LR14

Σ_{ε^0} does not need to be computed from the data ~~because~~ ^{training set} we can write it in terms of the computed solution and standards that we have already. More specifically

$$\begin{aligned}\Sigma_{\varepsilon^0} &= \frac{1}{m_2} \sum_{u=1}^{m_2} \underline{\varepsilon}^0[u] \underline{\varepsilon}^{0T}[u] = \\ &= \frac{1}{m_2} \sum_{u=1}^{m_2} (\underline{d}[u] - \underline{W}_{x_{ad}}^{0T} \underline{x}_{ad}[u]) (\underline{d}[u] - \underline{W}_{x_{ad}}^{0T} \underline{x}_{ad}[u])^T \\ &= \frac{1}{m_2} \sum_{u=1}^{m_2} \underline{d}[u] \underline{d}^{T}[u] + \underline{W}_{x_{ad}}^{0T} \sum_{u=1}^{m_2} \underline{x}_{ad}[u] \underline{x}_{ad}^{T}[u] \underline{W}_{x_{ad}}^0 \\ &\quad - \underline{W}_{x_{ad}}^{0T} \left(\frac{1}{m_2} \sum_{u=1}^{m_2} \underline{x}_{ad}[u] \underline{d}^{T}[u] - \frac{1}{m_2} \sum_{u=1}^{m_2} \underline{d}[u] \underline{x}_{ad}^{T}[u] \right) \underline{W}_{x_{ad}}^0\end{aligned}$$

Since $\underline{W}_{x_{ad}}^0 = \underline{R}_{x_{ad}}^{-1} \underline{R}_{x_{ad}}$

$$\Sigma_{\varepsilon^0} = \underline{R}_d + \underline{R}_{x_{ad}}^{T} \underline{R}_{x_{ad}}^{-1} \underline{R}_{x_{ad}} \underline{R}_{x_{ad}}^{-1} - \underline{R}_{x_{ad}}^{T} \underline{R}_{x_{ad}}^{-1} \underline{R}_{x_{ad}} - \underline{R}_{x_{ad}}^{T} \underline{R}_{x_{ad}}^{-1} \underline{R}_{x_{ad}}$$

$$\begin{aligned}\Sigma_{\varepsilon^0} &= \underline{R}_d - \underline{R}_{x_{ad}}^{T} \underline{R}_{x_{ad}}^{-1} \underline{R}_{x_{ad}} \\ \text{or} \\ \Sigma_{\varepsilon^0} &= \underline{R}_d - \underline{W}_{x_{ad}}^{0T} \underline{R}_{x_{ad}} = \underline{R}_d - \underline{R}_{x_{ad}}^{T} \underline{W}^0\end{aligned}$$

On the validation set the situation is (LR15)

different because the error covariance for the computed solution is

$$\sum_{\varepsilon^0}^V \triangleq \frac{1}{mV} \sum_{n=1}^{mV} (\underline{\varepsilon}^0[n] - \underline{\mu}_{\varepsilon^0}^V)(\underline{\varepsilon}^0[n] - \underline{\mu}_{\varepsilon^0}^V)^T$$

$\underline{\mu}_{\varepsilon^0}^V$ is not necessarily null, as pointed out above. We have

$$\underline{\varepsilon}^0[n] = \underline{o}[n] - \underline{W}^0 \underline{x}_a[n]$$

$$\underline{\mu}_o^V = \frac{1}{mV} \sum_{n=1}^{mV} \underline{o}[n]$$

$$\underline{\mu}_{x_a}^V = \underline{\mu}_o^V - \underline{W}^0 \underline{\mu}_{x_a}^V$$

$$\underline{\mu}_{x_a}^V = \frac{1}{mV} \sum_{n=1}^{mV} \underline{x}_a[n]$$

\sum_{ε^0} must be computed on the data and its elements

$$\sum_{\varepsilon^0}^V = \begin{pmatrix} \sigma_{\varepsilon_1^0}^2 & \sigma_{\varepsilon_1 \varepsilon_2}^0 & \dots & \sigma_{\varepsilon_1 \varepsilon_m}^0 \\ \sigma_{\varepsilon_2 \varepsilon_1}^0 & \sigma_{\varepsilon_2^0}^2 & & \\ \vdots & & \ddots & \\ \sigma_{\varepsilon_m \varepsilon_1}^0 & & & \sigma_{\varepsilon_m^0}^2 \end{pmatrix}$$

you tell how well the system works in generalization to \sum_{ε^0} obtained on the training set.

SOLUTION WITH THE PSEUDOINVERSE

BR16

Last sources solutions to linear problems have a long history in the literature and they have been formulated in all kinds of ways. One of these is to see the linear regression problem as the approximate solutions to an over-determined linear system of equations. Consider for now the case with $M=1$ and with \underline{x} already augmented

$$\begin{matrix} & \vdots & & \\ \begin{matrix} 1 & x_1[1] & x_2[1] & \dots & x_N[1] \\ 1 & x_1[2] & x_2[2] & \dots & x_N[2] \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1[m_2] & x_2[m_2] & \dots & x_N[m_2] \end{matrix} & \begin{matrix} w_0 \\ w_1 \\ \vdots \\ w_N \end{matrix} & = & \begin{matrix} d[1] \\ d[2] \\ \vdots \\ d[m_2] \end{matrix} \\ \underline{x}_m & \underline{w} & \approx & \underline{d} \end{matrix}$$

Clearly the match between the two sides of the equation cannot be obtained exactly. We assume $m_2 \gg N$ and the system is over-determined. Therefore we look for a solution that minimizes the sum of the squares of the differences between the two members of the equation.

$$\|\underline{x}\underline{w} - \underline{d}\|_2^2 = \text{square } \frac{1}{m_2} \|\underline{x}\underline{w} - \underline{d}\|^2$$

This is exactly the cost function $E(w)$ used in the previous sections, written as a function of the large (full) matrix \underline{x}_m and the full vector \underline{d} .

$$\mathcal{E}(\underline{\underline{w}}) = \frac{1}{m_2} (\underline{\underline{X}} \underline{\underline{w}} - \underline{\underline{d}})^T (\underline{\underline{X}} \underline{\underline{w}} - \underline{\underline{d}}) = \frac{1}{m_2} (\underline{\underline{w}}^T \underline{\underline{X}}^T \underline{\underline{X}} \underline{\underline{w}} - \underline{\underline{d}}^T \underline{\underline{X}} \underline{\underline{w}} - \underline{\underline{w}}^T \underline{\underline{X}} \underline{\underline{d}} + \underline{\underline{d}}^T \underline{\underline{d}}) \quad (LR17)$$

$$= \frac{1}{m_2} (\underline{\underline{w}}^T \underline{\underline{X}}^T \underline{\underline{X}} \underline{\underline{w}} - 2 \underline{\underline{d}}^T \underline{\underline{X}} \underline{\underline{w}} + \underline{\underline{d}}^T \underline{\underline{d}})$$

Taking the gradient with respect to $\underline{\underline{w}}$ (*), and setting it to zero, we have

$$2 \underline{\underline{X}}^T \underline{\underline{X}} \underline{\underline{w}} - 2 \underline{\underline{X}}^T \underline{\underline{d}} = 0$$

The solution is then

$$\underline{\underline{w}}^* = (\underline{\underline{X}}^T \underline{\underline{X}})^{-1} \underline{\underline{X}}^T \underline{\underline{d}}$$

We can write the solution using the definition of pseudo-inverse (pinv(.) in Matlab)

$$\underline{\underline{X}}^{\#} = (\underline{\underline{X}}^T \underline{\underline{X}})^{-1} \underline{\underline{X}}^T$$

$$\underline{\underline{w}}^* = \underline{\underline{X}}^{\#} \underline{\underline{d}}$$

To recognize that the solution is exactly the same derived in the previous sections, note that

$$\underline{\underline{X}}^T \underline{\underline{X}} = \begin{bmatrix} 1 & & & \\ X_1[1] X_1[1] & \cdots & X_1[m_2] \\ \vdots & & \vdots \\ X_N[1] X_N[1] & & X_N[m_2] \end{bmatrix} \begin{bmatrix} 1 & X_1[1] X_2[1] \dots X_N[1] \\ 1 & X_1[2] X_2[2] \dots X_N[2] \\ \vdots & \vdots \\ 1 & X_1[m_2] X_2[m_2] \dots X_N[m_2] \end{bmatrix} = m_2 R_{k_m}$$

(*) We have used the well-known results from differential matrix calculus

$$\nabla_{\underline{\underline{x}}} \underline{\underline{b}}^T = \underline{\underline{b}} ; \quad \nabla_{\underline{\underline{x}}} \underline{\underline{x}}^T \underline{\underline{A}} \underline{\underline{x}} = 2 \underline{\underline{A}} \underline{\underline{x}}$$

$$\underline{X}_n^T \underline{d} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ X_1[1] & X_1[2] & \dots & X_1[m] \\ \vdots & & & \vdots \\ X_N[1] & X_N[2] & \dots & X_N[m] \end{bmatrix} \begin{bmatrix} \underline{d}[1] \\ \underline{d}[2] \\ \vdots \\ \underline{d}[m] \end{bmatrix} = n_2 \underline{\Sigma}_{x_{ad}}$$

(LR18)

and there is a perfect match with equation (*)
The total squared error in optimal conditions on
the learning set is also

$$\begin{aligned} E(\underline{w}^o) &= \frac{1}{n_2} (\underline{w}^o)^T \underline{X}_a^T \underline{X}_a \underline{w}^o - 2 \underline{d}^T \underline{X}_a \underline{w}^o + \underline{d}^T \underline{d} \\ &= \frac{1}{n_2} \left(\underline{d}^T \underline{X}_a \left(\underline{X}_a^T \underline{X}_a \right)^{-1} \underline{X}_a \underline{d} - 2 \underline{d}^T \underline{X}_a \left(\underline{X}_a^T \underline{X}_a \right)^{-1} \underline{w}^o \right. \\ &\quad \left. + \underline{d}^T \underline{d} \right) = \frac{1}{n_2} \underline{d}^T \underline{d} - \frac{1}{n_2} \underline{d}^T \underline{X}_a \left(\underline{X}_a^T \underline{X}_a \right)^{-1} \underline{w}^o \\ &= E_d - \frac{1}{n_2} \underline{d}^T \underline{X}_a \underline{w}^o = E_d - \underline{\Sigma}_{x_{ad}}^T \underline{R}_{x_{ad}}^{-1} \underline{x}_{ad} (\underline{x}) \end{aligned}$$

obviously the same obtained in the previous section.
Note again that the total error on the validation
set is different and cannot be computed as in (*),
but using the general expression

$$E(\underline{w}^o) = \frac{1}{n_v} \left(\underline{w}^o^T \underline{X}_v^T \underline{X}_v \underline{w}^o - 2 \underline{d}^T \underline{X}_v \underline{w}^o + \underline{d}^T \underline{d} \right)$$

where \underline{X}_v and \underline{d}^v contain the validation set
examples.

The extension to the multi-output case is straight forward
considering the system.

(LR19)

$$\begin{bmatrix} 1 & X_1[1] & X_2[1] & \dots & X_N[1] \\ 1 & X_1[2] & X_2[2] & \dots & X_N[2] \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_1[m] & X_2[m] & \dots & X_N[m] \end{bmatrix} \begin{bmatrix} W_1 & W_2 & \dots & W_M \end{bmatrix} = \begin{bmatrix} d_1[1] & d_2[1] & \dots & d_N[1] \\ d_1[2] & d_2[2] & \dots & d_N[2] \\ \vdots & \vdots & \ddots & \vdots \\ d_1[m] & d_2[m] & \dots & d_N[m] \end{bmatrix}$$

$$X \quad W \quad D$$

Where we have arranged also the desired outputs in
a unique matrix D

$$W^o = X^{\#} D = (X^T X)^{-1} X^T D$$

To recover the correspondence to the solution found in the
previous section is again immediate if we
observe that

$$R_{x_1} = \frac{1}{m_2} X^T X \quad \text{and} \quad R_{x_{ad}} = \frac{1}{m_2} X^T D$$

The output ~~error covariance~~, that contains the
mean squared errors on the components in the
diagonal, can be written for the training set in
optimal condition directly as a function of X and D .

$$\sum \varepsilon^o = \frac{1}{m_2} D^T D - \frac{1}{m_2} D^T X (X^T X)^{-1} X^T D$$

or $\boxed{\boxed{}}$

$$\sum \varepsilon^o = \frac{1}{m_2} D^T D - \frac{1}{m_2} D^T W^o = \frac{1}{m_2} D^T D - \frac{1}{m_2} W^T D$$

On the validation set, the error mean vector
 (not necessarily null) is

LR20

$$\hat{\mu}_{\varepsilon^o}^v = \frac{1}{n_v} \left(\underline{X}^v \underline{W}^o - \underline{\hat{D}}^v \right)^T \underline{e}_{n_v}$$

where $\underline{e}_{n_v} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}_{n_v}$

and the error covariance

$$\hat{\Sigma}_{\varepsilon^o}^v = \frac{1}{n_v} (\underline{X}^v \underline{W}^o - \underline{\hat{D}}^v)(\underline{X}^v \underline{W}^o - \underline{\hat{D}}^v)^T - \hat{\mu}_{\varepsilon^o}^v \hat{\mu}_{\varepsilon^o}^{v^T}$$

The compact matrix expressions in terms of $\underline{X}, \underline{D}$,
 $\underline{X}^v \underline{D}^v$, can be quite useful because they provide
 or closed-form immediate solution for
 regression. This is provided that the training
 set and the validation set can be computed in
 large matrices.

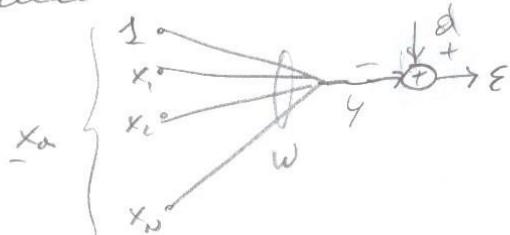
LEARNING THE LINEAR REGRESSOR WITH A GRADIENT SEARCH

LR21

We have seen in the previous sections how the solution to the minimization of the sum of squares (MSE) for the linear network can be found in closed-form. It is however very instructive to consider gradient search also for this case for two main reasons:

- (1) Data sets may not be available in a batch (as in linear filtering)
- (2) Most importantly, the nature of the network already contains some elements that are important to understand gradient searches for more complex architectures.
- (3) We may use loss functions different from the MSE for which closed-form solutions do not exist.

With reference to the scalar case with $M=1$



we seek to find

$$\underline{w}^* = \underset{\underline{w}}{\operatorname{arg\,min}} E(\underline{w}) = \operatorname{arg\,min}_{\underline{w}} \frac{1}{M} \sum_{m=1}^M \Psi(d[m] - \underline{w}^T \underline{x}_o[m])$$

where $\Psi(\cdot)$ is one of the loss functions defined in one of the previous lectures.
(The most typical being $\Psi(s) = s^2$)

A gradient search follows the recurrence LR22

$$\underline{w}[k] = \underline{w}[k-1] - \mu(k) \nabla_{\underline{w}} E(\underline{w}[k-1])$$

where

$$E(\underline{w}[k-1]) = \frac{1}{m} \sum_{n=1}^m \Psi(d[n] - \underline{w}^T[k-1] \underline{x}_a[n])$$

$$\nabla_{\underline{w}} E(\underline{w}[k-1]) = \frac{1}{m} \sum_{n=1}^m \Psi'(d[n] - \underline{w}^T[k-1] \underline{x}_a[n]) (-\underline{x}_a[n])$$

When $\Psi(\xi) = \xi^2$, we have

$$\underline{w}[k] = \underline{w}[k-1] + \mu(k) \frac{2}{m} \sum_{n=1}^m (d[n] - \underline{w}^T[k-1] \underline{x}_a[n]) \underline{x}_a[n]$$

Since the factor $\frac{2}{m}$ can be absorbed into $\mu(k)$

the Batch-LMS algorithm is written as

$$\underline{w}[k] = \underline{w}[k-1] + \mu(k) \sum_{n=1}^m (d[n] - \underline{w}^T[k-1] \underline{x}_a[n]) \underline{x}_a[n]$$

Note that
 $\underline{w}^T[k-1] \underline{x}_a[n]$
is the output
 $y[n]$ computed
with the temperary
value $\underline{w}[k-1]$

More generally for other loss functions the

$$\underline{w}[k] = \underline{w}[k-1] + \mu(k) \sum_{n=1}^m \Psi'(d[n] - \underline{w}^T[k-1] \underline{x}_a[n]) \underline{x}_a[n]$$

Batch algorithm is written as

$$\underline{w}[k] = \underline{w}[k-1] + \mu(k) \sum_{n=1}^m \operatorname{sgn}[d[n] - \underline{w}^T[k-1] \underline{x}_a[n]] \underline{x}_a[n]$$

For example for $\Psi(\xi) = |\xi|$ $\Psi'(\xi) = \operatorname{sgn}(\xi)$ and

$$\underline{w}[k] = \underline{w}[k-1] + \mu(k) \sum_{n=1}^m \operatorname{sgn}[d[n] - \underline{w}^T[k-1] \underline{x}_a[n]] \underline{x}_a[n]$$

the algorithm is named

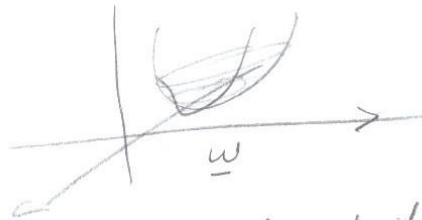
Batch-sign algorithm. Any other cost function

in the table X can be used. The question is:

Does the algorithm converge?

(LR23)

In the LMS case, we know that the cost function $E(\underline{w})$ is quadratic and therefore convex in \underline{w} .



The minimum is unique (we find it in closed-form) and the gradient search points towards the bottom of the bowl. There may be an issue at the bottom of the bowl because unless the updated value $w[k]$ is ~~such that~~ exactly

$$\sum_{n=1}^{m^2} (\underline{d}[n] - \underline{w}^T[\underline{k}] \underline{x}_a[n]) \underline{x}_a[n] = 0,$$

the search keeps wandering around the minimum.

It is therefore ~~unstable~~ there comes to a decreasing

$\mu(k)$, or ADAM, to help conclude the search.

This is usually a minor issue in linear problems.

All the considerations of using mini-batch

searches or on-line searches of section X

clearly applies here too.

Noteworthy is the on-line (instantaneous) version of the algorithm

$$w[k] = w[k-1] + \mu(k) x_a[k] (\underline{d}[k] - \underline{w}^T[k-1] \underline{x}_a[k])$$

which is the famous LMS algorithm. Note the use of the time index here because every time

an example is presented to the network, there is an update. This algorithm is usually very noisy and to get convergence we have to use a small stepsize parameter and pass the training set many times. (LR24)

A CONVERGENCE ANALYSIS

In linear systems there is a vast literature on studying the convergence of these algorithms, mostly for adaptive filters. Even if such analyses may not apply to more complex architectures, it is instructive for the student to be aware of the main ideas, also because some features may be incorporated into larger architectures. Consider the batch update with a constant stepsize parameter $\mu[K] = \mu$.

Convergence for the weight vector $w[k]$

$$w[k] = w[k-1] + \mu \sum_{n=1}^{N^2} (d[n] - w^T[k-1] x_a[n]) x_a[n]$$

that we can write as

$$w[k] = w[k-1] + \mu \sum_{n=1}^{N^2} x_a[n] d[n] - \mu \left(\sum_{n=1}^{N^2} x_a[n] x_a[n]^T \right) w[k-1]$$

$$w[k] = w[k-1] + \mu \underbrace{x_{\text{grad}}}_{=x_a} - \mu \underbrace{R_{x_a}}_{=R} w[k-1]$$

We know that the optimal solution is

$w^0 = R^{-1} x_{\text{grad}}$, therefore we consider the difference

$$w[k] - w^0 = w[k-1] - w^0 + \mu x_{\text{grad}} - \mu R x_a w[k-1]$$

Define $w[k] - w^0 = v[k]$ to be the weight error vector

(LR25)

$$V[k] = V[k-1] + \mu R_{x_a} - \mu R_{x_a} (V[k-1] + W^2)$$

$$V[k] = V[k-1] + \cancel{\mu R_{x_a}} - \mu R_{x_a} V[k-1] - \cancel{\mu R_{x_a} R_{x_a}^{-1} R_{x_a}}$$

$$V[k] = (I - \mu R_{x_a}) V[k-1]$$

The sequence $\{V[k]\}$ must converge to zero.
 This is not easily seen in this recurrence because
 all the components are coupled (R_{x_a} is not diagonal).
 Therefore we use the spectral decomposition of R_{x_a}

$$R_{x_a} = Q \Lambda Q^T$$

$Q = [q_1, q_2, \dots, q_N]$ eigenvectors

$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_N)$ eigenvalues.

Substituting

$V[k] = (I - \mu Q \Lambda Q^T) V[k-1]$
 and projecting the error vector in the eigenspace
 of R_{x_a} , $V[k] = Q^T V[k]$,

$$Q^T V[k] = Q^T (I - \mu Q \Lambda Q^T) V[k-1]$$

$$V[k] = (Q^T - \mu Q \Lambda Q^T) V[k-1]$$

$$V[k] = (I - \mu \Lambda) V[k-1]$$

Now all the components of $V[k]$ are "decoupled" (N modes)

$$\begin{cases} V_1[k] = (1 - \mu \lambda_1) V_1[k-1] \\ V_2[k] = (1 - \mu \lambda_2) V_2[k-1] \\ \vdots \\ V_N[k] = (1 - \mu \lambda_N) V_N[k-1] \end{cases}$$

The condition for convergence of all the components is that

(LR26)

$$|1 - \mu \lambda_i| < 1 \quad \forall i = 1, \dots, N$$

i.e.

$$-1 < 1 - \mu \lambda_i < 1$$

$$-2 < -\mu \lambda_i < 0$$

$$0 < \mu \lambda_i < 2$$

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

A condition that works for all is

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

this means that μ must be positive and not too large.

Each component $v^{[k]}$ converges to the optimal value with a factor $(1 - \mu \lambda_i)$. If we want to associate a convergence time to each mode, we write

$$1 - \mu \lambda_i = e^{-\frac{1}{T_i}}$$

$$T_i = -\frac{1}{\ln(1 - \mu \lambda_i)}$$

T_i is the time that $v^{[k]}$ takes to reach $\frac{1}{e}$ of $v_i^{[0]}$. If μ is small, $1 - \mu \lambda_i$ is around 1 and negative, arising $\ln x \approx x - 1$ $\ln(1 - \mu \lambda_i) \approx -\mu \lambda_i - 1$

$$T_i \approx \frac{1}{\mu \lambda_i}$$

A more detailed analysis of the evolution of $\underline{w}^{[k]}$ can be obtained observing that

$$v^{[k]} = Q v^{[k]}$$

or

$$\underline{w}^{[k]} = \underline{w}^0 + Q v^{[k]}$$

$$\underline{w}^{[k]} = \underline{w}^0 + \sum_{i=1}^N q_i v_i^{[k]}$$

Each element of $\underline{w}^{[k]}$ can be written as

$$w_i^{[k]} = w_i^0 + \sum_{j=1}^N q_{ij} v_j^{[0]} (1 - \mu \lambda_j)^k \quad i = 1, \dots, N$$

The overall convergence of $v^{[k]}$ is the combination of N modes. The slowest mode is the one corresponding to the largest eigenvalue and the fastest by the smallest.

The time to reach the solution will be between the bounds

$$-\frac{1}{\ln(1 - \mu \lambda_{\max})} < T_0 < -\frac{1}{\ln(1 - \mu \lambda_{\min})}$$

CONVERGENCE FOR THE SQUARED ERROR

In one of our previous chapters, we have shown that the cost function for the mean squared error, can be expressed in the following form

$$E(\underline{w}^{[k]}) = E(\underline{w}^0) + (\underline{w}^{[k]} - \underline{w}^0)^T \underline{R}_{x_n} (\underline{w}^{[k]} - \underline{w}^0)$$

(easy to verify expanding the second term)

Using the spectral decomposition for P_{ext} , and $\mathcal{D}[k] = \frac{W[k]}{w^0}$

$$\mathcal{E}(W[k]) = \mathcal{E}(w^0) + \underbrace{V[k]^T Q X Q^T}_{\Delta V[k]} W[k]$$

(LR28)

$$= \mathcal{E}(w^0) + V[k]^T \Delta V[k]$$

$$= \mathcal{E}(w^0) + \sum_{i=1}^N V_i[k] \lambda_i$$

$$= \mathcal{E}(w^0) + \sum_{i=1}^N \lambda_i (1 - \mu \lambda_i)^{-1} V_i[0]$$

Under the same conditions $(1 - \mu \lambda_i)^{-1}$ we have

that

$$\mathcal{E}(W[k]) \rightarrow \mathcal{E}(w^0)$$

Now since the error is the superposition of N modes, just as before, we can write that the time to reach the solution for each mode of $\mathcal{E}(W[k])$ is about

$$\tau_{i,\epsilon} \approx -\frac{1}{2 \ln(1 - \mu \lambda_i)} \approx \frac{1}{2 \mu \lambda_i}$$

Note that the time it takes for the sum to reach convergence is about a half of the time it takes the coefficients to converge.

This is a behavior that we find also in more complex neural networks.

OTHER CONVERGENCE ANALYSIS

As mentioned above, there are many contributions in the literature that analyze the convergence of more algorithms for the linear filters.

A more expanded version of these ideas goes beyond the scope of these notes but see

directed towards more complex neural network architectures.
(LRZ9)

We just would like to mention a popular analysis of the ~~instantaneous~~ stochastic algorithm.

$$w[k] = w[k-1] + \mu x_o[k] (d[k] - \hat{w}[k-1] x_o[k])$$

where we consider $w[k]$ to be a stochastic process and study the evolution of its expected value.

$$E[w[k]] = E[w[k-1]] + \mu E[x_o[k] d[k]] - \mu E[x_o[k] x_o[k]]_{w[k-1]} / E[w[k-1]]$$

$$\text{with } E[x_o[k] d[k]] = r_{x_o d} \text{ and } E[x_o[k] x_o[k]] = R_{x_o}$$

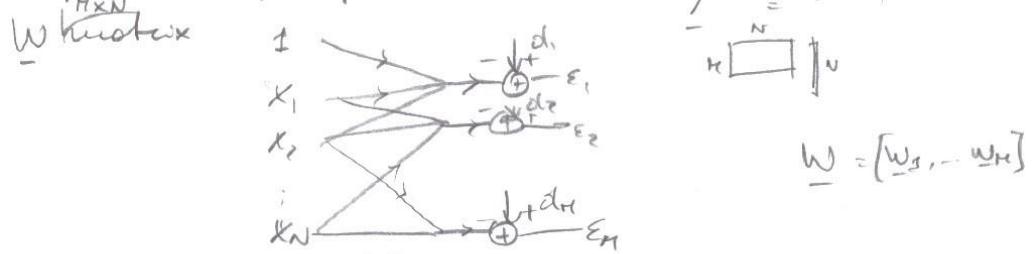
The same analysis reported above can be immediately applied to the recursion

$$E[\underline{w}[k]] = E[w[k-1]] - \mu R_{x_o d} - \mu R_{x_o} E[w[k-1]]$$

We refer the interested student to the vast literature on the topic [Huglin book] where simplifying assumptions are carefully detailed.

GRADIENT SEARCH FOR THE MULTI-OUTPUT CASE (LR3D)

The extension of the gradient search for the linear case when there are multiple outputs is quite straightforward. Now $y = \underline{W}^T \underline{x}_o$ with



The cost function $E(\underline{W})$ to minimize is

$$E(\underline{W}) = \frac{1}{M} \sum_{m=1}^M \sum_{j=1}^N \Psi'_j (d_j[m] - \underline{W}_j^T \underline{x}_o[m])^2$$

Where the losses for the M outputs are simply added. Now we can compute the gradients with respect to each one of the vectors w_j as

$$\frac{\partial}{\partial w_j} = \frac{-1}{M} \sum_{m=1}^M \Psi'_j (d_j[m] - \underline{W}_j^T \underline{x}_o[m]) \underline{x}_o[m] \quad j = 1, \dots, M$$

and use the gradient update

$$\underline{w}_j[k] = \underline{w}_j[k-1] + \mu[k] \sum_{m=1}^M \Psi'_j (d_j[m] - \underline{W}_j[k-1]^T \underline{x}_o[m]) \underline{x}_o[m]$$

$$j = 1, \dots, M$$

However it is quite appealing to use directly a matrix formulation and consider a gradient flow, i.e. the derivative with respect to a matrix. Let us see how.

The cost function is written as

(LR31)

$$E(\underline{W}) = \frac{1}{m_2} \sum_{n=1}^{m_2} L(\underline{d}[n], \underline{W}^T \underline{x}_a[n])$$

And the updates are

$$\underline{W}[k] = \underline{W}[k-1] - \mu[k] \sum_{n=1}^{m_2} \nabla_{\underline{W}} L(\underline{d}[n], \underline{W}^T \underline{x}_a[n])$$

where $\nabla_{\underline{W}}$ is the matrix with all its elements being the derivatives of L with respect to each element of \underline{W} . For example if we consider $L(\underline{d}, \underline{W}^T \underline{x}_a)$ to be the sum of the squares of the differences, i.e.

$$L(\underline{d}, \underline{W}^T \underline{x}_a) = \sum_{j=1}^{m_2} (\underline{d}_j - \underline{W}_j^T \underline{x}_a)^2$$

we can write it in compact matrix form as

$$L(\underline{d}, \underline{W}^T \underline{x}_a) = (\underline{d} - \underline{W}^T \underline{x}_a)^T (\underline{d} - \underline{W}^T \underline{x}_a)$$

which can be expanded to give

$$\begin{aligned} L(\underline{d}, \underline{W}^T \underline{x}_a) &= \underline{d}^T \underline{d} + \underline{x}_a^T \underline{W} \underline{W}^T \underline{x}_a - \underline{x}_a^T \underline{W} \underline{d} - \underline{d}^T \underline{W}^T \underline{x}_a \\ &= \underline{d}^T \underline{d} + \underline{x}_a^T \underline{W} \underline{W}^T \underline{x}_a - 2 \underline{d}^T \underline{W}^T \underline{x}_a \end{aligned}$$

Now we can use matrix differential calculus to compute the gradient with respect to \underline{W} [gradient flow (Hogins and Wundtke in Matrix cookbook)]

We need the following results

$$\nabla_X \text{tr}(AX) = A^T \quad (1)$$

(LR32)
 $\text{tr}(\cdot)$ is the trace

$$\nabla_X \text{tr}(XAX^T B) = B^T X A^T + B X A \quad (2)$$

Now, since the trace of a scalar is a scalar and $\text{tr}(AB) = \text{tr}(BA)$, we can write the loss as

$$L(d, \underline{w}^T \underline{x}_a) = d^T d + \text{tr}[\underline{x}_a^T W W^T \underline{x}_a] - 2 \text{tr}[d^T W^T \underline{x}_a]$$

$$= d^T d + \text{tr}[W W^T \underline{x}_a \underline{x}_a^T] - 2 \text{tr}[\underline{x}_a d^T W]$$

we can use (1) setting $X = W^T$ and $A = \underline{x}_a d^T$, to get

$$\nabla_{\underline{w}^T} \text{tr}[\underline{x}_a d^T W^T] = d \underline{x}_a^T \Rightarrow \nabla_{\underline{w}} = \underline{x}_a d^T$$

we can use (2) setting $\begin{cases} X = W \\ A = I \\ B = \underline{x}_a \underline{x}_a^T \end{cases}$, to get

$$\nabla_{\underline{w}} \text{tr}[W W^T \underline{x}_a \underline{x}_a^T] = \underline{x}_a \underline{x}_a^T W + \underline{x}_a \underline{x}_a^T W = 2 \underline{x}_a \underline{x}_a^T W$$

Therefore

$$\nabla_{\underline{w}} (d - \underline{w}^T \underline{x}_a)^T (d - \underline{w}^T \underline{x}_a) = 2 \underline{x}_a \underline{x}_a^T W - 2 \underline{x}_a d^T$$

the gradient of $E(W)$ is clearly null when

$$\left(\sum_{u=1}^{m_a} \underline{x}_a[u] \underline{x}_a[u]^T \right) W = \sum_{u=1}^{m_a} \underline{x}_a[u] d^T u$$

$$R_{\underline{x}_a} W = R_{\underline{x}_a d}$$

the we had already derived. the update (batch) is also immediately the obvious matrix version of the one derived for the individual w_j

$$W[k] = W[k-1] + \mu[k] \sum_{n=1}^{m^2} \left(x_a[n] d[n] - \underbrace{x_a[n] x_a^T[n] W[k-1]}_{e[n]} \right)$$

(LR33)

which can be written as

$$W[k] = W[k-1] + \mu[k] \sum_{n=1}^{m^2} x_a[n] (d[n] - W[k-1] x_a[n])^T$$

$$= W[k-1] + \mu[k] \sum_{n=1}^{m^2} x_a[n] (d[n] - \underbrace{y[n, k-1]}_{\text{out put computed with the old value } W[k-1]})^T$$

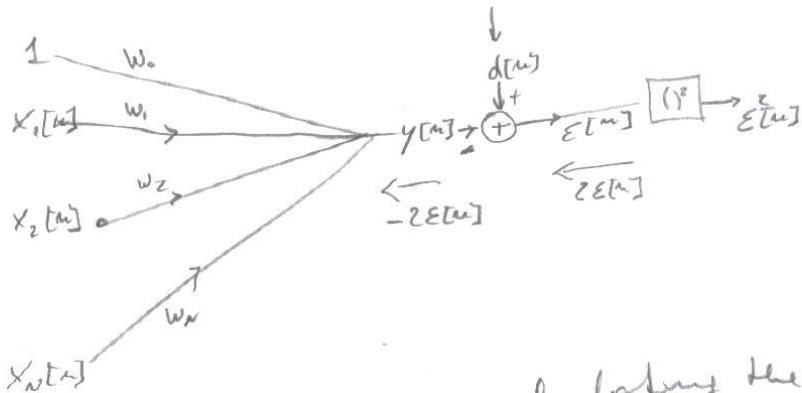
out put computed with
the old value $W[k-1]$.

$\overbrace{\text{Error vector}}$
with the old value $W[k-1]$
and the current $d[n]$,
example $(d[n], x_a[n])$

To write the recursions in matrix form may
simplify greatly the programming of the algorithm.

A FIRST LOOK AT THE GRADIENT ALGORITHM AS BACKPROPAGATION

LR34



Note that in the objective of calculating the derivative of $\epsilon^2[m]$ with respect to w_i we can use the chain rule

$$\frac{\partial \epsilon^2}{\partial w_i} = \frac{\partial \epsilon^2}{\partial \epsilon} \frac{\partial \epsilon}{\partial w_i} = \underbrace{\frac{\partial \epsilon^2}{\partial \epsilon}}_{2\epsilon} \underbrace{\frac{\partial \epsilon}{\partial y}}_{-1} \underbrace{\frac{\partial y}{\partial w_i}}_{x_i}$$

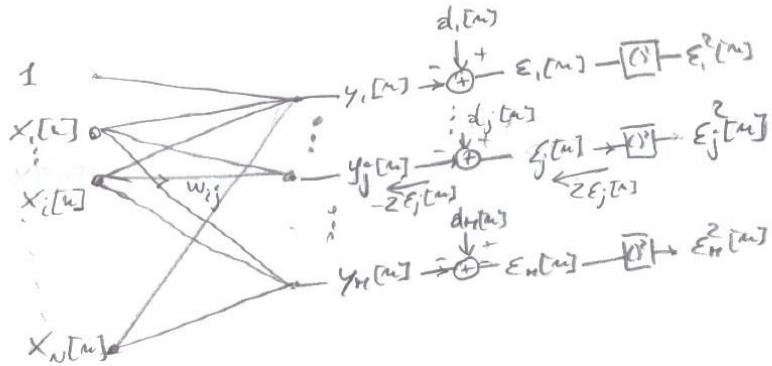
(1) At each example backpropagate $\epsilon^2[m]$ on $y[m]$

- (2) Multiplies by $x_i[m]$
- (3) Stores this product and use it to update w_i at the end of the epoch (or a subsection of it for minibatch and instantaneous versions)

Backpropagation through the block $(\cdot)^2$ is simply the multiplication by 2. Back propagation through the non-linear node is simply multiplication by -1 .

The same for multiple outputs

(LR35)



$$\frac{\partial}{\partial w_{ij}} \sum_{j=1}^n \epsilon_j^2 = \frac{\partial \epsilon_j^2}{\partial w_{ij}} = \underbrace{\frac{\partial \epsilon_j^2}{\partial \epsilon_j}}_{2\epsilon_j} \underbrace{\frac{\partial \epsilon_j}{\partial y_{ij}}}_{-1} \underbrace{\frac{\partial y_{ij}}{\partial w_{ij}}}_{x_i}$$

Note that in this architecture the gradients do not intersect. For each ϵ_j and w_{ij} the updates are independent. This will not be the case in multi-layer structures where the gradients in the previous layers will intersect.

COMPRESSING THE INPUT SPACE

(LR36)

It is quite common in the applications of signal processing that the data vector $x[m]$ presented to the network is largely redundant.

This may be the consequence that many of its components carry the same information on different scales and modes. In a linear system this ~~redundancy~~^{independence} is captured by the covariances among the various components. knowing that the means have ~~been removed~~^{we can look at} the structure of the autocorrelation matrix R_x , that is we have seen previously, plays an essential role in providing least squares solutions to regression problems.

Therefore ~~we consider~~^{we} the spectral decomposition

$$R_x = Q \Lambda Q^T$$

where Λ is a diagonal matrix containing the eigenvalues

$$\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_N)$$

and Q the unitary matrix containing on its columns the corresponding eigenvectors

$$Q = [q_1, q_2, \dots, q_N]$$

It often happens that some of the eigenvalues are very small, or even null. This means that the input space is degenerate and only a smaller number of components may be sufficient to represent the whole input dataset.

LR3Z

Assume that the eigenvalues are ranked in a descending order

$$\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_N$$

We could consider a division of the set into two parts

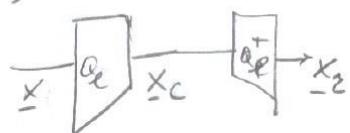
$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_l \geq \lambda_{l+1} \geq \dots \geq \lambda_N$$

where the first l eigenvalues are the most relevant and the others may be considered too small or null to be relevant. The spectral factorization can be written in a partitioned form as

$$R_{X_k} = Q \Lambda Q^T = \begin{bmatrix} Q_e & Q_{N-e} \\ Q_e & Q_{N-e} \end{bmatrix} \begin{bmatrix} \Lambda_e & 0 \\ 0 & \Lambda_{N-e} \end{bmatrix} \begin{bmatrix} Q_e^T \\ Q_{N-e}^T \end{bmatrix}$$

and consider the projection of the input vector onto the first l components on the principal vectors

$$x_c = Q_e^T x$$



It is easy to demonstrate that the reconstruction of x by the inverse Q_e^T (remember that Q is unitary)

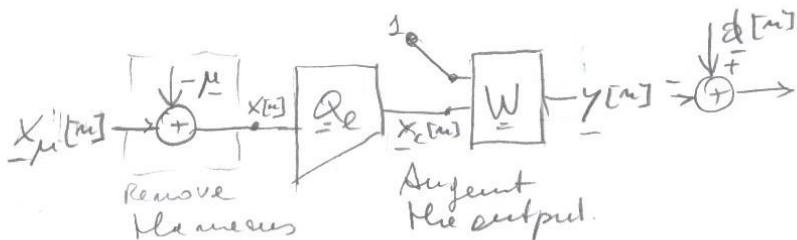
$$\underline{x}_c = \underline{Q}_c \underline{x}_c$$

$$\frac{1}{m} \sum_{n=1}^m \| \underline{x}_c[n] - \underline{x}[n] \|^2 = \sum_{i=l+1}^N \lambda_i$$

(LR38)

So that if the ignored eigenvalues are small or null, the new representation of $\underline{x}[n]$ may be sufficiently accurate. This would reduce the computational complexity of the following step.

For a linear regression problem the scheme may be the following:



Note that even if the means have been removed from the original vector, and $\underline{x}[n]$ and $\underline{x}_c[n]$ have zero means, the bias before the linear regressor may be necessary to take care of the means of d .

L.R. 39

Approaches to digital (FIR) filters

Wiener Filters

Adaptive filters

CMS

etc...