

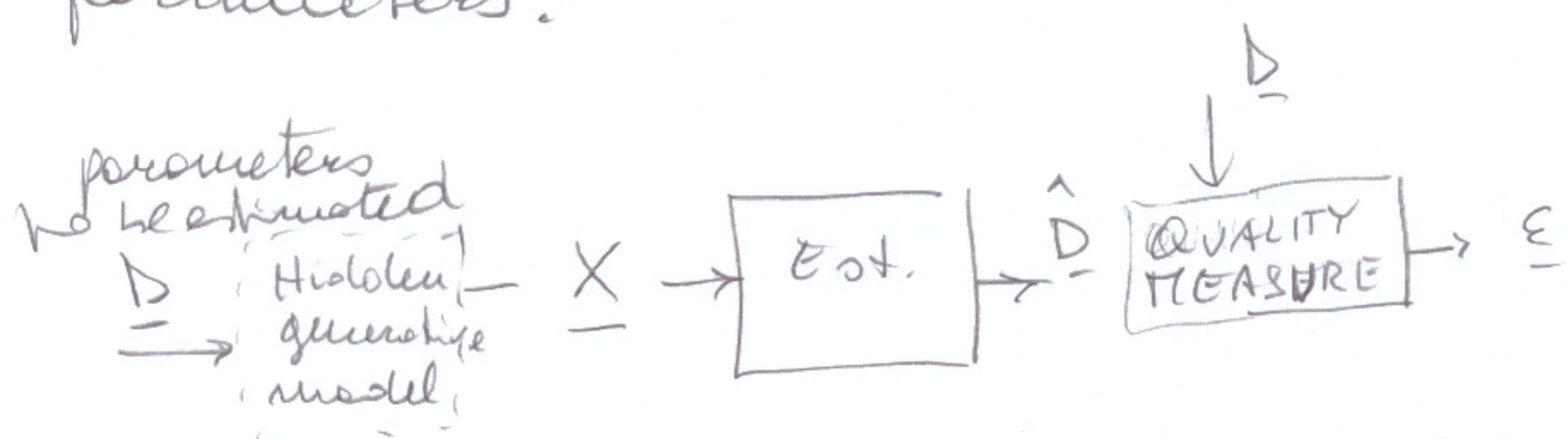
# MODEL-BASED REGRESSION

lezioni di Telecomunicazioni

Prof. FRANCESCO A.N. PACIUBRI  
(Università della Campania)

NOTE DAL CORSO DI "SIGNAL PROCESSING  
AND DATA FUSION" AA. 2022-23

In a general estimation problem, we assume that our observations are the available data to estimate unknown parameters.



A "regression" problem is posed when the parameters to be estimated are defined in the continuous space

$$D \in D \subseteq \mathbb{R}^d$$

The filtering problem belongs to this class of problems. This is to be distinguished from a "classification" problem in which  $D$  belongs to a discrete finite set.

The classification problem will be treated separately.

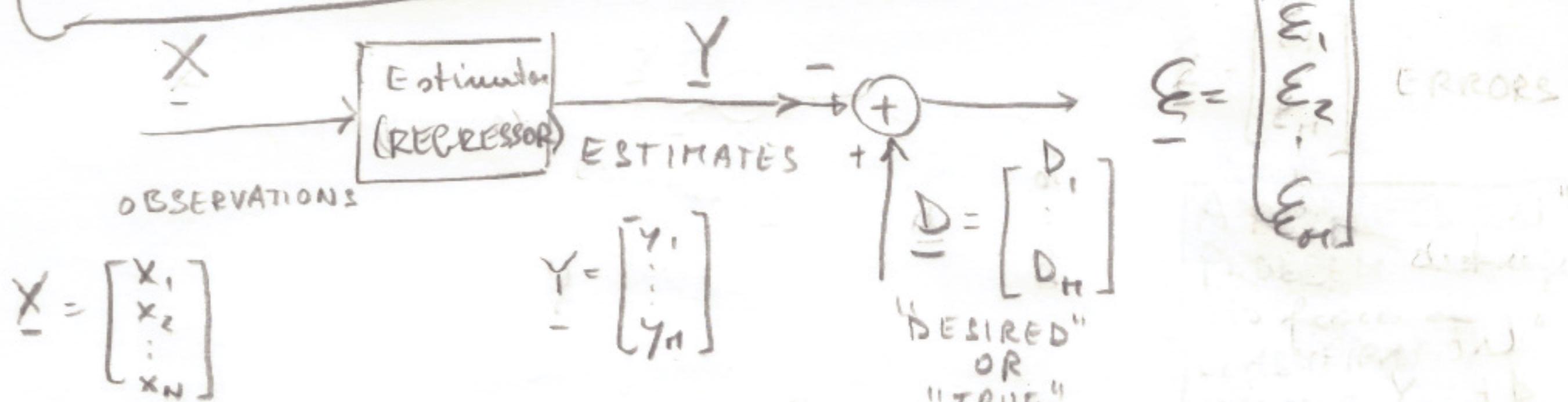
In model-based regression we rely on knowledge of the "generative model" following the scheme, to be distinguished from data-driven methods, in which we only have examples  $\{(X(n), D(n))\}$



Data can be sometimes used to estimate the model.

# THE GENERAL LINEAR MODEL

## (MODEL-BASED) LINEAR REGRESSION



$\underline{X}$  and  $\underline{D}$  are two statistically dependent multidimensional random variables with joint density  $p_{X,D}(x,d)$ . Also  $\underline{Y} \in \mathbb{Y}$

What is the function  $f(x)$  that minimizes a measure of the error  $e$ ?

Let us define a "Risk" as a function of the errors that has to measure the "borders" of the estimate.

$$R = E[\mathcal{C}(\underline{D} - \underline{Y})] = E[\mathcal{C}(\underline{D} - f(\underline{X}))] = \int \mathcal{C}(\underline{d} - f(\underline{x})) p_{X,D}(\underline{x}, \underline{d}) d\underline{x} d\underline{d} \quad (1)$$

The most typical choice for  $\mathcal{C}$  is the square of the euclidean norm.

$$\mathcal{C}_{\text{ns}}(\underline{D} - \underline{Y}) = (\underline{D} - \underline{Y})^T (\underline{D} - \underline{Y}) \quad (2)$$

If different components at the output are to be considered differently, a more general measure may be,

$$\mathcal{C}_{\text{ns}}(\underline{Y}, \underline{D}) = (\underline{D} - \underline{Y})^T Q (\underline{D} - \underline{Y}), \quad (3)$$

where  $Q$  is an  $M \times M$  positive definite matrix.

In the sequel we will assume that  $Q = I$ . All the results can be easily generalized

including Q.

Using the latter measure the risk (mean squared error) is:

$$\begin{aligned} R_{ms} &= \int_{X \times D} (\underline{d} - \underline{y})^T (\underline{d} - \underline{y}) P_{\underline{X}, \underline{D}}(\underline{x}, \underline{d}) d\underline{x} d\underline{d} \\ &= \int_X P_{\underline{X}}(\underline{x}) \int_D (\underline{d} - \underline{y})^T (\underline{d} - \underline{y}) P_{D|\underline{X}}(\underline{d}|\underline{x}) d\underline{d} d\underline{x} \quad (4) \end{aligned}$$

The result can be obtained in two steps.  
Observing that the inner integral is  
the conditional risk

$$R_{ms}(\underline{x}) = E[\mathbb{E}(\underline{d} - \underline{y}) | X = \underline{x}] = \int_D (\underline{d} - \underline{y})^T (\underline{d} - \underline{y}) P_{D|\underline{X}}(\underline{d}|\underline{x}) d\underline{d}, \quad (5)$$

which is a function of  $\underline{x}$  (not a number!).

Expanding the product we have

$$\begin{aligned} R_{ms}(\underline{x}) &= \int_D \underline{d}^T \underline{d} P_{D|\underline{X}}(\underline{d}|\underline{x}) d\underline{d} - 2 \underline{y}^T \int_D \underline{d} P_{D|\underline{X}}(\underline{d}|\underline{x}) d\underline{d} \\ &\quad + \underline{y}^T \underline{y} \int_D P_{D|\underline{X}}(\underline{d}|\underline{x}) d\underline{d} \quad (6) \end{aligned}$$

which can be re-written as (CANONICAL FORM)

$$\begin{aligned} R_{ms}(\underline{x}) &= \left( \underline{y}^T - \int_D \underline{d}^T P_{D|\underline{X}}(\underline{d}|\underline{x}) d\underline{d} \right) \left( \underline{y} - \int_D \underline{d} P_{D|\underline{X}}(\underline{d}|\underline{x}) d\underline{d} \right)^T \\ &\quad + \int_D \underline{d}^T \underline{d} P_{D|\underline{X}}(\underline{d}|\underline{x}) d\underline{d} - \left\| \int_D \underline{d} P_{D|\underline{X}}(\underline{d}|\underline{x}) d\underline{d} \right\|^2 \quad (7) \end{aligned}$$

(Verifiable by expanding the last expression)

Since the last two terms do not depend on  $\underline{R}_3$ , the only minimum is achieved when the first term goes to zero, namely when

$$\underline{y}_0 = \underline{f}_0(\underline{x}) = \int_{\underline{D}}^T \underline{d}^T \underline{P}_{\underline{D}|\underline{X}}(\underline{d}|\underline{X}) d\underline{d} = E[\underline{D}|\underline{X}=\underline{x}] . \quad (8)$$

The above expression <sup>is the conditional mean of  $\underline{D}$  given</sup> minimizes the conditional risk  $R_{ms}(\underline{x})$  which in optimal estimator is

$$R_{ms}(\underline{x}) = E[\|\underline{d}\|^2 | \underline{X}=\underline{x}] - \|\underline{y}_0\|^2 . \quad (9)$$

(Notice how this quantity still depends on  $\underline{x}$ .)

Now we want to prove that the conditional mean estimator, not only minimizes the conditional risk  $R_{ms}(\underline{x})$ , but also  $R_{ms} = E[R_{ms}(\underline{x})]$ .

Recall the following properties of the conditional expectation operator:

$$(a) E_{\underline{X}|\underline{Y}}[h(\underline{X}, \underline{Y}) | \underline{Y}=\underline{y}] = E_{\underline{X}|\underline{Y}}[h(\underline{X}, \underline{y}) | \underline{Y}=\underline{y}] \quad (10)$$

$$(b) E_{\underline{Y}}[E_{\underline{X}|\underline{Y}}[h(\underline{X}, \underline{Y}) | \underline{Y}=\underline{y}]] = E_{\underline{X}|\underline{Y}}[h(\underline{X}, \underline{Y})] \quad (11)$$

(The proof is straightforward through direct substitution of the density functions.)

Now, for the result proved above on the conditional risk, we have that

$$E_{\underline{D}|\underline{X}}[\|\underline{D}-\underline{f}_0(\underline{x})\|^2 | \underline{X}=\underline{x}] \leq E_{\underline{D}|\underline{X}}[\|\underline{D}-\underline{f}_1(\underline{x})\|^2 | \underline{X}=\underline{x}] \quad (12)$$

for any other function of  $\underline{x}$ ,  $f_1(\underline{x})$

The inequality is rewritten using (a) as

$$\underset{\underline{D}|\underline{x}}{E}\left[\|\underline{D} - f_0(\underline{x})\|^2 \mid \underline{X} = \underline{x}\right] \leq \underset{\underline{D}|\underline{x}}{E}\left[\|\underline{D} - f_1(\underline{x})\|^2 \mid \underline{X} = \underline{x}\right]. \quad (13)$$

Taking the expected value with respect to  $\underline{X}$  on both sides and using (b) we have

$$\underset{\underline{D}, \underline{X}}{E}\left[\|\underline{D} - f_0(\underline{x})\|^2\right] \leq \underset{\underline{D}, \underline{X}}{E}\left[\|\underline{D} - f_1(\underline{x})\|^2\right], \quad (14)$$

which proves the claim.

We can now say that the conditional expected value

$$f_0(\underline{x}) = \underset{\underline{D}|\underline{x}}{E}\left[\underline{D} \mid \underline{X} = \underline{x}\right] \quad (15)$$

is the estimator that minimizes the risk

$$R_{\text{ms}} = \underset{\underline{D}}{E}\left[\|\underline{D} - f(\underline{x})\|^2\right]. \quad (16)$$

The minimum value of the risk is

$$\begin{aligned} R_{\text{mso}} &= \underset{\underline{D}}{E}\left[\underline{D}^T \underline{D}\right] + \underset{\underline{x}}{E}\left[\underline{f}_0^T(\underline{x}) \underline{f}_0(\underline{x})\right] - 2 \underset{\underline{D}, \underline{x}}{E}\left[\underline{D}^T \underline{f}_0(\underline{x})\right] \\ &= \underset{\underline{D}}{E}\left[\underline{D}^T \underline{D}\right] + \underset{\underline{x}}{E}\left[\underset{\underline{D}|\underline{x}}{E}\left[\underline{D}^T \mid \underline{X}\right] \underline{f}_0(\underline{x})\right] - 2 \underset{\underline{D}, \underline{x}}{E}\left[\underline{D}^T \underline{f}_0(\underline{x})\right] \\ &= \underset{\underline{D}}{E}\left[\underline{D}^T \underline{D}\right] + \underset{\underline{x}}{E}\left[\underset{\underline{D}|\underline{x}}{E}\left[\underline{D}^T \underline{f}_0(\underline{x}) \mid \underline{X}\right]\right] - 2 \underset{\underline{D}, \underline{x}}{E}\left[\underline{D}^T \underline{f}_0(\underline{x})\right] \\ &= \underset{\underline{D}}{E}\left[\underline{D}^T \underline{D}\right] - \underset{\underline{D}}{E}\left[\underline{D}^T \underline{f}_0(\underline{x})\right] \end{aligned} \quad (17)$$

## ALTERNATIVE PROOF

R.5

The minimization of  $R_{\text{ms}}$  in (4) could be achieved by computing the gradient with respect to  $\underline{y}$  and setting it to zero.

$$\nabla_{\underline{y}} R_{\text{ms}} = \int_{\underline{x}} p_{\underline{x}}(\underline{x}) \nabla_{\underline{y}} \int_{\underline{d}} (\underline{d} - \underline{y})^T (\underline{d} - \underline{y}) P_{D|\underline{x}}(\underline{d}|\underline{x}) d\underline{d} d\underline{x}$$

Since the outer integral  $p_{\underline{x}}(\underline{x})$  is a positive function,  $\nabla_{\underline{y}} R_{\text{ms}} = 0$  needs

$$\nabla_{\underline{y}} \int_{\underline{d}} (\underline{d} - \underline{y})^T (\underline{d} - \underline{y}) P_{D|\underline{x}}(\underline{d}|\underline{x}) d\underline{d} = 0$$

Expanding and taking the gradient inside the integral we get

$$-2 \int_{\underline{D}} \underline{d}^T P_{D|\underline{x}}(\underline{d}|\underline{x}) d\underline{d} + 2 \underline{y}^T \int_{\underline{D}} P_{D|\underline{x}}(\underline{d}|\underline{x}) d\underline{d} = 0$$

and

$$\underline{y}_0 = f_0(\underline{x}) = E[D|X=\underline{x}] .$$

To prove that this is the minimum point we could take the second derivatives, in the multidimensional case, the Hessian which is <sup>(\*)</sup>

$$\mathcal{H}_{\underline{y}} = (\nabla_{\underline{y}} \nabla_{\underline{y}}^T) R_{\text{ms}}(\underline{x}) = 2 \mathbb{I}$$

where  $\underline{I}$  is the identity matrix. The Hessian is obviously positive definite, therefore the solution is a minimum. R.6

### ERROR COVARIANCE

Since the estimation is multi-dimensional it is also useful to look at the error covariance matrix, namely to

$$\sum_{\underline{\varepsilon}} = E_{\underline{X} \underline{D}} [(\underline{D} - f(\underline{X})) - E[\underline{D} - f(\underline{X})]] (\underline{D} - f(\underline{X})) - E[(\underline{D} - f(\underline{X}))]^T]$$

which in optimal conditions is

$$\sum_{\underline{\varepsilon}_0} = E[(\underline{D} - f_0(\underline{X}))(\underline{D} - f_0(\underline{X}))^T],$$

$$\text{since } E[\underline{D} - f_0(\underline{X})] = E[\underline{D}] - E_x[E[\underline{D}|\underline{X}]] = \underline{0}$$

Clearly the mean square error is

$$R_{ms} = \text{tr} \sum_{\underline{\varepsilon}_0} \quad \text{(the cost function in (3))},$$

which in optimal conditions is

$$R_{mso} = \text{tr} \sum_{\underline{\varepsilon}_0} \quad \bullet$$

PROBLEM: Formulate all the results in the problem regarding the cost function in (2)

(\*) Note that the gradient operator  $\nabla_{\underline{x}} z$  applied to the scalar

$z$  is  $\left[ \frac{\partial z}{\partial x_1} \frac{\partial z}{\partial x_2} \dots \frac{\partial z}{\partial x_N} \right]^T$ . The Hessian ~~matrix~~ is obtained

$$\text{applying the operator } \nabla_{\underline{x}} \nabla_{\underline{x}}^T = \begin{bmatrix} \frac{\partial}{\partial x_1} \\ \vdots \\ \frac{\partial}{\partial x_N} \end{bmatrix} \left[ \frac{\partial}{\partial x_1} \dots \frac{\partial}{\partial x_N} \right]^T = \begin{bmatrix} \frac{\partial^2}{\partial x_1^2} & \dots & \frac{\partial^2}{\partial x_1 \partial x_N} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial x_N \partial x_1} & \dots & \frac{\partial^2}{\partial x_N^2} \end{bmatrix}$$

to the scalar  $z$ .

## EXAMPLE 1: $\underline{X}$ and $\underline{D}$ jointly gaussian

R.7

If the observation  $\underline{X}$  and the "desired" output  $\underline{D}$  are jointly gaussian with

$$E[\underline{X}] = \bar{x}, E[\underline{D}] = \bar{d}, \text{cov}[\underline{X}] = \Sigma_x, \text{cov}[\underline{D}] = \Sigma_d$$

$E[(\underline{X} - \bar{x})(\underline{D} - \bar{d})^T] = \Sigma_{xD}$ , we have, using the results about gaussian variables that the minimum mean square error filter is given by

$$f_o(x) = E[\underline{D} | \underline{X} = x] = \bar{d} + \sum_{Dx} \Sigma_x^{-1} (x - \bar{x})$$

The error covariance matrix is

$$\Sigma_{eo} = E[(\underline{D} - f_o(\underline{X}))(\underline{D} - f_o(\underline{X}))^T] = \Sigma_d - \sum_{Dx} \Sigma_x^{-1} \Sigma_{xD}$$

with the only minimum mean square error given by

$$R_{mse} = \text{tr} [\Sigma_d - \sum_{Dx} \Sigma_x^{-1} \Sigma_{xD}]$$

Note how the optimum estimator  $f_o(x)$  is an affine estimator (linear + constant).

If both  $\underline{X}$  and  $\underline{D}$  have zero mean, the estimator becomes purely linear. I would like to emphasize that even if one of the two variables,  $\underline{X}$  or  $\underline{D}$  has non-zero mean, the estimator remains affine.

the estimator can be re-written as

$$f_0(\underline{x}) = \underbrace{\sum_{Dx} \sum_x^{-1} \underline{x}}_{\text{Linear part}} + \underbrace{\left( \bar{d} - \sum_{Dx} \sum_x^{-1} \bar{x} \right)}_{\text{biases.}} \quad (*)$$

to emphasize the linear part and the biases.

Note that in a practical estimation problem the means are usually removed beforehand if necessary. However, the ~~means are removed~~ variables are ~~remained~~ ~~remaining~~ ~~variables~~ ~~are remaining~~ ~~are removed~~.

let us consider here in some detail some special cases of the above solution.

~~EXAMPLE 1.1~~

$$N=1, L=1, E[\underline{X}] = 0, E[D] = 0.$$

The solution is

$$f_0(\underline{x}) = \frac{E[DX]}{E[X^2]} \underline{x} = \underline{c}_0 \underline{x}, \quad c_0 = \frac{E[DX]}{E[X^2]}$$

$$R_{mso} = E[D^2] - \frac{E[DX]^2}{E[X^2]}$$

EXAMPLE 1.2

$$N > 1, L = 1, E[\underline{X}] = 0, E[D] = 0$$

$$f_0(\underline{x}) = E[DX^T] E[\underline{X}\underline{X}^T]^{-1} \underline{x} = \underline{c}_0^T \underline{x} \quad (\text{linear})$$

$$\text{with } \underline{c}_0 = E[\underline{X}\underline{X}^T]^{-1} E[DX]$$

$$R_{mso} = E[D^2] - \underbrace{E[DX^T] E[\underline{X}\underline{X}^T]^{-1} E[DX]}_{\substack{\text{CORR. MATRIX} \\ \text{CROSS-CORR.} \\ \text{VECTOR}}}$$

$$R_{mso} = E[D^2] - E[DX^T] E[\underline{X}\underline{X}^T]^{-1} E[DX]$$

EXAMPLE 1.3

$$N \geq 1, L \geq 1 \quad E[\underline{X}] = \underline{0}, E[\underline{D}] = \underline{0}$$

$$\underline{f}_o(\underline{x}) = E[\underline{D}\underline{X}^T]E[\underline{X}\underline{X}^T]^{-1}\underline{x} = \underline{\underline{c}}_o^T \underline{x} \quad \begin{matrix} N \\ L \end{matrix} \times \begin{matrix} N \\ N \end{matrix}$$

$$\underline{\underline{c}}_o = E[\underline{X}\underline{X}^T]^{-1}E[\underline{X}\underline{D}^T] \quad \begin{matrix} N \\ N \end{matrix}$$

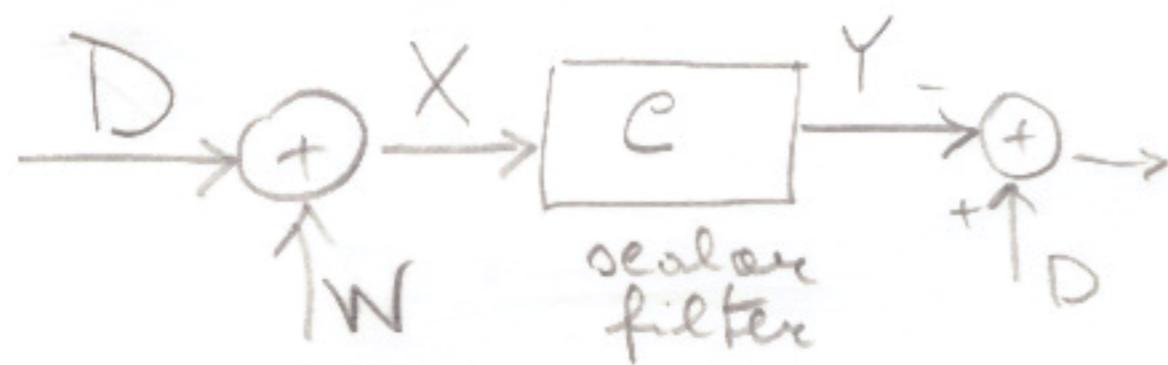
$$E[(\underline{D} - \underline{f}_o(\underline{x}))(\underline{D} - \underline{f}_o(\underline{x}))^T] = E[\underline{D}\underline{D}^T] - E[\underline{D}\underline{X}^T]E[\underline{X}\underline{X}^T]^{-1}E[\underline{X}\underline{D}^T]$$

EXAMPLE 1.4

Suppose the observation  $\underline{X}$  is a scalar and is the result of superposition of two <sup>independent</sup> Gaussian zero mean random variable  $D$  and  $W$ .

We want to estimate  $D$  and assume that  $W$  is just noise.

The statistics are  $E[D] = 0, E[W] = 0,$   
 $E[D^2] = \sigma_D^2, E[W^2] = \sigma_W^2, E[DW] = 0$ .



This case follows in the category of Example 1.1 since  $X$  and  $D$  are jointly Gaussian, with

$$E[\underline{D}\underline{X}] = E[\underline{D}(\underline{D} + \underline{W})] = E[\underline{D}^2]$$

$$E[\underline{X}^2] = E[(\underline{D} + \underline{W})^2] = E[\underline{D}^2] + E[\underline{W}^2].$$

The optimum estimator is linear, scalar and has coefficient

$$c_o = \frac{E[\underline{D}^2]}{E[\underline{D}^2] + E[\underline{W}^2]} = \frac{\sigma_D^2}{\sigma_D^2 + \sigma_W^2}$$

and minimum error

$$\begin{aligned} R_{\text{MSO}} &= E[\underline{D}^2] - \frac{E[\underline{D}^2]^2}{E[\underline{D}^2] + E[\underline{W}^2]} \\ &= \frac{E[\underline{D}^2] + E[\underline{W}^2]E[\underline{D}^2] - E[\underline{D}^2]^2}{E[\underline{D}^2] + E[\underline{W}^2]} \\ &= \frac{E[\underline{W}^2]E[\underline{D}^2]}{E[\underline{D}^2] + E[\underline{W}^2]} \end{aligned}$$

### EXAMPLE 1.5

Consider the multi-dimensional case of EXAMPLE 1.4, namely  $\underline{X} = \underline{D} + \underline{W}$ ,

where  $\underline{D}$  and  $\underline{W}$  are two  $N$ -dimensional independent random variables with zero means. The statistics are  $E[\underline{D}] = E[\underline{W}] = \underline{0}$ ,  $E[\underline{W}] = \sigma_w^2 \underline{I}$ ,  $E[\underline{D}\underline{D}^T] = \underline{R}_D$ ,  $E[\underline{D}\underline{W}] = \underline{0}$ .

$$E[\underline{D}\underline{X}^T] = E[\underline{D}(\underline{D}^T + \underline{W}^T)] = E[\underline{D}\underline{D}^T]$$

$$\begin{aligned} E[\underline{X}\underline{X}^T] &= E[(\underline{D} + \underline{W})(\underline{D}^T + \underline{W}^T)] = E[\underline{D}\underline{D}^T] + E[\underline{W}\underline{W}^T] \\ &= \underline{R}_D + \sigma_w^2 \underline{I} \end{aligned}$$

The optimum estimator for the whole vector  $\underline{D}$  is

$$\underline{f}_0(\underline{x}) = \underline{R}_D \underbrace{\left( \underline{R}_D + \sigma_w^2 \underline{I} \right)^{-1}}_{C_0} \underline{x} = \underline{C}_0^T \underline{x}$$

$$\underline{C}_0 = \underline{R}_D \left( \underline{R}_D + \sigma_w^2 \underline{I} \right)^{-1}$$

The error covariance in optimal condition is

$$\sum_{\underline{\varepsilon}_0} = \underline{R}_D - \underline{R}_D \left( \underline{R}_D + \sigma_w^2 \underline{I} \right)^{-1} \underline{R}_D^T$$

In the special case of  $E[\underline{D}\underline{D}^T] = \sigma_D^2 \underline{I}$ , we have

$$\underline{c}_0 = \frac{\sigma_D^2}{\sigma_D^2 + \sigma_w^2} \underline{I},$$

which is the same as the case in example 1.6. It is not necessary in this case to use a matrix filter, but a scalar estimator on all the components would suffice.

In the general hypothesis it is often desired to estimate just one element of  $\underline{D}$  from  $\underline{x}$ , say  $D_j$  ( $j \in \{1, \dots, N\}$ ).

The cross-correlations necessary are only those in

$$E[D_j \underline{x}^T] = E[D_j (\underline{D}^T + \underline{W}^T)] = E[D_j \underline{D}^T]$$

$$\hat{d}_j = f_0(\underline{x}) = E[D_j \underline{D}^T] (\underline{R}_D + \sigma_w^2 \underline{I})^{-1} \underline{x} = \underline{c}_0^T \underline{x}$$

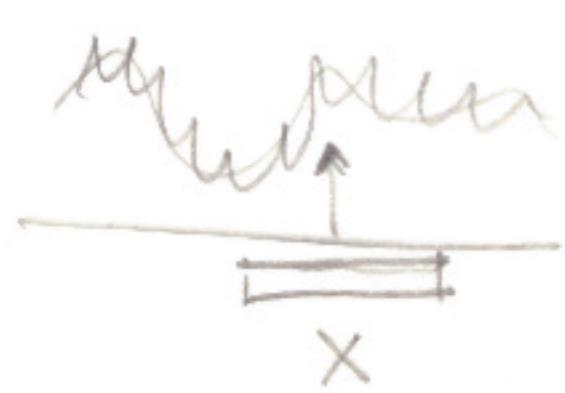
$$\underline{c}_0 = (\underline{R}_D + \sigma_w^2 \underline{I})^{-1} E[\underline{D}; \underline{D}]$$

The minimum mean square error is

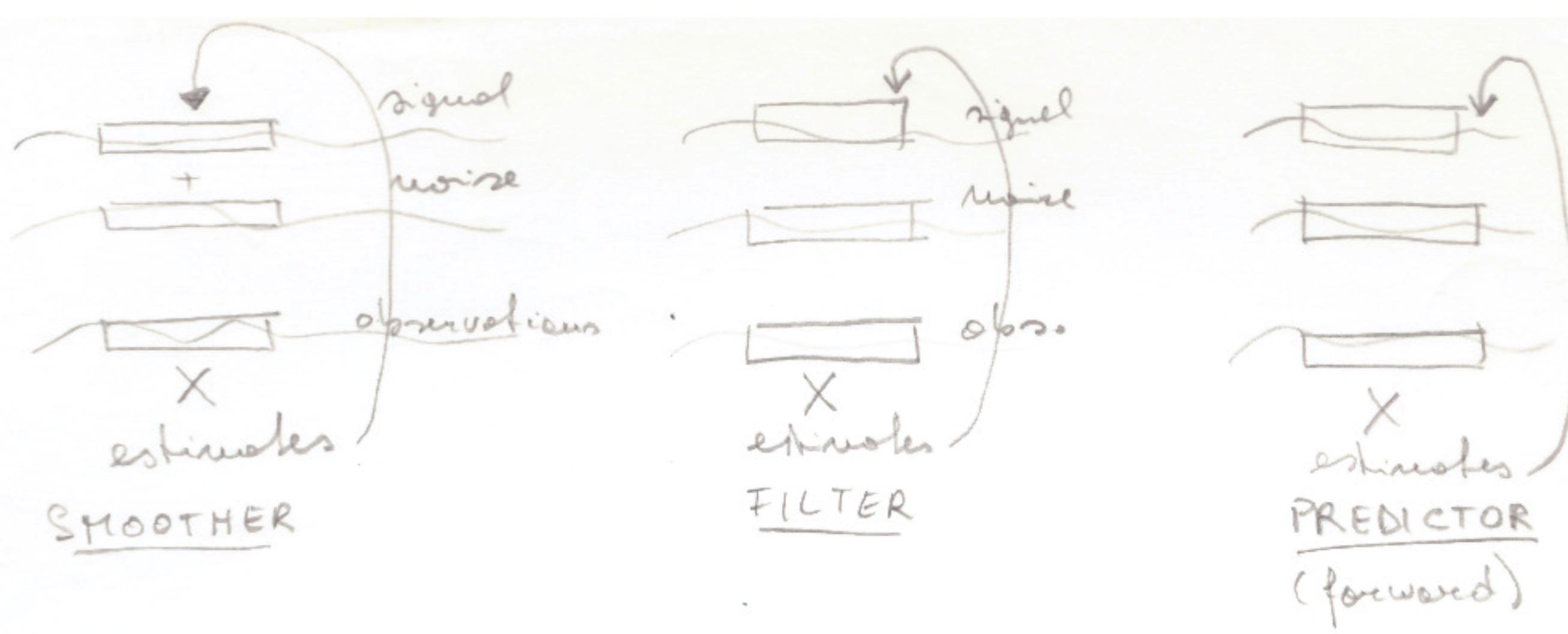
$$E[(d_j - f_0(\underline{x}))^2] = E[D_j^2] - E[D_j \underline{D}^T] (\underline{R}_D + \sigma_w^2 \underline{I})^{-1} E[\underline{D}; \underline{D}]$$

This last case is common in signal processing when we have a stationary Gaussian signal buried in Gaussian noise. A filter estimates the

center value of the signal from our observation window  $\underline{x}$



The usual nomenclature is the following one:



### EXAMPLE 1.6

Suppose that the observation vector  $\underline{X}$  is the result of the superposition of a desire (the desired variable (gaussian) and an independent gaussian error vector.

$$\underline{X} = D \begin{bmatrix} \underline{e} \\ \underline{w} \end{bmatrix} + \underline{W} = D\underline{e} + \underline{W}$$

We want to estimate  $D$ , this is a so-called location estimation problem.

$$E[D\underline{X}^T] = E[D(D\underline{e} + \underline{W})^T] = E[D^T] \underline{e}^T + \underline{\sigma}_w^2 I$$

$$E[\underline{X}\underline{X}^T] = E[(D\underline{e} + \underline{W})(D\underline{e} + \underline{W})^T] = E[D^2] \underline{e}\underline{e}^T + \underline{\sigma}_w^2 I$$

$$f_o(\underline{x}) = E[D^2] \underline{e}^T (E[D^2] \underline{e}\underline{e}^T + \underline{\sigma}_w^2 I)^{-1} \underline{x}$$

Use the matrix inversion lemma (reviewed in one of the appendices) to get a more explicit result.

Matrix  $f_o(\underline{x})$  is called a smoothing filter.

$$\left( \frac{\sigma_d^2}{\sigma_w^2} I + \frac{e}{c} \frac{\sigma_d^2}{c^T} e^T \right)^{-1} = \sigma_w^{-2} I - \sigma_w^{-2} I e \left( \frac{1}{\sigma_d^2} + \frac{e^T \sigma_w^{-2} I e}{\sigma_d^2} \right)^{-1} e^T \sigma_w^{-2} I$$

R.13

$$= \sigma_w^{-2} I - \sigma_w^{-4} e \left( \frac{1}{\sigma_d^2} + \frac{N}{\sigma_w^2} \right)^{-1} e^T$$

$$f_0(\underline{x}) = \left( \sigma_d^2 e^T \sigma_w^{-2} - \sigma_d^2 \sigma_w^{-4} N \left( \frac{1}{\sigma_d^2} + \frac{N}{\sigma_w^2} \right)^{-1} e^T \right) \underline{x}$$

$$= \left( \frac{\sigma_d^2}{\sigma_w^2} - \frac{\sigma_d^2}{\sigma_w^2} N \frac{\cancel{\sigma_d^2 \sigma_w}}{\sigma_w^2 + N \sigma_d^2} \right) e^T \underline{x}$$

$$= \frac{\sigma_d^2}{\sigma_w^2} \left( 1 - \frac{N \sigma_d^2}{\sigma_w^2 + N \sigma_d^2} \right) e^T \underline{x}$$

$$= \frac{\sigma_d^2}{\sigma_w^2} \frac{\cancel{\sigma_w^2 + N \sigma_d^2 - N \sigma_d^2}}{\sigma_w^2 + N \sigma_d^2} e^T \underline{x}$$

$$= \frac{\sigma_d^2}{\sigma_w^2 + N \sigma_d^2} e^T \underline{x}$$

This is a modified sample mean, to the sample mean.

Rewriting  $f_0(\underline{x}) = \frac{1}{N + \frac{\sigma_w^2}{\sigma_d^2}}$ , for large  $N$  it tends to the sample mean.  
As we will verify in the future this is also the MAP estimator as expected since symmetric of the density function.

The minimum risk is

$$R_0 = E[D^2] - E[D f_0(\underline{x})] = \sigma_d^2 - E[D \frac{\sigma_d^2}{\sigma_w^2 + N \sigma_d^2} e^T (D e + w)]$$

$$= \sigma_d^2 - \frac{\sigma_d^2}{\sigma_w^2 + N \sigma_d^2} N \sigma_d^2 = \frac{\sigma_d^2 \sigma_w^2 + N \sigma_d^4 - \sigma_d^4 N}{\sigma_w^2 + N \sigma_d^2}$$

$$= \frac{\sigma_d^2 \sigma_w^2}{\sigma_w^2 + N \sigma_d^2}$$

It is in general difficult to derive the exact R.14 expression for the conditional expectation estimator for a generic situation.

Let us look first at a simple problem.

### EXAMPLE 2 (A Non-gaussian example)

An observation  $X = D + W$  with  $D$  and  $W$  independent and both distributed exponentially. Namely

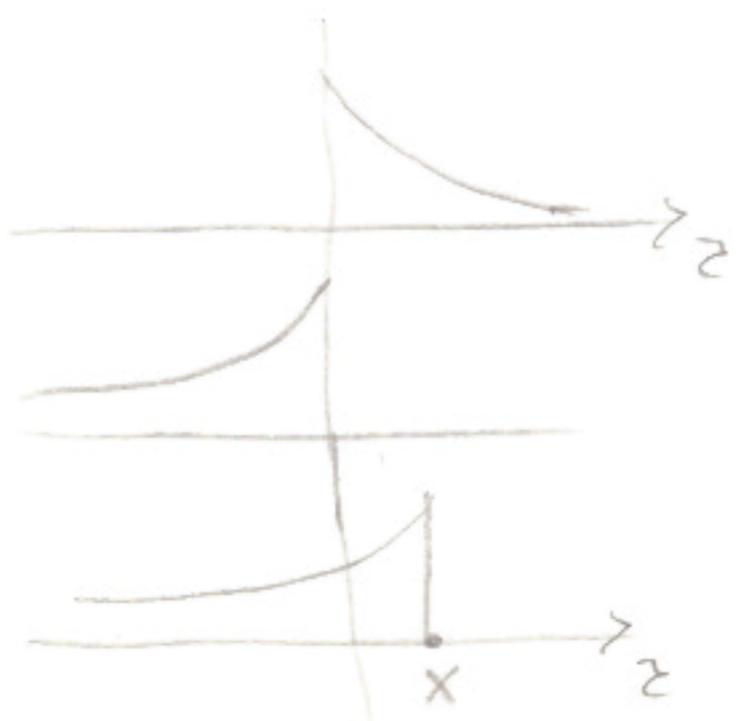
$$p_D(d) = \begin{cases} \lambda_D e^{-\lambda_D d} & d \geq 0 \\ 0 & \text{else} \end{cases} \quad p_W(w) = \begin{cases} \lambda_W e^{-\lambda_W w} & w \geq 0 \\ 0 & \text{else} \end{cases}$$

To compute the conditional expectation  $E[D|X=x]$  we need to compute the conditional density  $p_{D|X}(d|x)$ .

$$p_{D|X}(d|x) = \frac{p_{X|D}(x|d) p_D(d)}{p_X(x)}$$

Let us compute  $p_X(x)$  first:

$$p_X(x) = p_D(x) * p_W(x) = \int_{-\infty}^{+\infty} \lambda_D \lambda_W e^{-\lambda_W z} u(z) e^{-\lambda_D (x-z)} u(x-z) dz$$



$$p_X(x) = 0 \quad x < 0$$

$$p_X(x) = \int_0^x \lambda_D \lambda_W e^{-\lambda_W z} e^{-\lambda_D (x-z)} dz \quad x \geq 0$$

$$= \lambda_D \lambda_W e^{-\lambda_D x} \int_0^x e^{-(\lambda_W - \lambda_D)z} dz$$

$$= \lambda_D \lambda_W e^{-\lambda_D x} \left[ \frac{e^{-(\lambda_W - \lambda_D)x}}{-(\lambda_W - \lambda_D)} \right]^x = \frac{\lambda_D \lambda_W}{\lambda_W - \lambda_D} e^{-\lambda_D x} \left( 1 - e^{-(\lambda_W - \lambda_D)x} \right)$$

$$P_X(x) = \frac{\lambda_D \lambda_W}{\lambda_W - \lambda_D} \left( e^{-\lambda_D x} - e^{-\lambda_W x} \right) \quad x \geq 0$$

$$P_{X|D}(x|d) = P_W(x-d)$$

$$P_{D|x}(D|x) = \frac{\lambda_W e^{-\lambda_W(x-d)} u(x-d) \lambda_D e^{-\lambda_D d} u(d)}{\frac{\lambda_D \lambda_W}{\lambda_W + \lambda_D} (e^{-\lambda_D x} - e^{-\lambda_W x}) u(x)}$$

$$= (\lambda_W - \lambda_D) \frac{e^{-\lambda_W x} e^{-(\lambda_D - \lambda_W)d} u(d) u(x-d)}{(e^{-\lambda_D x} - e^{-\lambda_W x}) u(x)}$$

$$E[D|x] = \frac{(\lambda_W - \lambda_D)}{(e^{-\lambda_D x} - e^{-\lambda_W x}) u(x)} \int_0^\infty u(z) u(x-z) e^{-(\lambda_D - \lambda_W)z} z dz$$

$$= \int_0^x z e^{-(\lambda_D - \lambda_W)z} dz = \left[ \frac{e^{-(\lambda_D - \lambda_W)z}}{-(\lambda_D - \lambda_W)} \right]_0^x - \int_0^x \frac{e^{-(\lambda_D - \lambda_W)z}}{-(\lambda_D - \lambda_W)} dz$$

$$= \frac{e^{-(\lambda_D - \lambda_W)x}}{(\lambda_D - \lambda_W)} - \frac{(e^{-(\lambda_D - \lambda_W)x} - 1)}{(\lambda_D - \lambda_W)^2}$$

Since  $x$  is always  $> 0$ , the estimator is

$$f_0(x) = \frac{e^{-\lambda_D x}}{e^{-\lambda_D x} - e^{-\lambda_W x}} - \frac{(e^{-\lambda_D x} - 1)}{(\lambda_D - \lambda_W)(e^{-\lambda_D x} - e^{-\lambda_W x})}$$

This is  
non linear  
and memoryless.  
The minimum  
risk is easily  
computed

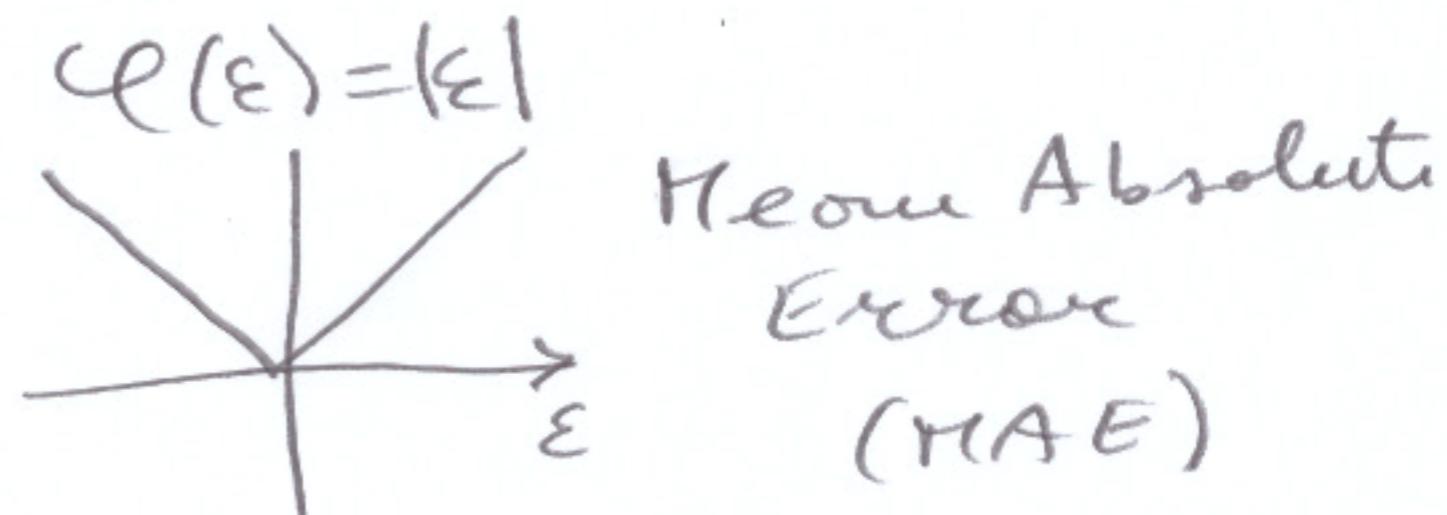
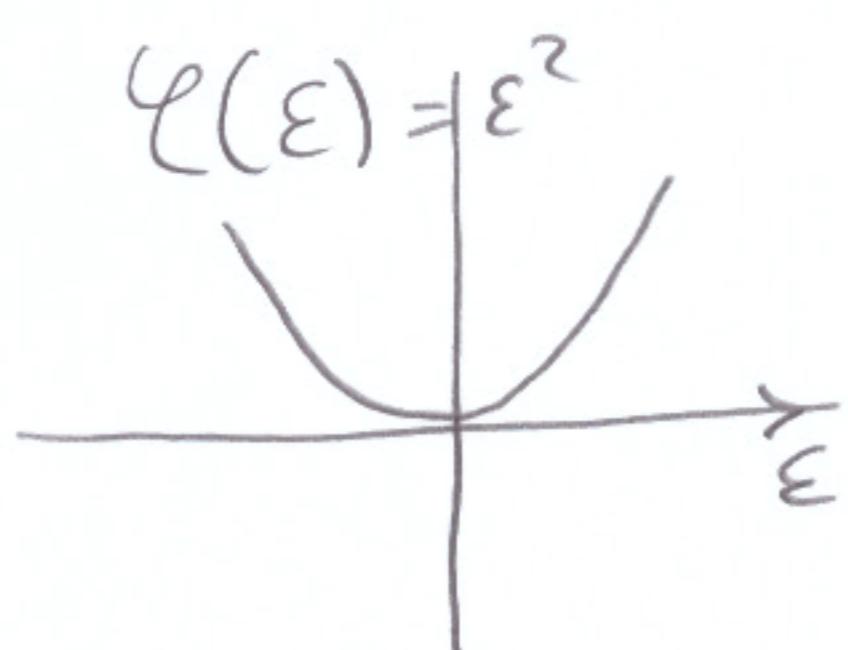
## OTHER COST FUNCTIONS

Minimum Mean Squared Error (MMSE)  
estimators are based on the cost function

$$E_{\text{MS}} = E[\varepsilon^2]$$

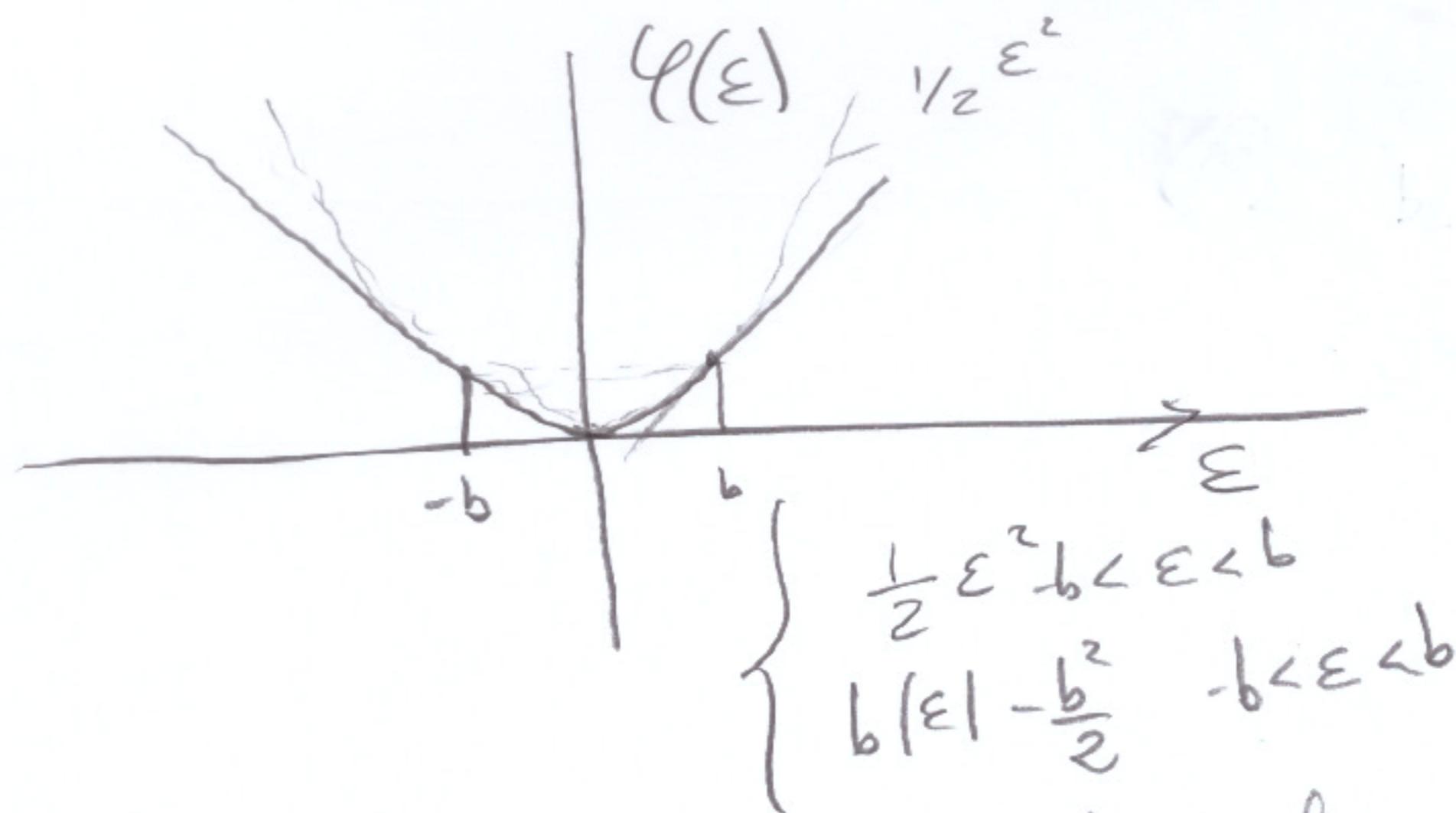
It is possible to use different measures  
of deviations (Loss Functions) other than the quadratic.

$$\mathcal{E} = E[\ell(\varepsilon)]$$

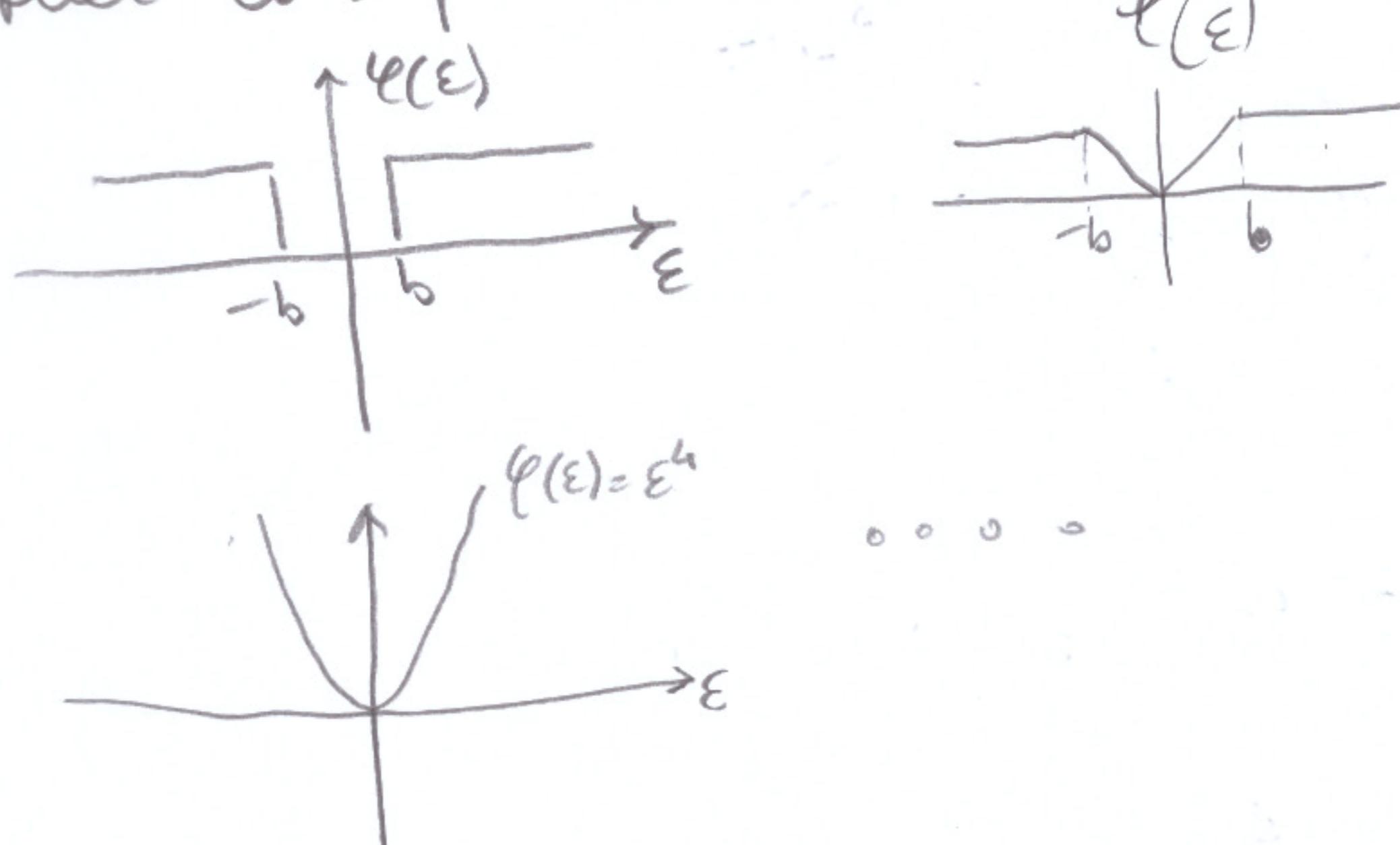


MAE sometimes is preferable to MMSE because it grows linearly versus the quadratic function. The quadratic function weights more strongly large error values. Unfortunately the function  $|\varepsilon|$  is discontinuous in its derivative in  $\varepsilon=0$  and does not lead to analytical solutions. In any case, sometimes it is used in search algorithms.

To combine the benefits of MAE and MMSE in estimation problems sometimes we use Huber loss composed by a quadratic loss for small values and an absolute deviation for large values

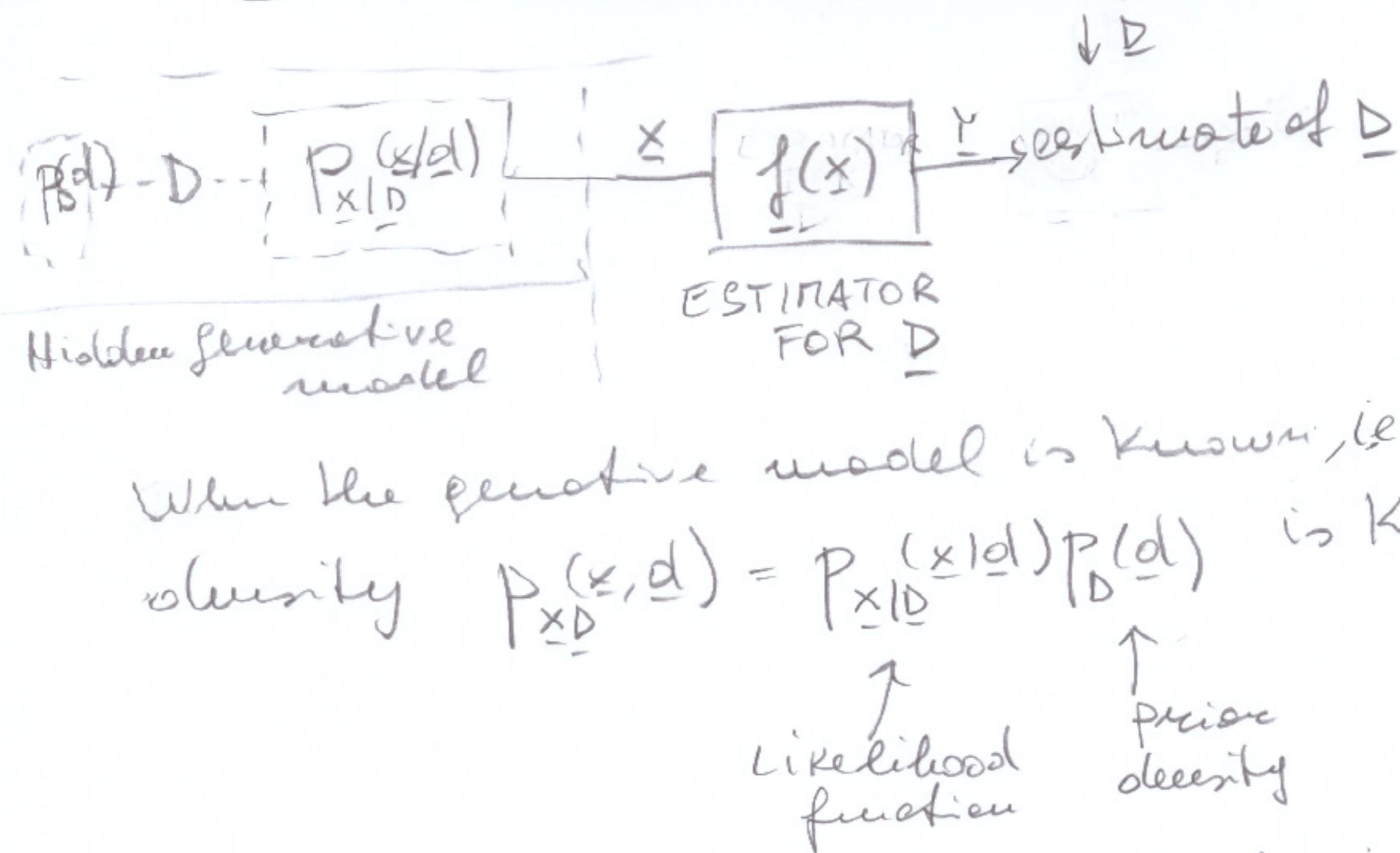


Note the continuity of the derivative for  $\varepsilon = \pm b$ .  
Other loss functions are



# MAP AND ML ESTIMATORS

Rold



An estimator for  $\underline{D}$  can be found by taking the value of  $\underline{D}$  which maximizes the posterior probability

$$\hat{\underline{D}} = f(\underline{x}) = \underset{\underline{d}}{\operatorname{argmax}} P_{\underline{D}|\underline{X}}^{(\underline{d}|\underline{x})}$$

In other words, given the observation  $\underline{x}$ , we

take the most probable  $\underline{D}$ :

MAXIMUM A POSTERIORI ESTIMATE (MAP)

Note that using Bayes theorem, the posterior can be written as

$$P_{\underline{D}|\underline{X}}(\underline{d}|\underline{x}) = \frac{P_{\underline{X}|\underline{D}}(\underline{x}|\underline{d})P_{\underline{D}}(\underline{d})}{P_{\underline{X}}(\underline{x})}$$

and the MAP estimator