

# S.I.D.A. review

F. A. Ramponi, 2013/5/27.

## 1 Least squares

### 1.1 Introduction and motivation

We know, suspect, or assume a relation between a variable  $x \in \mathbb{R}^m$  and a variable  $y \in \mathbb{R}$  in the form of a function that further depends on an unknown parameter  $\theta$ :  $y = f(x, \theta)$ . *Interpolation problem*: given a sequence of pairs  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$  that satisfy such a functional relation, find the particular  $\theta^o$  that corresponds to the function. In general the problem does not admit an exact solution. The issue is that the functional relation is an *assumption*, but measures are always corrupted by measurement errors  $\varepsilon_i$ :

$$y_i = f(x_i, \theta) + \varepsilon_i, \quad i = 1, \dots, N.$$

Hence, we search an *approximate* function  $f(\cdot, \hat{\theta})$ , instead of the “true” one  $f(\cdot, \theta^o)$ . There are at least three reasons why such an approximation is useful: 1) to investigate some internal properties of the mechanism that generates the data; 2) to measure the parameter  $\theta^o$ , which may correspond to some physical quantity of interest; 3) to *predict* a future value  $y_{N+1}$ , when a future value  $x_{N+1}$  will be available. In approximating  $f(\cdot, \theta^o)$  with a certain  $f(\cdot, \theta)$ , we commit approximation errors  $\epsilon_i(\theta)$  called *residuals*:

$$\epsilon_i(\theta) := y_i - f(x_i, \theta), \quad i = 1, \dots, N.$$

Least squares method: given  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ , find the estimate

$$\hat{\theta}_{\text{LS}} := \arg \min_{\theta} \sum_{i=1}^N \epsilon_i(\theta)^2 = \arg \min_{\theta} \sum_{i=1}^N (y_i - f(x_i, \theta))^2.$$

If  $f$  is *linear in the parameter*  $\theta$ , then the least squares problem has an analytic solution.

### 1.2 Linearity in the parameter $\theta$

Let  $x_i = (x_i^1, \dots, x_i^m) \in \mathbb{R}^m$  and  $y_i \in \mathbb{R}$  for  $i = 1, \dots, N$ . If  $y_i$  is a function of  $x_i$ , we call  $x_i$  “explanatory variable”, and we say that  $y_i$  are “explained” by  $x_i$ .  $y_i$  may not be linked to  $x_i$  directly, but through  $p$  functions  $\varphi_1 : \mathbb{R}^m \rightarrow \mathbb{R}, \dots, \varphi_p : \mathbb{R}^m \rightarrow \mathbb{R}$  called *regressors*, which can be nonlinear. The data generation model to which we will apply least squares is the following:

$$y_1 = \theta_1 \varphi_1(x_1^1, \dots, x_1^m) + \theta_2 \varphi_2(x_1^1, \dots, x_1^m) + \dots + \theta_p \varphi_p(x_1^1, \dots, x_1^m) + \varepsilon_1;$$

⋮

$$y_N = \theta_1 \varphi_1(x_N^1, \dots, x_N^m) + \theta_2 \varphi_2(x_N^1, \dots, x_N^m) + \dots + \theta_p \varphi_p(x_N^1, \dots, x_N^m) + \varepsilon_N.$$

Defining the vectors  $\theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_p \end{bmatrix}$  and  $\varphi(x_i^1, \dots, x_i^m) = \begin{bmatrix} \varphi_1(x_i^1, \dots, x_i^m) \\ \vdots \\ \varphi_p(x_i^1, \dots, x_i^m) \end{bmatrix}$  and recalling that  $x_i = [x_i^1 \dots x_i^m]^\top$ ,

we write the model in compact form:  $y_i = f(x_i, \theta) + \varepsilon_i = \varphi(x_i)^\top \theta + \varepsilon_i$  for  $i = 1, \dots, N$ . Note that the function  $f$  is linear in the parameter  $\theta$ , although not necessarily in the explanatory variable  $x$ . We shorten notation even more letting  $\varphi_i := \varphi(x_i)$ , and the model finally reads  $y_i = \varphi_i^\top \theta + \varepsilon_i$ ,  $i = 1, \dots, N$ . The least squares method asks to find

$$\hat{\theta}_{\text{LS}} := \arg \min_{\theta \in \mathbb{R}^p} \sum_{i=1}^N (y_i - \varphi_i^\top \theta)^2.$$

To find the minimum, we set its derivative w.r.t.  $\theta$  equal to zero:  $\frac{\partial}{\partial \theta} \sum_{i=1}^N (y_i - \varphi_i^\top \theta)^2 = -2 \sum_{i=1}^N \varphi_i (y_i - \varphi_i^\top \theta) = 0$ . We obtain the *normal equations*:

$$\left( \sum_{i=1}^N \varphi_i \varphi_i^\top \right) \theta = \sum_{i=1}^N \varphi_i y_i. \tag{1}$$

Any  $\theta$  that solves (1) is a minimum point for the sum of squares and a solution to the least squares problem; in particular, if  $R = \sum_{i=1}^N \varphi_i \varphi_i^\top$  is invertible (in practical applications this is generally the case), then the only solution to the normal equations reads explicitly:

$$\hat{\theta}_{\text{LS}} = \left( \sum_{i=1}^N \varphi_i \varphi_i^\top \right)^{-1} \sum_{i=1}^N \varphi_i y_i. \quad (2)$$

### 1.3 Existence and uniqueness of a solution of the normal equations

**Lemma 1.3.1**

$$\text{range} \left( \sum_{i=1}^N \varphi_i \varphi_i^\top \right) = \text{span} \{ \varphi_1, \dots, \varphi_N \}.$$

**Proof.** Let  $R = \sum_{i=1}^N \varphi_i \varphi_i^\top$  and  $S = \text{span} \{ \varphi_1, \dots, \varphi_N \}$ . Since  $R$  is symmetric,  $\text{range } R^\top = \text{range } R$  and  $\text{null } R = (\text{range } R^\top)^\perp = (\text{range } R)^\perp$ . Suppose that  $v \in S^\perp$ . Then  $\varphi_i \perp v$  for all  $i$ . Then  $Rv = \sum_{i=1}^N \varphi_i (\varphi_i^\top v) = 0$ , hence  $v \in \text{null } R = (\text{range } R)^\perp$ . Suppose, on the other hand, that  $v \in (\text{range } R)^\perp = \text{null } R$ . Then  $Rv = 0$ , hence  $0 = v^\top Rv = \sum_{i=1}^N (\varphi_i^\top v)^2$ . Since the last expression is a sum of nonnegative quantities, for it to be zero it must hold  $\varphi_i^\top v = 0$ , that is  $\varphi_i \perp v$  for all  $i$ , hence  $v \in S^\perp$ . We conclude that  $(\text{range } R)^\perp = S^\perp$ . Since  $\text{range } R$  and  $S$  are subspaces of  $\mathbb{R}^p$  (finite-dimensional), taking the orthogonal complement on both sides we obtain  $\text{range } R = S$ .  $\square$

**Corollary 1.3.1** *The normal equations have at least one solution.*

Such solution may not be unique; this happens precisely when  $R$  is singular. Possible reasons: 1) there are not enough data ( $N < p$ , pathological); 2) the explanatory data  $x_i$  do not carry enough information.

### 1.4 Interpretation in terms of projections

**Theorem 1.4.1** *Let  $V$  be a subspace of  $\mathbb{R}^N$ , and  $y \in \mathbb{R}^N$ . If there exists a vector  $v_m \in V$  such that  $\|y - v_m\| \leq \|y - v\|$  for all  $v \in V$ , then  $v_m$  is unique. Moreover, a necessary and sufficient condition for  $v_m$  to be the unique minimizing vector is that  $y - v_m$  is orthogonal to  $V$ .*

Any such  $v_m$  is called the *orthogonal projection* of  $y$  on  $V$ ; such orthogonality is precisely what the normal equations are asking for. Consider again the model  $y_i = \varphi_i^\top \theta + \varepsilon_i$ ,  $i = 1, \dots, N$ .

Let us define  $Y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$ ,  $\Phi = \begin{bmatrix} \varphi_1^\top \\ \vdots \\ \varphi_N^\top \end{bmatrix}$ ,  $E = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_N \end{bmatrix}$ , where  $Y, E \in \mathbb{R}^N$  and  $\Phi \in \mathbb{R}^{N \times p}$ . The model

then reads  $Y = \Phi\theta + E$ . Let  $v_1, \dots, v_p$  be the *columns* of  $\Phi$  (whereas the regressors are its *rows*); then  $V = \text{span} \{ v_1, \dots, v_p \}$  is a subspace of  $\mathbb{R}^N$ , and  $v = \Phi\theta$  is a vector in  $V$ . The least squares problem asks to minimize  $\|Y - \Phi\theta\|^2 = \|Y - v\|^2$ , that is, to find  $v_m = \Phi\hat{\theta}$  such that  $\|Y - v_m\|^2$  is minimal. Now we apply Theorem 1.4.1: if  $Y - v_m \perp V$ , then  $v_m$  is a minimizing vector. Explicitly, if

$$\begin{aligned} Y - \Phi\hat{\theta} &\perp v_i, \quad \text{that is,} \\ v_i^\top (Y - \Phi\hat{\theta}) &= 0 \quad \text{for all } i = 1, \dots, p, \end{aligned} \quad (3)$$

then  $\Phi\hat{\theta}$  is a minimizing vector, and  $\hat{\theta} = [\hat{\theta}_1, \dots, \hat{\theta}_p]^\top$  is the vector of coefficients such that  $\Phi\hat{\theta} = \hat{\theta}_1 v_1 + \dots + \hat{\theta}_p v_p$  is the orthogonal projection of  $Y$  on the space spanned by the columns  $v_1, \dots, v_p$  of  $\Phi$ . Stacking the rows  $v_i^\top$  on each other we get  $\Phi^\top$ , hence stacking the equations (3) on each other we obtain:

$$\Phi^\top (Y - \Phi\hat{\theta}) = 0,$$

which finally yields

$$\Phi^\top \Phi \hat{\theta} = \Phi^\top Y. \quad (4)$$

Equation (4) is just another way to write the normal equations (1), because  $\Phi^\top \Phi = \sum_{i=1}^N \varphi_i \varphi_i^\top$  and  $\Phi^\top Y = \sum_{i=1}^N \varphi_i y_i$ . Since we know that a solution  $\hat{\theta}_{\text{LS}}$  to the normal equations exists, now we also know that at least one minimizer of  $\|Y - \Phi\theta\|^2$  exists, or that there exists a solution to the least squares problem. Note that, according to Theorem 1.4.1, the minimizer  $v_m = \Phi \hat{\theta}_{\text{LS}}$  is unique. This does *not* mean that  $\hat{\theta}_{\text{LS}}$  is also unique! The uniqueness of  $\hat{\theta}_{\text{LS}}$  holds if  $\Phi$  has full rank (linearly independent columns), and this in turn is the case if  $R = \Phi^\top \Phi$  is invertible.

## 1.5 Goodness of fit

Given a collection of numbers  $y_1, \dots, y_N$ , the minimization of its square deviation from a number  $\theta$ :

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N (y_i - \theta)^2,$$

is a particular instance of least squares problem, with the only regressor  $\varphi = 1$ . The solution is the *sample average* of the  $\{y_i\}$ :  $\hat{\theta}_{\text{LS}} = \frac{1}{N} \sum_{i=1}^N y_i = M[y]$ . Hence, the *sample average* is the number from which is minimal the squared deviation of the  $\{y_i\}$ . The attained minimum deviation is the *sample variance* of the  $\{y_i\}$ :  $S.\text{Var}[y] = \frac{1}{N} \sum_{i=1}^N (y_i - M[y])^2$ .

Since the regressor  $\varphi(\cdot) = 1$  is always available irrespective of the nature of the explanatory data, it is often included in least squares problems. Now note that if  $\varphi_1 = 1, \varphi_2(x_i), \dots, \varphi_p(x_i)$  are the regressors, then there is always a parameter  $\hat{\theta}_m$  such that the corresponding sum of squares is exactly the sample variance of  $\{y_i\}$ , namely  $\hat{\theta}_m = [M[y], 0, \dots, 0]^\top$ . It follows that if the explanatory data are really “explaining”, the least squares solution must behave *better* than  $\hat{\theta}_m$ :

$$\frac{1}{N} \sum_{i=1}^N (y_i - \varphi_i^\top \hat{\theta}_{\text{LS}})^2 < \frac{1}{N} \sum_{i=1}^N (y_i - \varphi_i^\top \hat{\theta}_m)^2 = \frac{1}{N} \sum_{i=1}^N (y_i - M[y])^2 = S.\text{Var}[y],$$

We call *residual variance* the quantity  $RV := \frac{1}{N} \sum_{i=1}^N (y_i - \varphi_i^\top \hat{\theta}_{\text{LS}})^2$  and *explained variance* the quantity  $EV := S.\text{Var}[y] - RV$ . Then  $0 \leq EV \leq S.\text{Var}[y]$ . Dividing by  $S.\text{Var}[y]$  and letting  $\rho^2 := \frac{EV}{S.\text{Var}[y]}$ , we obtain finally the relation  $0 \leq \rho^2 \leq 1$ . On one extreme, when  $\rho^2 = 0$  (or equivalently  $RV = S.\text{Var}[y]$ ), the explanatory variables are not adding any information, and the least squares solution is not any more useful than a sample average; on the other extreme, when  $\rho^2 = 1$  (equivalently  $RV = 0$ ), the  $y_i$  are explained *perfectly* by the  $x_i$ , that is, the equations  $y_i = \varphi_i^\top \hat{\theta}_{\text{LS}}$  are all verified *exactly*.

One usually expresses  $\rho^2$  as a percentage, and says something like “the model explains 90% of the variance of the data”. Thus, in general, the closer  $\rho^2$  is to 100%, the better. But be warned: one possible reason why this happens is that *there are just too many regressors*, that is  $p \simeq N$ . This situation is called *over-fitting*. Remember: *the spirit of the least squares method is to use many measures to average out noise, not to identify many parameters!*

## 1.6 Consistency

So far, the functional relation  $y = \varphi(x)^\top \theta^o + \varepsilon$  has been a *model* that we have used to explain the data  $\{(x_i, y_i)\}_{i=1}^N$  through regressor functions, but we have not really pretended that a “true”  $\theta^o$  actually exists. Up to now  $\{x_i\}$ , and consequently  $\{\varphi_i\}$ , have been just vectors, and the results were “algorithmic” in the sense that they guarantee the existence of a solution to the least squares problem, and tell us how to compute it; the numbers  $\{\varepsilon_i\}$  have been there just to model “disturbances”, without any particular requirement other than they be preferably small. Now we assume that there is indeed a  $\theta^o$ , and that the data conform to the model  $y = \varphi(x)^\top \theta^o + \varepsilon$ . We assume, moreover, that  $\{\varepsilon_i\}$  and possibly  $\{\varphi_i\}$  are *random variables*, and examine the consequences on the asymptotic behavior of  $\hat{\theta}_{\text{LS}}$  (its behavior “for large  $N$ ”). Our results will justify the claim that  $\hat{\theta}_{\text{LS}}$  is a good estimator of  $\theta^o$ .

Let the data  $\{(\varphi_i, y_i)\}$  be generated by the model

$$y_i = \varphi_i^\top \theta^o + \varepsilon_i \tag{5}$$

where  $\theta^o \in \mathbb{R}^p$  is understood as the “true” parameter, deterministic but unknown. To see how  $\hat{\theta}_{\text{LS}}$  is related to  $\theta^o$ , substitute (5) in the normal equations (1) and divide by  $N$ :

$$\left( \frac{1}{N} \sum_{i=1}^N \varphi_i \varphi_i^\top \right) \hat{\theta}_{\text{LS}} = \left( \frac{1}{N} \sum_{i=1}^N \varphi_i \varphi_i^\top \right) \theta^o + \frac{1}{N} \sum_{i=1}^N \varphi_i \varepsilon_i$$

In practical situations, we expect that:

1. for big  $N$ , the matrix  $\frac{1}{N} \sum_{i=1}^N \varphi_i \varphi_i^\top$  becomes invertible due to the fact that the  $\varphi_i$  carry more and more information; then the solution  $\hat{\theta}_{\text{LS}}$  is unique, and  $\hat{\theta}_{\text{LS}} = \theta^o + \left( \frac{1}{N} \sum_{i=1}^N \varphi_i \varphi_i^\top \right)^{-1} \frac{1}{N} \sum_{i=1}^N \varphi_i \varepsilon_i$ ;
2. that  $\frac{1}{N} \sum_{i=1}^N \varphi_i \varepsilon_i \rightarrow 0$  as  $N \rightarrow \infty$ , yielding  $\hat{\theta}_{\text{LS}} \rightarrow \theta^o$  ( $\hat{\theta}_{\text{LS}}$  is a *consistent estimator* of  $\theta^o$ ). This idea is in conformity with the principle that underlies least squares optimization: *noise tends to get averaged out and lose importance as the number of measures increases; consequently, it is usually better to take many noisy measures, than to take few precise measures.*

**Theorem 1.6.1** Suppose that  $\{y_i\}_{i=1}^\infty$  are generated by the model  $y_i = \varphi_i^\top \theta^o + \varepsilon_i$ , and that

1.  $\{\varphi_i\}_{i=1}^\infty$  are i.i.d. random vectors, with covariance matrix  $\Sigma := \mathbb{E} [\varphi_i \varphi_i^\top] > 0$ ;
2.  $\{\varepsilon_i\}_{i=1}^\infty$  are independent and identically distributed random variables, with mean  $\mathbb{E} [\varepsilon_i] = 0$ .
3.  $\varepsilon_i$  is independent from  $\varphi_i$  for all  $i$ .

Then  $\hat{\theta}_{\text{LS}} \rightarrow \theta^o$  almost surely as  $N \rightarrow \infty$ .

Sketch of proof:

- use a strong law of large numbers to show that, almost surely,  $\frac{1}{N} \sum_{i=1}^N \varphi_i \varphi_i^\top$  becomes invertible for  $N \geq \bar{N}$  (and converges to the invertible  $\Sigma$ );
- note that  $\mathbb{E} [\varphi_i \varepsilon_i] = 0$ ;
- consequently, by a strong law of large numbers,  $\frac{1}{N} \sum_{i=1}^N \varphi_i \varepsilon_i \rightarrow 0$  almost surely.

**Theorem 1.6.2** Suppose that  $\{y_i\}_{i=1}^\infty$  are generated by the model  $y_i = \varphi_i^\top \theta^o + \varepsilon_i$ , and that

1.  $\{\varphi_i\}_{i=1}^\infty$  are deterministic vectors, such that  $\frac{1}{N} \sum_{i=1}^N \varphi_i \varphi_i^\top \geq aI$  for all  $N \geq \bar{N}$ , and such that  $\|\varphi_i\|^2 \leq A$  for all  $i$ , for certain positive constants  $a, A$ ;
2.  $\{\varepsilon_i\}_{i=1}^\infty$  are independent random variables, not necessarily identically distributed, but such that, for all  $i$ ,  $\mathbb{E} [\varepsilon_i] = 0$  and  $\mathbb{E} [\varepsilon_i^2] \leq c$  for a certain constant  $c \in \mathbb{R}$ .

Then  $\hat{\theta}_{\text{LS}} \rightarrow \theta^o$  almost surely as  $N \rightarrow \infty$ .

Sketch of proof:

- use the first hypothesis to show that, for all  $N \geq \bar{N}$ ,  $\left( \frac{1}{N} \sum_{i=1}^N \varphi_i \varphi_i^\top \right)^{-1}$  exists and is bounded;
- use the bound on  $\|\varphi_i\|^2$ , the second hypothesis, and a strong law of large numbers, to show that  $\frac{1}{N} \sum_{i=1}^N \varphi_i \varepsilon_i \rightarrow 0$  almost surely.

If the  $\varphi_i$  are deterministic and the  $\varepsilon_i$  are i.i.d. with mean 0 and variance  $\sigma^2$ , then  $\hat{\theta}_{\text{LS}} \rightarrow \theta^o$  also in the mean-square sense. Indeed, with compact notation,  $Y = \Phi \theta^o + E$ ; for big  $N$ , we expect that  $\Phi^\top \Phi = \sum_{i=1}^N \varphi_i \varphi_i^\top$  becomes invertible; then

$$\hat{\theta}_{\text{LS}} = (\Phi^\top \Phi)^{-1} \Phi^\top Y = (\Phi^\top \Phi)^{-1} \Phi^\top (\Phi \theta^o + E) = \theta^o + (\Phi^\top \Phi)^{-1} \Phi^\top E.$$

It follows

$$\begin{aligned} \mathbb{E}[\hat{\theta}_{\text{LS}}] &= \theta^o + (\Phi^\top \Phi)^{-1} \Phi^\top \mathbb{E}[E] = \theta^o; \\ \text{Var}[\hat{\theta}_{\text{LS}}] &= (\Phi^\top \Phi)^{-1} \Phi^\top \mathbb{E}[EE^\top] \Phi (\Phi^\top \Phi)^{-1} = (\Phi^\top \Phi)^{-1} \Phi^\top \sigma^2 I \Phi (\Phi^\top \Phi)^{-1} = \sigma^2 (\Phi^\top \Phi)^{-1}. \end{aligned}$$

If  $(\Phi^\top \Phi)^{-1}$  tends to 0, as we expect, then  $\text{Var}[\hat{\theta}_{\text{LS}}] \rightarrow 0$ .

## 2 Linear-algebraic aspects of the least squares method

### 2.1 The Moore-Penrose pseudoinverse

#### 2.1.1 Definition

Given any matrix  $A \in \mathbb{R}^{m \times n}$ , a matrix  $A^+ \in \mathbb{R}^{n \times m}$  is called a *Moore/Penrose pseudo-inverse* of  $A$ , or a pseudoinverse of  $A$  for short, if it satisfies the following properties:

1. The matrix  $AA^+ \in \mathbb{R}^{m \times m}$  is symmetric;
2. The matrix  $A^+A \in \mathbb{R}^{n \times n}$  is symmetric;
3.  $AA^+A = A$ ;
4.  $A^+AA^+ = A^+$ .

**Theorem 2.1.1** *For any matrix  $A \in \mathbb{R}^{m \times n}$ , a pseudoinverse of  $A$  exists, it is unique, and it is uniquely determined by the above four properties.*

To show that a matrix  $B$  is the pseudoinverse of  $A$  it suffices to show that  $B$  satisfies the four properties. For example:

- If  $A$  is square and invertible, then  $A^+ = A^{-1}$ . Indeed, 1)  $AA^{-1} = I$ , which is symmetric; 2)  $A^{-1}A = I$ , which is symmetric; 3)  $AA^{-1}A = IA = A$ ; 4)  $A^{-1}AA^{-1} = IA^{-1} = A^{-1}$ . Since  $A^{-1}$  satisfies the four properties, it is the pseudoinverse of  $A$ . This is where the name *pseudo-inverse* comes from: it is a generalization of the inverse of a square, full rank matrix.
- If  $A$  is a “tall” matrix, meaning that  $m > n$ , and if its columns are linearly independent, then  $A^+$  is a *left inverse* of  $A$ , that is a matrix such that  $A^+A = I$ . More precisely, we have  $A^+ = (A^\top A)^{-1}A^\top$ . As before, one just needs to check that  $(A^\top A)^{-1}A^\top$  satisfies the four properties. For example, the last one reads:

$$((A^\top A)^{-1}A^\top)A((A^\top A)^{-1}A^\top) = (A^\top A)^{-1}(A^\top A)(A^\top A)^{-1}A^\top = (A^\top A)^{-1}A^\top.$$

This is precisely the case of the least squares method, when  $\Phi$  has full rank:  $\Phi^+ = (\Phi^\top \Phi)^{-1}\Phi^\top$

In the general case (singular, non-square matrices), the pseudoinverse can be computed with the aid of a useful matrix decomposition, called the Singular Value Decomposition (SVD).

### 2.2 Least squares solution of a linear system

Given  $A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$ , the linear equation  $Ax = b$  may have no solution, and if a solution exists, it may not be unique. Note that, if  $x$  is a solution, then  $Ax - b = 0$ , hence  $\|Ax - b\|_2 = 0$ .

**Definition 2.2.1** *Given  $A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$ , consider the set*

$$S(A, b) = \{x \in \mathbb{R}^n \mid x \text{ minimizes } \|Ax - b\|_2\}$$

*(A vector  $x \in S(A, b)$  is “almost” a solution, in the sense that it attains the closest possible result to  $\|Ax - b\|_2 = 0$ .) Then any vector  $x^* \in S(A, b)$  of minimum norm, that is*

$$x^* = \arg \min_{x \in S(A, b)} \|x\|_2$$

*is called a least squares solution of the linear system  $Ax = b$ .*

**Theorem 2.2.1** A least squares solution of  $Ax = b$  always exists, it is unique, and it is given by

$$x^* = A^+b.$$

In the least squares method, the problem is to minimize  $\|\Phi\theta - Y\|_2$ . The solution is  $\hat{\theta}_{LS} = \Phi^+Y$  (it is the solution with minimum norm  $\|\theta\|_2$ , if there are more than one). Example: interpolation with a third-degree polynomial  $y = \theta_0 + \theta_1x + \theta_2x^2 + \theta_3x^3$  given the measures  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ :

1. stack  $x_i$  and  $y_i$  in two separate columns  $X$  and  $Y$  of the same dimension;
2. build the “data matrix”  $\Phi$  having  $\varphi_j(x_i) = (x_i)^j$ :

$$\Phi = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ 1 & x_3 & x_3^2 & x_3^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_N & x_N^2 & x_N^3 \end{bmatrix}$$

3. compute  $\hat{\theta}_{LS} = \Phi^+Y$ .

Matlab code:

```
Phi = [ones(length(X),1), X, X.^2, X.^3];
thetaLS = pinv(Phi)*Y;
```

### 3 Identification of dynamical systems

#### 3.1 Model classes

An *autoregressive* (AR) process  $\{y(t)\}_{t=0}^\infty$  is described by the following model:

$$y(t) + a_1y(t-1) + \cdots + a_ny(t-n) = e(t)$$

where  $a_1, \dots, a_n \in \mathbb{R}$ , and  $\{e(t)\}_{t=0}^\infty$  is a sequence of i.i.d. random variables (“process noise”), with mean 0 and unknown variance (not necessarily a disturbance: just an input invisible to the experimenter). We suppose that  $e(t)$  is independent from  $y(s)$  for all  $s < t$ . Let us denote with  $z$  and  $z^{-1}$  the anticipation operator and the delay operator respectively ( $zy(t) = y(t+1)$  and  $z^{-1}y(t) = y(t-1)$ ). Then the model reads  $(1 + a_1z^{-1} + \cdots + a_nz^{-n})y(t) = e(t)$ ; hence the process is the output of a linear system with transfer function

$$W(z) = \frac{z^n}{z^n + a_1z^{n-1} + \cdots + a_n}$$

and with a white noise as an input. If the poles of  $W(z)$ , i.e. the roots of the polynomial  $z^n + a_1z^{n-1} + \cdots + a_n$ , lie inside the interior of the unit circle ( $\{z \in \mathbb{C} \mid |z| < 1\}$ ), then such system is BIBO-stable, and it can be shown that  $\{y(t)\}$  is asymptotically wide-sense stationary.

The process  $\{y(t)\}_{t=0}^\infty$  is *autoregressive with exogenous input* (ARX) if it conforms to the following model:

$$y(t) + a_1y(t-1) + \cdots + a_ny(t-n) = b_0u(t) + b_1u(t-1) + \cdots + b_mu(t-m) + e(t)$$

where  $m \leq n$ ,  $\{e(t)\}$  is a white noise as above, and  $\{u(t)\}_{t=0}^\infty$  is another signal, either random or deterministic, but known to the experimenter. We suppose that  $e(t)$  is independent from  $y(s)$  for all  $s < t$ , and independent from  $u(s)$  for all  $s$ . If  $u(t)$  is stationary and if the roots of  $z^n + a_1z^{n-1} + \cdots + a_n$  lie inside the unit circle, then  $y(t)$  is asymptotically stationary as well.

With respect to both AR and ARX models, the system identification problem is the task of reconstructing  $W(z)$ , that is, finding the parameters  $a_1, \dots, a_n$  and possibly  $b_0, \dots, b_m$ , from measures, respectively  $(y(0), y(1), \dots, y(N))$  or  $((u(0), y(0)), (u(1), y(1)), \dots, (u(N-1), y(N-1)), y(N))$ . Examples: the identification of AR models is used in telecommunications to save bandwidth in the transmission of speech through digital channels. The identification of ARX models may be used to find the transfer function of a plant. AR and ARX models are interesting because their identification can be done with the least squares method.

### 3.2 Application of the least squares method

Consider the ARX model

$$y(t) + \sum_{i=1}^n a_i y(t-i) = \sum_{i=0}^m b_i u(t-i) + e(t) \quad (6)$$

A linear *predictor* of  $y(t)$ , given the past values  $y(t-1)$  ad  $u(t-1)$  is a function of the form

$$\begin{aligned} \hat{y}(t|t-1) &= -\sum_{i=1}^n \hat{a}_i y(t-i) + \sum_{i=0}^m \hat{b}_i u(t-i) \\ &= \begin{bmatrix} -y(t-1) & \cdots & -y(t-n) & u(t) & \cdots & u(t-m) \end{bmatrix} \begin{bmatrix} \hat{a}_1 \\ \vdots \\ \hat{a}_n \\ \hat{b}_0 \\ \vdots \\ \hat{b}_m \end{bmatrix} = \varphi_t^\top \hat{\theta} \end{aligned}$$

used to estimate  $y(t)$  from past measures  $\{y(t-1), \dots, y(t-n), u(t), \dots, u(t-m)\}$ . Predicting  $y(t)$ , we commit an error  $\tilde{y}(t|t-1) = y(t) - \hat{y}(t|t-1)$ , a random quantity with unknown variance. The particular predictor

$$\hat{y}^o(t|t-1) = \begin{bmatrix} -y(t-1) & \cdots & -y(t-n) & u(t) & \cdots & u(t-m) \end{bmatrix} \begin{bmatrix} a_1 \\ \vdots \\ a_n \\ b_0 \\ \vdots \\ b_m \end{bmatrix} = \varphi_t^\top \theta^o$$

is optimal, in the sense that it minimizes such variance; the ideal goal of prediction is to compute  $\hat{y}^o$ , but we do not know  $\theta^o$ , therefore we go for an approximation based on data. The so-called *prediction-error-minimization* (PEM) method prescribes to estimate  $\theta^o$  finding a  $\hat{\theta}$  that minimizes the sum of the squares of the prediction errors (the *residuals*):

$$V = \sum_{t=1}^N \tilde{y}(t|t-1)^2 = \sum_{t=1}^N (y(t) - \hat{y}(t|t-1))^2.$$

We apply the method of least squares. Let

$$y_t := y(t), \quad \varphi_t := \begin{bmatrix} y(t-1) \\ \vdots \\ y(t-n) \\ u(t-1) \\ \vdots \\ u(t-m) \end{bmatrix}, \quad \theta^o := \begin{bmatrix} a_1 \\ \vdots \\ a_n \\ b_0 \\ \vdots \\ b_m \end{bmatrix}, \quad \varepsilon_t := e(t);$$

then (6) can be rewritten  $y_t = \varphi_t^\top \theta^o + \varepsilon_t$  for  $t = 1, \dots, N$ . Now, the method of least squares finds the  $\hat{\theta}_{LS}$  such that the sum  $\sum_{t=1}^N (\tilde{y}_t - \varphi_t^\top \hat{\theta})^2$  is minimal, and this is precisely what the PEM method asks. Note that there is nothing strange in having, among the regressors that explain the sample  $y(t)$ , some samples of the same process itself ( $y(t-1), y(t-2)$  etc.). Indeed this is precisely the reason for the name *autoregressive*: the process *regresses on itself*.

We pose now the question as whether or not, if (6) is the *true* model that generates  $\{y(t)\}$ ,  $\hat{\theta}_{LS}$  converges to  $\theta^o := [a_1 \ \cdots \ a_n \ b_0 \ \cdots \ b_m]^\top$ . Since  $e(t)$  is independent from  $u(\cdot)$  and from the past of  $y(\cdot)$ ,  $\varepsilon_t$  is independent from  $\varphi_t$ , and you could be tempted to apply Theorem 1.6.1. This cannot be done directly, because here the regressors  $\{\varphi_t\}$  are neither independent nor identically distributed. However, the crucial fact is that they are independent from  $\varepsilon_t$ .

**Theorem 3.2.1** Let  $\{y(t)\}_{t=0}^{\infty}$  be an AR process, and suppose that

1.  $e(t)$  is independent from  $y(t-1), y(t-2), \dots$ ;

2. the roots of the polynomial  $z^n + a_1 z^{n-1} + \dots + a_n$  lie inside the interior of the unit circle ( $\{z \in \mathbb{C} \mid |z| < 1\}$ );

then the least squares-estimate  $\hat{\theta}_{LS}$  converges almost surely to  $\theta^o$  as  $N \rightarrow \infty$ .

**Theorem 3.2.2** Let  $\{y(t)\}_{t=0}^{\infty}$  be an ARX process, and suppose that

1. the process  $\{u(t)\}_{t=0}^{\infty}$  is wide sense stationary, with correlation sequence  $r(\tau) = E[u(t)u(t+\tau)]$ ;

2. the Toeplitz matrix

$$M = \begin{bmatrix} r(0) & r(1) & \cdots & r(m-1) \\ r(1) & r(0) & \cdots & r(m-2) \\ \vdots & \vdots & \ddots & \vdots \\ r(m-1) & r(m-2) & \cdots & r(0) \end{bmatrix}$$

is positive definite (this hypothesis is called persistent excitation of the input, and ensures that it carries enough information. Recall that the only real-world issue that can prevent uniqueness in least squares was “the regressors  $\varphi_t$  do not carry enough information”);

3.  $e(t)$  is independent from  $y(t-1), y(t-2), \dots$  and from  $u(s)$  for all  $s$  (this implies that there is no feedback between  $y$  and  $u$ );

4. the roots of the polynomial  $z^n + a_1 z^{n-1} + \dots + a_n$  lie inside the interior of the unit circle;

then the least squares-estimate  $\hat{\theta}_{LS}$  converges almost surely to  $\theta^o$  as  $N \rightarrow \infty$ .

### 3.3 Instrumental variables

Consider now a different situation, in which the process  $y(\cdot)$  is generated by the linear system

$$y(t) = ay(t-1) + bu(t-1) \quad (7)$$

without process noise, with  $u(t)$  independent from  $y(s)$  for all  $s \leq t$  (no feedback). The output  $y(\cdot)$  is not accessible anymore; instead, the experimenter has access to measures of the output corrupted by noise:

$$y_m(t) = y(t) + e(t) \quad (8)$$

where  $\{e(t)\}$  are i.i.d. random variables with mean 0, independent from  $u(s)$  and  $y(s)$  for all  $s \leq t$ . Is it still possible to apply the method of least squares to estimate  $a$  and  $b$ ? It would seem so, because substituting (8) into (7) we obtain

$$\begin{aligned} y_m(t) - e(t) &= a(y_m(t-1) - e(t-1)) + bu(t-1) \\ y_m(t) &= ay_m(t-1) + bu(t-1) + (e(t) - ae(t-1)) \end{aligned}$$

and letting

$$y_t := y_m(t), \quad \varphi_t := \begin{bmatrix} y_m(t-1) \\ u(t-1) \end{bmatrix}, \quad \theta^o := \begin{bmatrix} a \\ b \end{bmatrix}, \quad \varepsilon_t := e(t) - ae(t-1)$$

the model becomes, as before,  $y_t = \varphi_t^\top \theta^o + \varepsilon_t$ . Nevertheless, the regressor and the disturbance are correlated, therefore none of the theorems about almost sure convergence of  $\hat{\theta}_{LS}$  apply. There is no way out of this issue within the standard theory of least squares; but there is a remedy, called the method of *instrumental variables*. Consider the equation from which the least squares solution  $\hat{\theta}_{LS}$  has been derived:

$$\frac{1}{N} \sum_{t=1}^N \varphi_t (\hat{\theta}) = \frac{1}{N} \sum_{t=1}^N \varphi_t (\varphi_t^\top (\theta^o - \hat{\theta}) + \varepsilon_t) = 0.$$

The main idea of least squares is that, in the limit, this equation becomes  $E[\varphi_t(\varphi_t^\top(\theta^o - \hat{\theta}) + \varepsilon_t)] = 0$ , and if 1)  $E[\varphi_t\varphi_t^\top]$  is invertible and 2)  $E[\varphi_t\varepsilon_t] = 0$ , then its only solution is  $\hat{\theta} = \theta^o$ . In the above example this is not the case, because the second condition fails. But if we replace (somehow heuristically) the first occurrence of the regressor  $\varphi_t$  with *another* variable  $\psi_t$  such that 1)  $E[\psi_t\varphi_t^\top]$  is invertible, and 2)  $E[\psi_t\varepsilon_t] = 0$ , then the equation  $E[\psi_t(\varphi_t^\top(\theta^o - \hat{\theta}) + \varepsilon_t)] = 0$  has again the only solution  $\hat{\theta} = \theta^o$ ; therefore it makes sense to solve

$$\left( \sum_{t=1}^N \psi_t \varphi_t^\top \right) \hat{\theta} = \sum_{t=1}^N \psi_t y_t$$

instead of the normal equations. The variable  $\psi_t$ , correlated with  $\varphi_t$  but *not* with  $\varepsilon_t$ , is called *instrumental*; in our case the vector  $\psi_t := \begin{bmatrix} u(t-2) \\ u(t-1) \end{bmatrix}$  is a good instrumental variable, provided that  $E[u(t)^2] \neq 0$  for all  $t$ .

## 4 Machine learning

### 4.1 Statement of the problem

A sequence of i.i.d. random pairs  $(U_1, Y_1), \dots, (U_N, Y_N)$  is observed, where  $U_i \in \mathbb{R}$  and  $Y_i \in \{0, 1\}$ . We wish to study the behavior of functions  $\hat{f} : \mathbb{R} \rightarrow \{0, 1\}$  in providing an estimate  $\hat{Y} = \hat{f}(U)$ . The objective is to find the “best” possible classifier. Following the least squares principle, “best” means that it minimizes

$$\bar{J} = E[(Y_i - \hat{f}(U_i))^2]$$

The best classifier  $\hat{f}$  is called *Bayesian*, but it cannot be computed since it requires knowledge of the data-generation mechanism. Moreover, we restrict the choice to a family of functions  $\mathcal{F} = \{\hat{f}_c\}_{c \in C}$  parameterized by some set  $C$  of indices, and the Bayesian classifier does not necessarily belong to  $\mathcal{F}$ .

**Definition 4.1.1** *Error function:*

$$\bar{J}(c) := E[(Y_i - \hat{f}_c(U_i))^2] = E[\mathbb{1}_{Y_i \neq \hat{f}_c(U_i)}] = P[Y_i \neq \hat{f}_c(U_i)]$$

The best classifier in the family is the one that minimizes  $\bar{J}$ :

$$\bar{c} = \arg \min_{c \in C} \bar{J}(c)$$

**Definition 4.1.2** *Empirical error function, based on the data  $(U_1, Y_1), \dots, (U_N, Y_N)$ :*

$$\hat{J}_N(c) := \frac{1}{N} \sum_{i=1}^N (Y_i - \hat{f}_c(U_i))^2 = \frac{1}{N} \sum_{i=1}^N \mathbb{1}(Y_i \neq \hat{f}_c(U_i))$$

The empirical error is our approximation of the “true” error  $\bar{J}(c)$ . With the sole knowledge of the data we can minimize only this one:

$$\hat{c}_N = \arg \min_{c \in C} \hat{J}_N(c)$$

Three quantities of interest:

- $\bar{J}(\bar{c})$ , the minimum *true* error choosing among all classifiers in  $\mathcal{F}$ ;
- $\hat{J}_N(\hat{c}_N)$ , obtained from the minimization of the empirical error;
- $\bar{J}(\hat{c}_N)$ , the true error obtained from the choice that is optimal for the empirical error.

We expect that:

- For big  $N$ ,  $\hat{J}_N(c)$  is close to  $\bar{J}(c)$  with high probability;

- For  $N \rightarrow \infty$ ,  $\hat{J}_N$  tends to  $\bar{J}$ ;
- For big  $N$ ,  $\hat{J}_N(\hat{c}_N)$  is close to  $\bar{J}(\hat{c}_N)$  and  $\bar{J}(\bar{c})$  with high probability.

In general, however, these statements are false (example with needle functions). The problem arises when  $\mathcal{F}$  is too complex, and it is advisable to stay away from complexity (recall pictures on blackboard). For these properties to hold, we need that as  $N \rightarrow \infty$ , the function  $\hat{J}_N$  tends to  $\bar{J}$  uniformly with respect to  $c$ . Program to follow:

1. If  $\hat{J}_N \rightarrow \bar{J}$  uniformly,  $\min \hat{J}_N \rightarrow \min \bar{J}$ ;
2. Convergence of the empirical distribution of a real random variable to its true distribution;
3. Apply results on empirical distributions to the family of threshold classifiers; show that  $\hat{J}_N \rightarrow \bar{J}$  uniformly; conclude that the minimum of  $\hat{J}_N$  indeed tends to the minimum of  $\bar{J}$ .

## 4.2 Lemma on uniform convergence

**Lemma 4.2.1** Suppose that

1.  $\hat{c}_N := \arg \min_{c \in C} \hat{J}_N(c)$  exists for all  $N$ ;
2.  $\bar{c} := \arg \min_{c \in C} \bar{J}(c)$  exists;
3.  $\hat{J}_N \rightarrow \bar{J}$  uniformly;

Then  $\lim_{N \rightarrow \infty} \bar{J}(\hat{c}_N) = \bar{J}(\bar{c})$ .

## 4.3 Convergence of the empirical distribution

Let  $X_1, \dots, X_N$  be i.i.d. random variables with distribution

$$F(x) := \mathbb{P}[X \leq x] = \mathbb{E}[\mathbb{1}_{X \leq x}]$$

The distribution  $F$  is unknown; estimate it with the *empirical distribution*:

$$\hat{F}_N(x) := \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{X_i \leq x}$$

$\hat{F}_N$  is a good distribution function: it is monotone non-decreasing, continuous from the right, and its limits at  $-\infty$  and  $+\infty$  are 0 and 1 respectively.  $\hat{F}_N(x)$  is an unbiased estimator of  $F(x)$ . Moreover, it is a consistent estimator of  $F(x)$ , because the strong law of large numbers implies:

**Lemma 4.3.1** For each  $x$ ,  $\hat{F}_N(x)$  converges to  $F(x)$  almost surely.

More:  $\hat{F}_N(x)$  converges in probability to  $F(x)$  with exponential rate, due to the following

**Theorem 4.3.1 (Hoeffding's inequality).** Let  $Z_i \in [a_i, b_i]$  be independent, bounded random variables for  $i = 1 \dots N$ , and let  $S_N = \frac{1}{N} \sum_{i=1}^N Z_i$ . Then

$$\mathbb{P}[|S_N - \mathbb{E}[S_N]| \geq \varepsilon] \leq 2 \exp\left(\frac{-2N^2\varepsilon^2}{\sum_{i=1}^N (b_i - a_i)^2}\right).$$

(Example on tossing a coin 100 times: Hoeffding's bound is much tighter than Čebyšev's one.) Apply Hoeffding's inequality to the empirical distribution. For a fixed  $x$ , let

$$\begin{aligned} Z_i &= \mathbb{1}_{X_i \leq x} &\Rightarrow \mathbb{E}[Z_i] &= F(x) \\ S_N &= \hat{F}_N(x) &\Rightarrow \mathbb{E}[S_N] &= F(x) \\ [a_i, b_i] &= [0, 1] \end{aligned}$$

Hoeffding's inequality yields

$$\mathbb{P} \left[ |\hat{F}_N(x) - F(x)| \geq \varepsilon \right] \leq 2e^{-2N\varepsilon^2}$$

Hence, at all  $x$  the empirical distribution converges *in probability* to the true distribution, with exponential rate of convergence. Exponential convergence is so strong that it is enough to re-establish *almost surely* convergence as well. Not yet enough. More:  $\hat{F}_N$  converges *uniformly* to  $F$ , due to the following

**Theorem 4.3.2 (Glivenko/Cantelli).** *Let  $\{X_i\}$  be i.i.d. variables with distribution  $F(x)$ . Then*

$$\lim_{N \rightarrow \infty} \sup_{x \in \mathbb{R}} |\hat{F}_N(x) - F(x)| = 0 \quad \text{almost surely.}$$

Probabilistic bounds for fixed  $\varepsilon, N$ , applying Hoeffding's inequality to the proof of Glivenko/Cantelli:

**Lemma 4.3.2** *Fix  $\varepsilon > 0, N$ . Then*

$$\mathbb{P} \left[ \sup_{x \in \mathbb{R}} |\hat{F}_N(x) - F(x)| \geq \varepsilon \right] \leq \frac{4}{\varepsilon} e^{-N\varepsilon^2/2}$$

(The exponent now has  $1/2$  instead of  $2$ , because the proof of Glivenko/Cantelli's theorem is based on splitting in subintervals of length  $\varepsilon/2$ .)

#### 4.4 Extension to two variables, one of which is binary

Let  $\{(X_i, Y_i)\}$  be i.i.d., where  $X_i \in \mathbb{R}$  and  $Y_i \in \{0, 1\}$ . Let  $\mathbb{P}[Y_i = 0] = \alpha$ . Let

$$F^0(x) = \mathbb{P}[X_i \leq x, Y_i = 0] \quad (\text{tends to } \alpha \text{ for } x \rightarrow \infty)$$

$$F^1(x) = \mathbb{P}[X_i \leq x, Y_i = 1] \quad (\text{tends to } 1 - \alpha \text{ for } x \rightarrow \infty)$$

$$\hat{F}_N^0(x) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}\{X_i \leq x, Y_i = 0\} \quad (\text{empirical version of } F^0(x))$$

$$\hat{F}_N^1(x) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}\{X_i \leq x, Y_i = 1\} \quad (\text{empirical version of } F^1(x))$$

Glivenko/Cantelli-like theorem:

$$\sup_{x \in \mathbb{R}} |\hat{F}_N^0(x) - F^0(x)| \rightarrow 0 \quad \text{almost surely;}$$

$$\sup_{x \in \mathbb{R}} |\hat{F}_N^1(x) - F^1(x)| \rightarrow 0 \quad \text{almost surely.}$$

Hoeffding-like probabilistic bound for fixed  $\varepsilon, N$ , like Lemma 4.3.2:

**Lemma 4.4.1** *Let  $\{(X_i, Y_i)\}$  be i.i.d. where  $X_i \in \mathbb{R}$  and  $Y_i \in \{0, 1\}$ . Let  $\mathbb{P}[Y_i = 0] = \alpha$ . Then, for fixed  $\varepsilon > 0, N$ ,*

$$\mathbb{P} \left[ \sup_{x \in \mathbb{R}} |\hat{F}_N^0(x) - F^0(x)| \geq \varepsilon \right] \leq \frac{4\alpha}{\varepsilon} e^{-N\varepsilon^2/2}$$

$$\mathbb{P} \left[ \sup_{x \in \mathbb{R}} |\hat{F}_N^1(x) - F^1(x)| \geq \varepsilon \right] \leq \frac{4(1-\alpha)}{\varepsilon} e^{-N\varepsilon^2/2}$$

## 4.5 Threshold classifiers

Consider the family of *threshold classifiers*:

$$\mathcal{F} = \{\mathbb{1}_{(-\infty, c]}(\cdot) \mid c \in \mathbb{R}\}.$$

Existence of  $\hat{c}_N$ : always. Uniqueness: never. Existence and uniqueness of  $\bar{c}$ : It depends on the data-generation rule. Assume as hypothesis.

**Theorem 4.5.1** Let  $(U_1, Y_1), \dots, (U_N, Y_N)$  be i.i.d. where  $U_i$  has continuous distribution  $F(u)$  and  $Y_i \in \{0, 1\}$ . Assume that  $\bar{c}$  exists. Then:

1. Almost surely  $\hat{J}_N \rightarrow \bar{J}$  uniformly;
2. Almost surely  $\bar{J}(\hat{c}_N) \rightarrow \bar{J}(\bar{c})$ ;
3. For fixed  $\varepsilon > 0$  and  $N$ , it holds

$$P \left[ \sup_{c \in C} |\hat{J}_N(c) - \bar{J}(c)| \geq \varepsilon \right] \leq \frac{8}{\varepsilon} e^{-N\varepsilon^2/8}.$$

**Proof.** For each  $U_i$  define  $V_i = -U_i$  and denote  $G(u)$  the distribution of  $V_i$ . Define as before

$$\begin{aligned} F^0(u) &= P[U_i \leq u, Y_i = 0] & \hat{F}_N^0(u) &= \frac{1}{N} \sum_{i=1}^N \mathbb{1}\{U_i \leq u, Y_i = 0\} \\ G^1(u) &= P[V_i \leq u, Y_i = 1] & \hat{G}_N^1(u) &= \frac{1}{N} \sum_{i=1}^N \mathbb{1}\{V_i \leq u, Y_i = 1\} \end{aligned}$$

We can show that:

$$\bar{J}(c) = F^0(c) + G^1(-c)$$

And for its empirical counterpart:

$$\hat{J}_N(c) = \hat{F}_N^0(c) + \hat{G}_N^1(-c) - \frac{1}{N} \sum_{i=1}^N \mathbb{1}(V_i = -c, Y_i = 1)$$

It follows:

$$|\hat{J}_N(c) - \bar{J}(c)| \leq \left| \hat{F}_N^0(c) - F^0(c) \right| + \left| \hat{G}_N^1(-c) - G^1(-c) \right| + \left| \frac{1}{N} \sum_{i=1}^N \mathbb{1}(V_i = -c, Y_i = 1) \right|$$

As  $N \rightarrow \infty$ , the supremum over  $C$  of the first two terms tends almost surely to zero by Glivenko/Cantelli, and the third one tends almost surely to zero by the strong law of large numbers. This establishes point 1. Apply the lemma on uniform convergence. This establishes point 2. Now fix  $\varepsilon$  and  $N$ . By Lemma 4.4.1,

$$\begin{aligned} P[\mathcal{A}] &= P \left[ \sup_{c \in C} \left| \hat{F}_N^0(c) - F^0(c) \right| \geq \varepsilon \right] \leq \frac{4\alpha}{\varepsilon} e^{-N\varepsilon^2/2} \\ P[\mathcal{B}] &= P \left[ \sup_{c \in C} \left| \hat{G}_N^1(-c) - G^1(-c) \right| \geq \varepsilon \right] \leq \frac{4(1-\alpha)}{\varepsilon} e^{-N\varepsilon^2/2} \end{aligned}$$

On the other hand, the event  $\left| \frac{1}{N} \sum_{i=1}^N \mathbb{1}(V_i = -c, Y_i = 1) \right| \geq \varepsilon$  has probability 0 for all  $\varepsilon > 0$ , since  $V_i$  has continuous distribution. Putting together the inequalities,

$$P \left[ \sup_{c \in C} |\hat{J}_N(c) - \bar{J}(c)| \geq 2\varepsilon \right] \leq \frac{4}{\varepsilon} e^{-N\varepsilon^2/2}$$

Point 3 of the claim follows from a change of variable  $2\varepsilon = \varepsilon'$ .  $\square$

One can establish bounds for the minima as well (they follow from the proof of the lemma on uniform convergence): for fixed  $\varepsilon > 0$  and  $N$ , both the events  $\bar{J}(\hat{c}_N) \leq J_N(\hat{c}_N) + \varepsilon$  and  $\bar{J}(\hat{c}_N) \leq \bar{J}(\bar{c}) + 2\varepsilon$  have probability at least  $1 - \frac{8}{\varepsilon} e^{-N\varepsilon^2/8}$ .

Almost sure convergence: good. Probabilistic bounds coming from Hoeffding's inequality: not impressive unless there are many data available. Not the tightest bound available. In practice, the empirical error behaves better than how the bound says. However, this bound is remarkable because it is *distribution-free*: it does not depend on either the distribution of the data or the mechanism that links input to output.

## 5 The LSCR method

### 5.1 Introduction and motivation

Suppose that  $\{y_1, \dots, y_N\}$  are i.i.d. random variables drawn from a Gaussian distribution  $N(\theta^o, \sigma^2)$  whose mean  $\theta^o$  is not known. We want to extract information about  $\theta^o$  from the data. The sample can be seen as  $N$  measures  $y_i = \theta^o + \varepsilon_i$ ,  $i = 1, \dots, N$ , to be “explained” by the “true” parameter  $\theta^o$ , corrupted by errors  $\varepsilon_i \sim N(0, \sigma^2)$ . LS approach: provide a *point estimate* of  $\theta^o$  minimizing the sum of the squared deviation. Result:  $\hat{\theta} = \frac{1}{N} \sum_{i=1}^N y_i$ , the *sample average*, which is an unbiased estimator of  $\theta^o$  ( $E[\hat{\theta}] = \theta^o$ ). A “guarantee” on the usefulness of  $\hat{\theta}$  should be a “certificate” that  $\hat{\theta}$  and  $\theta^o$  are *probably* close. Such “certificate” usually comes in two forms: 1) the variance of  $\hat{\theta}$  (the smaller, the better); 2) a statement like “the probability  $P[|\hat{\theta} - \theta^o| > d]$  is  $\alpha$ ”, for a certain threshold  $d$  and significance  $\alpha$ , say  $\alpha = 5\%$ . The second “certificate” can also be stated: the interval

$$I_\alpha = [\hat{\theta} - d, \hat{\theta} + d]$$

contains  $\theta^o$  with *probability*  $1 - \alpha$ . Be careful about the use of the word “probability”. *Before the sample is drawn*, the interval  $I_\alpha$  is random, and whether or not it will contain  $\theta^o$  depends on the outcome of the experiment; namely, it will happen with *probability*  $1 - \alpha$ . But *after the sample has been drawn* the interval is deterministic. Either it contains  $\theta^o$  or it does not. Thus,  $1 - \alpha$  is not anymore the “probability” of anything; the usual name for it is *confidence*. For example, in applied statistics the following interval

$$I_\alpha = \left[ \hat{\theta} - \frac{\bar{s}t_\alpha}{\sqrt{N}}, \hat{\theta} + \frac{\bar{s}t_\alpha}{\sqrt{N}} \right]$$

is often used as a  $(1 - \alpha)$ -confidence interval for  $\theta^o$ . Here,  $\bar{s}$  is the square root of the *sample variance*  $\bar{s}^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \hat{\theta})^2$ , which is an unbiased estimator of  $\sigma^2$ , and  $t_\alpha$  is a “percentile” such that  $\int_{-t_\alpha}^{t_\alpha} f(x) dx = 1 - \alpha$ , where  $f(x)$  is the density of a Student's  $t(N-1)$  random variable<sup>1</sup>. This is a beautiful *pivotal* result: it does not depend on the knowledge of any parameter  $\theta^o$  or  $\sigma^2$ , only on the number of data  $N$ . But it has a fundamental drawback: *it still depends crucially on Gaussianity*. Is it possible to provide a confidence interval for  $\theta^o$  without restricting the distribution of the errors to some parametric family?

### 5.2 Confidence intervals with the LSCR method

An *Abelian group*  $(G, +)$  is a set  $G$  endowed with an operation  $+$  such that: 1) for all  $a, b, c \in G$ ,  $((a + b) + c) = (a + (b + c))$  (associativity); 2) there exist an element  $0 \in G$  (zero) such that  $0 + a = a + 0 = a$  for all  $a \in G$ ; 3) for all  $a \in G$ , there exist an element  $-a \in G$  (inverse of  $a$ ) such that  $a + (-a) = (-a) + a = 0$ ; 4) for all  $a, b \in G$ ,  $a + b = b + a$  (commutativity). A *subgroup* of  $(G, +)$  is a set  $H \subseteq G$  which is closed with respect to the operation  $+$  (if  $a, b \in H$ , then also  $a + b \in H$ ). The simplest nontrivial *finite* Abelian group is  $B = (\{\circ, \bullet\}, +)$ , where addition is defined in the following way:

+	○	•
○	○	•
•	•	○
○	•	○

<sup>1</sup>The shape of a  $t(N-1)$  density resembles that of a Gaussian  $N(0, 1)$ , namely it is symmetric around 0 and is close to the Gaussian one, but it has *fatter tails*. As  $N \rightarrow \infty$ , the  $t(N-1)$  density converges to the  $N(0, 1)$  density.

(its zero is  $\circ$ ; recall modulo 2 arithmetics). And here is a group that will serve us for a “canonical” example:

	1	2	3	4	5	6	7
$G_7 =$	$I_1$	•	•	○	•	•	○
	$I_2$	•	○	•	•	○	•
	$I_3$	○	•	•	○	•	•
	$I_4$	•	•	○	○	○	•
	$I_5$	•	○	•	○	•	○
	$I_6$	○	•	•	•	○	•
	$I_7$	○	○	○	•	•	•
	$I_8$	○	○	○	○	○	○

It is a subgroup of  $(B^7, +)$  containing 8 elements, i.e. the rows  $I_1, \dots, I_8$ ; addition is defined component-wise. Any two rows in the table have another row as their sum; the zero is  $I_8$ , and the inverse of a row is the row itself.

Suppose that we are given 7 (very few) measures  $y_i = \theta^\circ + \varepsilon_i$  for  $i = 1, \dots, 7$ , where  $\{\varepsilon_i\}$  are *continuous variables with a density centered around zero, independent but not necessarily identically distributed*. Consider for each measure an affine function:  $f_i(\theta) = y_i - \theta$ . For each of the first 7 elements  $I_i = I_1, \dots, I_7$  of the group, build a point-wise average of the 4 functions whose indices are marked • in the corresponding row:

$$g_1(\theta) = \frac{1}{4} \sum_{k \in \{1, 2, 4, 5\}} f_k(\theta) = \frac{1}{4} \sum_{k \in \{1, 2, 4, 5\}} y_k - \theta = \frac{1}{4} \sum_{k \in \{1, 2, 4, 5\}} \theta^\circ + \varepsilon_k - \theta = (\theta^\circ - \theta) + \frac{1}{4} \sum_{k \in \{1, 2, 4, 5\}} \varepsilon_k,$$

etc. Each of the  $g_i$  has exactly one intersection  $\theta_i$  with the  $\theta$ -axis. The  $\theta_i$ ’s are *random* numbers coming in some order; since the errors  $\varepsilon_i$  are independent and have continuous distributions, almost surely no two of the  $\theta_i$  coincide, and none of them equals  $\theta^\circ$ . Furthermore, let  $\bar{\theta}_1, \dots, \bar{\theta}_7$  denote the same seven numbers, but *sorted*, i.e. such that  $\bar{\theta}_1 < \bar{\theta}_2 < \dots < \bar{\theta}_7$ . They split the  $\theta$ -axis in 8 intervals (the outermost ones are semi-infinite). Now note that at  $\theta = \theta^\circ$  it holds

$$g_1(\theta^\circ) = \frac{1}{4} \sum_{k \in \{1, 2, 4, 5\}} \varepsilon_k, \quad \dots, \quad g_7(\theta^\circ) = \frac{1}{4} \sum_{k \in \{4, 5, 6, 7\}} \varepsilon_k.$$

Each  $g_i(\theta^\circ)$  is a sum of independent variables, *with density symmetric around 0*. Hence,  $g_i(\theta^\circ)$  has also a density symmetric around 0, and it has equal probabilities of being positive or negative. The fundamental idea of LSCR is to compare the signs of each  $g_i(\theta^\circ)$ , depending on which one of the 8 intervals contains  $\theta^\circ$ . If  $\theta^\circ \in (-\infty, \bar{\theta}_1)$ , then  $g_1(\theta^\circ) > 0, \dots, g_7(\theta^\circ) > 0$ ; if  $\theta^\circ \in (\bar{\theta}_7, \infty)$ , then  $g_1(\theta^\circ) < 0, \dots, g_7(\theta^\circ) < 0$ ; if  $\theta^\circ$  belongs to the second interval  $(\bar{\theta}_1, \bar{\theta}_2)$ , then one and only one among  $g_1(\theta^\circ), g_2(\theta^\circ), \dots, g_7(\theta^\circ)$  has negative sign, etc. *What is the probability of these events?* The first one ( $\theta^\circ \in (-\infty, \bar{\theta}_1)$ , hence  $g_1(\theta^\circ) > 0, \dots, g_7(\theta^\circ) > 0$ ) happens when

$$\begin{array}{ccccccccc} \varepsilon_1 & +\varepsilon_2 & & +\varepsilon_4 & +\varepsilon_5 & & & & > 0, \\ \varepsilon_1 & & +\varepsilon_3 & +\varepsilon_4 & & +\varepsilon_6 & & & > 0, \\ & \varepsilon_2 & +\varepsilon_3 & & +\varepsilon_5 & +\varepsilon_6 & & & > 0, \\ \varepsilon_1 & +\varepsilon_2 & & & & +\varepsilon_6 & +\varepsilon_7 & & > 0, \\ \varepsilon_1 & & +\varepsilon_3 & & +\varepsilon_5 & & +\varepsilon_7 & & > 0, \\ \varepsilon_2 & +\varepsilon_3 & +\varepsilon_4 & & & & +\varepsilon_7 & & > 0, \\ & & & \varepsilon_4 & +\varepsilon_5 & +\varepsilon_6 & +\varepsilon_7 & & > 0. \end{array} \tag{9}$$

If  $\theta^\circ \in (\bar{\theta}_1, \bar{\theta}_2)$ , then exactly one  $g_i(\theta^\circ)$  among  $g_1(\theta^\circ), \dots, g_7(\theta^\circ)$  has negative sign, hence all the other values at  $\theta^\circ$  (including  $g_8(\theta^\circ) = 0$ ) are greater than it. Suppose, without loss of generality, that such  $g_i$  is  $g_1$ . Then

$$\begin{array}{ccccccccc} & & & 0 & & & & & \\ \varepsilon_1 & & +\varepsilon_3 & +\varepsilon_4 & & +\varepsilon_6 & & & > \varepsilon_1 + \varepsilon_2 + \varepsilon_4 + \varepsilon_5, \\ & \varepsilon_2 & +\varepsilon_3 & & +\varepsilon_5 & +\varepsilon_6 & & & > \varepsilon_1 + \varepsilon_2 + \varepsilon_4 + \varepsilon_5, \\ \varepsilon_1 & +\varepsilon_2 & & & +\varepsilon_6 & +\varepsilon_7 & & & > \varepsilon_1 + \varepsilon_2 + \varepsilon_4 + \varepsilon_5, \\ \varepsilon_1 & & +\varepsilon_3 & & +\varepsilon_5 & & +\varepsilon_7 & & > \varepsilon_1 + \varepsilon_2 + \varepsilon_4 + \varepsilon_5, \\ \varepsilon_2 & +\varepsilon_3 & +\varepsilon_4 & & & +\varepsilon_7 & & & > \varepsilon_1 + \varepsilon_2 + \varepsilon_4 + \varepsilon_5, \\ & & & \varepsilon_4 & +\varepsilon_5 & +\varepsilon_6 & +\varepsilon_7 & & > \varepsilon_1 + \varepsilon_2 + \varepsilon_4 + \varepsilon_5. \end{array}$$

Simplifying and bringing every term to the left-hand side we obtain

$$\begin{array}{ccccccccc}
 -\varepsilon_1 & -\varepsilon_2 & \cdots & -\varepsilon_4 & -\varepsilon_5 & & & & > 0, \\
 & -\varepsilon_2 & +\varepsilon_3 & & -\varepsilon_5 & +\varepsilon_6 & & & > 0, \\
 -\varepsilon_1 & & +\varepsilon_3 & -\varepsilon_4 & & +\varepsilon_6 & & & > 0, \\
 & & & -\varepsilon_4 & -\varepsilon_5 & +\varepsilon_6 & +\varepsilon_7 & & > 0, \\
 & & -\varepsilon_2 & +\varepsilon_3 & -\varepsilon_4 & & +\varepsilon_7 & & > 0, \\
 -\varepsilon_1 & & +\varepsilon_3 & & -\varepsilon_5 & & +\varepsilon_7 & & > 0, \\
 -\varepsilon_1 & -\varepsilon_2 & & & & +\varepsilon_6 & +\varepsilon_7 & & > 0.
 \end{array} \tag{10}$$

The sets of inequalities (9) and (10) are very similar, the differences being that the inequalities appear in different orders, and that some of the signs are changed. (9) and (10) contain terms with the same indices due to the group structure, and they hold with the same probability because all the  $\varepsilon_i$  have symmetric densities. The same reasoning applies to the entire sets of inequalities. Hence the probabilities that  $\theta^o \in (-\infty, \bar{\theta}_1)$  and that  $\theta^o \in (\bar{\theta}_1, \bar{\theta}_2)$  are equal. But the very same procedure could have been applied to any of the intervals; hence  $\theta^o$  belongs to any of the intervals with the same probability  $\frac{1}{8}$ . And finally, the probability that  $\theta^o$  belongs to any of the two outermost intervals is  $\frac{1}{8} + \frac{1}{8} = \frac{1}{4}$ , hence:

**Theorem 5.2.1** (Campi, Weyer). Let  $y_i = \theta^o + \varepsilon_i$ ,  $i = 1, \dots, 7$ ; suppose that  $\{\varepsilon_i\}$  are continuous variables with a density centered around zero (with mean zero), independent but not necessarily identically distributed. Construct the functions  $g_1(\cdot), \dots, g_7(\cdot)$  and sort their intersections with the  $\theta$ -axis,  $\bar{\theta}_1 < \dots < \bar{\theta}_7$ , as above. Then  $I := [\bar{\theta}_1, \bar{\theta}_7]$  is a  $\frac{3}{4} = 75\%$ -confidence interval for  $\theta^o$ .

The applicability of the method is not limited to 7 measures. To cope with the general,  $N$ -measures case one must construct a suitable subgroup of  $(B^N, +)$ , “balanced” in the sense that it should be fairly small, the number of bullets ( $\bullet$ ) should be the same in each row, and it should be approximately half the size of a row. This construction is easy when  $N = 2^n - 1$  for some  $n$ . For  $N = 3 = 2^2 - 1$  and  $N = 7 = 2^3 - 1$  we have:

	1	2	3	4	5	6	7		1	2	3	4	5	6	7
$G_3 =$	$I_1$	$\bullet$	$\bullet$	$\circ$	$\bullet$	$\bullet$	$\circ$	$\circ$	$I_1$						$\circ$
	$I_2$	$\bullet$	$\circ$	$\bullet$	$\bullet$	$\circ$	$\bullet$	$\circ$	$I_2$	$G_3$		$G_3$			$\circ$
	$I_3$	$\circ$	$\bullet$	$\bullet$	$\circ$	$\bullet$	$\bullet$	$\circ$	$I_3$						$\circ$
	$I_4$	$\circ$	$\circ$	$\circ$					$I_4$						$\bullet$
$G_7 =$				$I_4$	$\bullet$	$\bullet$	$\circ$	$\circ$	$\circ$	$\bullet$					
				$I_5$	$\bullet$	$\circ$	$\bullet$	$\circ$	$\bullet$	$\circ$	$\bullet$				
				$I_6$	$\circ$	$\bullet$	$\bullet$	$\bullet$	$\circ$	$\circ$	$\bullet$				
				$I_7$	$\circ$	$\circ$	$\circ$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\circ$	$\circ$	$\circ$	$\bullet$
				$I_8$	$\circ$	$\circ$	$\circ$	$\circ$	$\circ$						

Balanced groups for the next powers of 2 are built recursively in this fashion. Since approximately  $\frac{N}{2}$  noise terms are averaged in each intersection, the confidence interval “shrinks” towards  $\theta^o$  as  $N$  increases. Since all of the  $N+1$  intervals have the same probability  $\frac{1}{N+1}$ , some more intervals can be discarded, other than the two outermost ones, to obtain a smaller interval (with of course a smaller confidence). Thus, the confidence is “tunable”, although not continuously as in the Gaussian case.

### 5.3 The case with inputs

LSCR allows to find a confidence interval for  $\theta^o$  also when an input is involved. Let  $y_i = \theta^o u_i + \varepsilon_i$  for  $i = 1, \dots, 7$ . Suppose that  $\{\varepsilon_i\}$  are independent variables with a density symmetric around zero. The inputs  $u_i$  can be deterministic or random, but if we want the above construction to be useful and the confidence interval to “shrink” towards  $\theta^o$  for big  $N$ , they should “stay away from zero”. If they are deterministic, they should be nonzero, and should not tend to zero as  $N \rightarrow \infty$ ; if they are random, they should be independent from  $\varepsilon_i$ , and have nonzero mean. Consider, for each measure, an affine function:  $f_i(\theta) = y_i - \theta u_i$ . Differently from before, these functions may have different slopes. Exactly as before, they intersect the  $\theta$ -axis at one point (unless  $u_i = 0$ , but we have excluded this case). For each  $I_i = I_1, \dots, I_7$  build an average of 4 different  $f_k$ :

$$g_1(\theta) = \frac{1}{4} \sum_{k \in \{1, 2, 4, 5\}} f_k(\theta) = \frac{1}{4} \sum_{k \in \{1, 2, 4, 5\}} y_k - \theta u_k = (\theta^o - \theta) \frac{1}{4} \sum_{k \in \{1, 2, 4, 5\}} u_k + \frac{1}{4} \sum_{k \in \{1, 2, 4, 5\}} \varepsilon_k,$$

etc. The reasoning of the above section applies without other changes, and  $[\bar{\theta}_1, \bar{\theta}_7]$  is a 75%-confidence interval for  $\theta^o$ . What happens if the inputs  $\{u_i\}$  are random, independent from the  $\{\varepsilon_i\}$ , but their mean is 0? Consider the intersection of (say)  $g_1(\theta)$  with the  $\theta$ -axis. It is the point  $\theta_1$  such that

$$0 = g_1(\theta_1) = (\theta^o - \theta_1) \left( \frac{1}{4} \sum_{k \in \{1,2,4,5\}} u_k \right) + \frac{1}{4} \sum_{k \in \{1,2,4,5\}} \varepsilon_k; \quad \theta_1 = \theta^o + \frac{\frac{1}{4} \sum_{k \in \{1,2,4,5\}} \varepsilon_k}{\frac{1}{4} \sum_{k \in \{1,2,4,5\}} u_k}.$$

As more measures come, the term  $\frac{1}{N+1} \sum_k u_k$  at the denominator of the last expression, which is also the coefficient of  $\theta$  in the expression of  $g_1(\theta)$ , tends to 0. The same happens for all the  $g_i$ ; hence the straight lines corresponding to the  $g_i$  are approximately horizontal, their intersections tend to be far from each other, and the confidence interval is practically useless. Here is the remedy: suppose, to fix ideas, that the inputs are i.i.d., that  $E[u_i] = 0$ , and  $E[u_i^2] > 0$ . We consider, instead of the affine functions above, the following ones:

$$f_i(\theta) = (y_i - \theta u_i) u_i = y_i u_i - \theta u_i^2,$$

$$g_1(\theta) = \frac{1}{4} \sum_{k \in \{1,2,4,5\}} f_k(\theta) = \frac{1}{4} \sum_{k \in \{1,2,4,5\}} \theta^o u_k^2 + \varepsilon_k u_k - \theta u_k^2 = (\theta^o - \theta) \left( \frac{1}{4} \sum_{k \in \{1,2,4,5\}} u_k^2 \right) + \frac{1}{4} \sum_{k \in \{1,2,4,5\}} \varepsilon_k u_k,$$

etc. Intersections with  $\theta$ -axis:

$$0 = g_1(\theta_1) = (\theta^o - \theta_1) \left( \frac{1}{4} \sum_{k \in \{1,2,4,5\}} u_k^2 \right) + \frac{1}{4} \sum_{k \in \{1,2,4,5\}} \varepsilon_k u_k; \quad \theta_1 = \theta^o + \frac{\frac{1}{4} \sum_{k \in \{1,2,4,5\}} \varepsilon_k u_k}{\frac{1}{4} \sum_{k \in \{1,2,4,5\}} u_k^2}.$$

The denominator is OK now: as more measures come, the term  $\frac{1}{N+1} \sum_k u_k^2$  tends to  $E[u_i^2]$ , and the functions  $g_i$  have slopes with comparable magnitude. On the other hand, as far as the  $\{u_i\}$  are independent from the  $\{\varepsilon_i\}$ , the terms  $\varepsilon_k u_k$  still have densities symmetric around 0, so that the fundamental idea of LSCR (sets of inequalities having same probability) applies. Thus,  $[\bar{\theta}_1, \bar{\theta}_7]$  is once again a 75%-confidence interval for  $\theta^o$ .

## 5.4 Leave-out Sign-dominant Correlation Regions

The LSCR method generalizes to system identification. Consider, for example, the simple AR process  $y_i = \theta^o y_{i-1} + \varepsilon_i$ . We cannot use anymore affine functions like  $f_i(\theta) = y_i - \theta y_{i-1}$ , because typically  $E[y_{i-1}] = 0$  (confidence intervals get larger and larger). Choose the following functions instead:

$$f_i(\theta) = (y_{i+1} - \theta y_i)(y_i - \theta y_{i-1})$$

$$g_1(\theta) = \frac{1}{4} \sum_{k \in \{1,2,4,5\}} f_k(\theta) = \frac{1}{4} \sum_{k \in \{1,2,4,5\}} ((\theta^o - \theta)y_k + \varepsilon_{k+1}) ((\theta^o - \theta)y_{k-1} + \varepsilon_k)$$

$$= \frac{(\theta^o - \theta)^2}{4} \sum_{k \in \{1,2,4,5\}} y_k y_{k-1} + \frac{\theta^o - \theta}{4} \sum_{k \in \{1,2,4,5\}} y_k \varepsilon_k + \frac{\theta^o - \theta}{4} \sum_{k \in \{1,2,4,5\}} y_{k-1} \varepsilon_{k+1} + \frac{1}{4} \sum_{k \in \{1,2,4,5\}} \varepsilon_{k+1} \varepsilon_k,$$

etc. The  $g_i$  are parabolas, their intersections with the  $\theta$ -axis are not anymore unique. However, they split the  $\theta$ -axis in 8 regions, that are unions of one or more disjoint intervals, where one out of 8 situations happen: 1) all the  $g_i$  are positive; 2) exactly 1 of the  $g_i$  is negative, 3) exactly 2 of the  $g_i$  are negative, ..., 8) all the  $g_i$  are negative. For each region we repeat the reasoning: what happens if  $\theta^o$  belongs to this region? Note that since  $\varepsilon_{k+1}$  and  $\varepsilon_k$  are independent, their product  $\varepsilon_{k+1} \varepsilon_k$  also has symmetric density. Hence  $g_1(\theta^o) = \frac{1}{4} \sum_{k \in \{1,2,4,5\}} \varepsilon_{k+1} \varepsilon_k$  is again an average of terms with symmetric density. Repeat the LSCR proof, which is based only on the symmetry and the group structure, to conclude that the regions have equal probability  $\frac{1}{8}$ . A function like

$$g_1(\theta) = \frac{1}{4} \sum_{k \in \{1,2,4,5\}} (y_{k+1} - \theta y_k)(y_k - \theta y_{k-1}) = \frac{1}{4} \sum_{k \in \{1,2,4,5\}} \varepsilon_k(\theta) \varepsilon_{k-1}(\theta)$$

is an empirical 1-lag correlation; this explains the name LSCR: Leave-out Sign-dominant Correlation Regions.

## 6 Interval predictor models

### 6.1 Prediction intervals without explanatory data

The least squares method tells us how to construct a model  $y = \hat{a}_{\text{LS}} + \hat{b}_{\text{LS}} x$  from measures  $(x_1, y_1), \dots, (x_N, y_N)$ . One of the objectives of such a model is to *predict* the value of a future variable  $y_{N+1}$  when the corresponding explanatory variable  $x_{N+1}$  will be available:  $\hat{y}_{N+1} = \hat{a}_{\text{LS}} + \hat{b}_{\text{LS}} x_{N+1}$ . But what does it mean that  $\hat{y}_{N+1}$  *predicts* the value of  $y_{N+1}$ ?  $\hat{y}_{N+1}$  is just a “plausible” value, that comes with the hope that the future observation be close to it. Our way to quantify “hope” is probability, hence the prediction is meaningful if we can show that the density of  $y_{N+1}$  is concentrated around  $\hat{y}_{N+1}$ . Usually, this means two things: 1) that the means of  $y_{N+1}$  and  $\hat{y}_{N+1}$  are equal, and the variance of their difference is small; 2) that there exists an interval  $[a, b]$ , preferably small, that contains both  $\hat{y}_{N+1}$  and  $y_{N+1}$  with known probability (the higher, the better). Here we will focus on the second one. Let us start without explanatory variables ( $x_i$ ). Goal: given  $y_1, \dots, y_N$ , construct an interval that contains  $y_{N+1}$  with certified probability. If  $y_1, \dots, y_N$  are Gaussian, it can be done with Student’s  $t$  ( $[a, b] = [\hat{\theta} - t_\alpha \bar{s} \sqrt{1 + \frac{1}{N}}, \hat{\theta} + t_\alpha \bar{s} \sqrt{1 + \frac{1}{N}}]$ ). However, here we want a *pivotal* (distribution-free) result.

Let us consider a “chain” of examples, that get closer and closer to the main point.

*Example.* We have a list of  $N + 1$  different numbers. We write each number on a ball, and mark the two balls with the minimum and the maximum number as “extreme point”. Now we put the balls in an urn, shake, and extract one ball. What is the probability that the extracted ball is an extreme? Of course,  $\frac{2}{N+1}$ .

*Example.* We extract *all* the balls in random order. What is the probability that the *last* ball is an extreme? If all the different  $(N + 1)!$  orders have the same probability, the selection of the last ball is equivalent to the extraction of a single ball. Hence the requested probability is again  $\frac{2}{N+1}$ . Put another way, the permutations of the balls that have either the minimum or the maximum at the last position are  $N! + N! = 2N!$ , hence the probability of extracting a permutation with an extreme at the last position is  $\frac{2N!}{(N+1)!} = \frac{2}{N+1}$ .

*Example.* We extract  $N + 1$  independent samples  $y_1, y_2, \dots, y_{N+1}$  of a random variable having density  $f(y)$ . What is the probability  $P$  that  $y_{N+1}$  is either the maximum or the minimum? Denote the event “the last number  $y_{N+1}$  is an extreme” as  $\mathcal{E}$ , and the event “any two numbers among  $y_1, \dots, y_{N+1}$  are equal” as  $\mathcal{N}$ . Due to the hypothesis that the random variables have a density, the probability of  $\mathcal{N}$  is 0, and the probability of  $\mathcal{E}$  conditioned to the extraction of  $N + 1$  different numbers is always the same ( $\frac{2}{N+1}$ ). Hence

$$\begin{aligned} P &= \int_{\mathbb{R}^{N+1}} \mathbb{P}[\mathcal{E} \mid y_1, \dots, y_{N+1}] f(y_1) \cdots f(y_{N+1}) dy_1 \cdots dy_{N+1} \\ &= \int_{\mathbb{R}^{N+1} \setminus \mathcal{N}} \mathbb{P}[\mathcal{E} \mid y_1, \dots, y_{N+1}] f(y_1) \cdots f(y_{N+1}) dy_1 \cdots dy_{N+1} \\ &= \int_{\mathbb{R}^{N+1} \setminus \mathcal{N}} \frac{2}{N+1} f(y_1) \cdots f(y_{N+1}) dy_1 \cdots dy_{N+1} \\ &= \frac{2}{N+1} \int_{\mathbb{R}^{N+1} \setminus \mathcal{N}} f(y_1) \cdots f(y_{N+1}) dy_1 \cdots dy_{N+1} = \frac{2}{N+1}. \end{aligned}$$

*Example.* We extract  $N$  independent samples  $y_1, y_2, \dots, y_N$ , from a random variable with density  $f(y)$ . Let  $a$  and  $b$  be their minimum and maximum. What is the probability that a “future” sample  $y_{N+1}$  falls outside the interval  $[a, b]$ ? The probability is  $\frac{2}{N+1}$ , because  $y_{N+1}$  is outside  $[a, b]$  exactly whenever it is either the maximum or the minimum of the whole sample. We have proven the following

**Lemma 6.1.1** *Let  $y_1, y_2, \dots, y_N$  be independent and identically distributed variables with density  $f(y)$ . Let  $a$  and  $b$  be their minimum and maximum respectively. Then  $[a, b]$  is a prediction interval for a future sample  $y_{N+1}$ , with probability  $1 - \frac{2}{N+1} = \frac{N-1}{N+1}$ .*

## 6.2 Convex problems and Helly's theorem

**Definition 6.2.1** We will call convex problem an optimization problem of the form:

$$\begin{aligned} & \text{minimize } f(\theta) \\ & \text{subject to } f_1(\theta) \leq 0, \\ & \quad \vdots \\ & \quad f_n(\theta) \leq 0, \\ & \quad \theta \in \Theta \subseteq \mathbb{R}^d, \end{aligned}$$

where  $f, f_1, \dots, f_n$  are convex functions and  $\Theta$  is a convex set. Each of the inequalities  $f_i(\theta) \leq 0$  is called a constraint. The problem is said to be feasible if there exists at least a point  $\bar{\theta} \in \Theta$  that satisfies all the constraints, that is, such that  $f_i(\bar{\theta}) \leq 0$  for all  $i = 1, \dots, n$ .

A convex problem may or may not have a solution, and a solution may not be unique. The fundamental fact about convex problems is that if a point is *locally* a minimum, then it is a minimum also *globally*, hence a solution to the problem. Convex problems are considered “easy”, because we have readily available efficient and robust algorithms to solve them numerically. We consider  $(d+1)$ -dimensional problems of this form:

$$\begin{aligned} & \text{minimize } \gamma \\ & \text{subject to } g_1(\theta) - \gamma \leq 0 \\ & \quad \vdots \\ & \quad g_n(\theta) - \gamma \leq 0, \\ & \quad \theta \in \mathbb{R}^d, \gamma \in \mathbb{R}, \end{aligned} \tag{11}$$

where  $g_1, \dots, g_n$  are convex functions. The goal function  $f(\theta, \gamma) = \gamma$  is trivially convex, and if  $g_i(\theta)$  is convex, then the constraint function  $f_i(\theta, \gamma) = g_i(\theta) - \gamma$  is also convex. We assume that all the problems under consideration are feasible and admit a unique solution  $(\theta^*, \gamma^*)$ .

**Definition 6.2.2** Let  $(\theta^*, \gamma^*)$  be the solution of (11), and consider one of its constraints,  $g_i(\theta) - \gamma \leq 0$ . We will call the latter a support constraint if the solution  $(\theta^{**}, \gamma^{**})$  of the problem obtained by removing the constraint is strictly “better” than  $(\theta^*, \gamma^*)$ , meaning that it attains  $\gamma^{**} < \gamma^*$ .

In other words, removing a support constraint the solution “falls”. If a constraint is *not* a support constraint, then it can be removed from the problem without any consequence (the solution is the same). If we add a constraint to a problem, and it happens to become a support constraint in the new problem, then the new  $\gamma^*$  of the solution must *increase* with respect to the old solution (indeed, removing it again  $\gamma^*$  must fall).

**Theorem 6.2.1 (Helly).** Let  $S_1, \dots, S_N$  be convex subsets of  $\mathbb{R}^n$ . If the intersection of any  $n+1$  of these subsets is nonempty, then the intersection of all of them is nonempty.

**Theorem 6.2.2 (Campi, Calafiori, Garatti; “scenario approach”).** Any feasible convex problem like (11) has at most  $d+1$  support constraints.

**Proof.** Let  $(\theta^*, \gamma^*)$  be the solution to the problem, and define

$$\begin{aligned} S_1 &= \{(\theta, \gamma) \in \mathbb{R}^{d+1} \mid g_1(\theta) \leq \gamma\}, \\ &\quad \vdots \\ S_n &= \{(\theta, \gamma) \in \mathbb{R}^{d+1} \mid g_n(\theta) \leq \gamma\}, \\ Z &= \{(\theta, \gamma) \in \mathbb{R}^{d+1} \mid \gamma < \gamma^*\}. \end{aligned}$$

The sets  $S_1, \dots, S_n \subset \mathbb{R}^{d+1}$  are the epigraphs of the convex functions  $g_1(\cdot), \dots, g_n(\cdot)$ , hence they are convex sets; a point belongs to  $S_i$  if and only if it satisfies the  $i$ -th constraint. The set  $Z \subset \mathbb{R}^{d+1}$  is an open half-plane,

hence of course another convex set. Any point belonging to  $Z$  is “sub-optimal” ( $\gamma < \gamma^*$ ); there cannot exist a point that satisfies all the constraints and belongs to  $Z$ , otherwise  $(\theta^*, \gamma^*)$  would not be the minimizing solution. For the sake of contradiction, assume now that the support constraints of the problem are (without loss of generality)  $d + 2$ . Extract an arbitrary collection of  $d + 2$  sets from  $S_1, \dots, S_n, Z$ . If  $Z$  happens to be in the collection, then the collection contains exactly  $d + 1$  (epigraphs of) constraints, of which at most  $d + 1$  are support constraints. Since the support constraints are assumed to be  $d + 2$ , at least one of them has been “removed”, and the solution “falls”, namely there exists a point  $(\theta^{**}, \gamma^{**})$  which satisfies the  $d + 1$  constraints (meaning that it belongs to the  $d + 1$  sets  $\{S_i\}$ ) and attains  $\gamma^{**} < \gamma^*$  (meaning that it belongs to  $Z$ ). On the other hand, if  $Z$  is not in the collection, the latter contains just  $d + 2$  constraints, and since the problem is feasible, there exists at least a point satisfying all of them. Summing up, for any choice of  $d + 2$  sets in  $S_1, \dots, S_n, Z$ , their intersection is non-empty. Applying Helly’s theorem we obtain that

$$S_1 \cap \dots \cap S_n \cap Z \neq \emptyset,$$

but this is clearly in contradiction with the hypothesis that  $(\theta^*, \gamma^*)$  is the solution to the problem. The contradiction stems from the assumption that the support constraints were more than  $d + 1$ , and is enough to establish the claim.  $\square$

### 6.3 Prediction intervals revisited

*Example.* An urn contains  $N + 1$  balls, of which  $d + 1$  are labeled “support constraint” and the others are white. We extract a ball from the urn. What is the probability that it has the label? Of course,  $\frac{d+1}{N+1}$ . Let  $y_1, y_2, \dots, y_N$  a random sample drawn from the density  $f(y)$ , and consider the following convex problem (the constraints are now random):

$$\begin{aligned} & \text{minimize } \gamma \\ & \text{subject to } |\theta - y_1| - \gamma \leq 0 \\ & \quad \vdots \\ & \quad |\theta - y_N| - \gamma \leq 0 \\ & \quad \theta \in \mathbb{R}, \gamma \in \mathbb{R}, \end{aligned}$$

Note that a constraint of the form  $g(\theta, \gamma) = |\theta - y| - \gamma \leq 0$  is convex in the variables  $(\theta, \gamma)$ . Indeed, any set of the form  $|\theta - y| - \gamma \leq 0$  is the intersection of two half-spaces:  $\theta - y - \gamma \leq 0, \theta - y + \gamma \geq 0$ . Hence it is the intersection of two convex sets, which is itself convex. The solution to the problem is the pair  $(\theta^*, \gamma^*)$  attaining the *minimum*  $\gamma$  such that  $\theta$  has distance at most  $\gamma$  from all the points  $\{y_i\}$ . It is such that  $\theta^* - \gamma^* = \min_i \{y_i\}$ ,  $\theta^* + \gamma^* = \max_i \{y_i\}$ . Recall that the support constraints of the problem are at most  $d + 1 = 2$ . It may actually be the case that they are less than 2; But since the random variables have a density, this happens with probability 0. Hence, *almost surely* the support constraints are exactly 2 (the ones with  $\min_i \{y_i\}$  and  $\max_i \{y_i\}$ ). The solution yields the prediction interval that we had obtained before:

$$[a, b] = [\theta^* - \gamma^*, \theta^* + \gamma^*] = \left[ \min_i \{y_i\}, \max_i \{y_i\} \right].$$

Does the solution to the above problem remain the same if another constraint, corresponding to a new measure  $y_{N+1}$  drawn from the same density, is added?

$$\begin{aligned} & \text{minimize } \gamma \\ & \text{subject to } |\theta - y_1| - \gamma \leq 0 \\ & \quad \vdots \\ & \quad |\theta - y_N| - \gamma \leq 0 \\ & \quad |\theta - y_{N+1}| - \gamma \leq 0 \\ & \quad \theta \in \mathbb{R}, \gamma \in \mathbb{R}, \end{aligned}$$

If  $y_{N+1} \in [a, b]$ , then the last constraint does not change anything, and the solutions to the two problems are identical. If  $y_{N+1} \notin [a, b]$ , the last constraint  $|\theta - y_{N+1}| - \gamma \leq 0$  becomes a support constraint for the new problem, and the solution must change.

What is the probability? Let the event  $\mathcal{E} = \{\text{the new constraint } |\theta - y_{N+1}| - \gamma \leq 0 \text{ is a support constraint}\}$ . Since the support constraints are almost surely  $d + 1$ , conditioning to the extraction  $y_1, y_2, \dots, y_N, y_{N+1}$ ,

$$\begin{aligned} P &= \int_{\mathbb{R}^{N+1}} \mathbb{P}[\mathcal{E} | y_1, \dots, y_{N+1}] f(y_1) \cdots f(y_{N+1}) dy_1 \cdots dy_{N+1} \\ &= \int_{\mathbb{R}^{N+1} \setminus \mathcal{N}} \mathbb{P}[\mathcal{E} | y_1, \dots, y_{N+1}] f(y_1) \cdots f(y_{N+1}) dy_1 \cdots dy_{N+1} \\ &= \int_{\mathbb{R}^{N+1} \setminus \mathcal{N}} \frac{d+1}{N+1} f(y_1) \cdots f(y_{N+1}) dy_1 \cdots dy_{N+1} \\ &= \frac{d+1}{N+1} \int_{\mathbb{R}^{N+1} \setminus \mathcal{N}} f(y_1) \cdots f(y_{N+1}) dy_1 \cdots dy_{N+1} = \frac{d+1}{N+1}. \end{aligned}$$

Summing up: we have considered a convex problem with  $N$  constraints, and its solution has yielded the prediction interval that we knew from before. A new measure  $y_{N+1}$  falls outside this prediction interval exactly when, adding to the problem a new constraint corresponding to  $y_{N+1}$ , this becomes a support constraint, in other words it violates the old solution. By virtue of the same reasoning on the ordering of measurements that we have considered in Section 6.1, the probability of violation is  $\frac{d+1}{N+1}$ . Finally, since in this example  $d = \text{dimension of } \theta = 1$ , the violation probability is  $\frac{2}{N+1}$ , as we had found before.

## 6.4 Interval predictor models

Let  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N), (x_{N+1}, y_{N+1})$  be random pairs, independent and identically distributed according to an unknown density, where  $x_i, y_i \in \mathbb{R}$ . Consider the optimization problem with the constraints corresponding to the first  $N$  measures:

$$\begin{aligned} &\text{minimize } \gamma \\ &\text{subject to } |\theta_1 + \theta_2 x_1 - y_1| - \gamma \leq 0 \\ &\quad \vdots \\ &\quad |\theta_1 + \theta_2 x_N - y_N| - \gamma \leq 0 \\ &\quad (\theta_1, \theta_2) \in \mathbb{R}^2, \gamma \in \mathbb{R}. \end{aligned}$$

A constraint like  $g(\theta_1, \theta_2, \gamma) = |\theta_1 + \theta_2 x - y| - \gamma \leq 0$  is convex in the variables  $(\theta_1, \theta_2, \gamma)$ , because any set of the form  $|\theta_1 + \theta_2 x - y| - \gamma \leq 0$  is the intersection of two half-spaces, as before. The solution to the problem is a certain triple  $(\theta_1^*, \theta_2^*, \gamma^*)$ . The parameters  $\theta_1^*, \theta_2^*$  yield a linear model  $y = \theta_1^* + \theta_2^* x$  which is “closest to the data  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ ” in the sense that it attains  $\min \max_{i=1, \dots, N} |\theta_1^* + \theta_2^* x_i - y_i| = \gamma^*$ . It is of course *not* the least squares model (the LS model attains  $\min \sum_{i=1}^N (\theta_1^* + \theta_2^* x_i - y_i)^2$  instead), although similar for well-behaved data. As the “next” explanatory variable  $x_{N+1}$  comes, we can predict the value of  $y_{N+1}$ :

$$\hat{y}_{N+1} = \theta_1^* + \theta_2^* x_{N+1}$$

and append a “certificate” to this prediction: it will hold  $|\hat{y}_{N+1} - y_{N+1}| = |\theta_1^* + \theta_2^* x_{N+1} - y_{N+1}| > \gamma^*$  with the same probability with which the new constraint  $|\theta_1 + \theta_2 x_{N+1} - y_{N+1}| - \gamma \leq 0$  would violate the solution  $(\theta_1^*, \theta_2^*, \gamma^*)$  if added to the problem. Since here  $d = \text{dimension of } \theta = 2$ , the maximum number of support constraints in this problem is  $d + 1 = 3$ . We conclude that, given  $x_{N+1}$ , the interval

$$[a, b] = [\theta_1^* + \theta_2^* x_{N+1} - \gamma^*, \theta_1^* + \theta_2^* x_{N+1} + \gamma^*]$$

is a prediction interval for  $y_{N+1}$  with probability  $1 - \frac{3}{N+1}$ .

Generalization:

$$\begin{aligned}
 & \text{minimize } \gamma \\
 & \text{subject to } |\varphi(x_1)^\top \theta - y_1| - \gamma \leq 0 \\
 & \quad \vdots \\
 & \quad |\varphi(x_N)^\top \theta - y_N| - \gamma \leq 0 \\
 & \theta \in \mathbb{R}^d, \gamma \in \mathbb{R};
 \end{aligned}$$

The support constraints are  $\leq d + 1$ . If the data  $(x_i, y_i)$  are distributed according to a density, then in general the support constraint are *almost surely*  $d + 1$ , hence the violation probability is *exactly*  $\frac{d+1}{N+1}$ . On the other hand, if the data are not distributed according to a density the support constraints may be *less* than  $d + 1$ , and the violation probability satisfies an inequality in the “good” direction:  $P[|\varphi(x_{N+1})^\top \theta^* - y_{N+1}| > \gamma^*] \leq \frac{d+1}{N+1}$ . The following:

$$[a, b] = [\varphi(x_{N+1})^\top \theta^* - \gamma^*, \varphi(x_{N+1})^\top \theta^* + \gamma^*]$$

is a prediction interval for  $y_{N+1}$  with probability *at least*  $1 - \frac{d+1}{N+1}$ .

## 7 Essential concepts to recall from maths

### 7.1 Linear algebra

Prototypical scalar product in  $\mathbb{R}^p$ :  $\langle x, y \rangle = x^\top y$ .

Corresponding norm (Euclidean norm):  $\|x\|_2 = \sqrt{\langle x, x \rangle}$ .

Two vectors  $v, w \in \mathbb{R}^p$  are *orthogonal* if  $\langle v, w \rangle = 0$ . This is denoted  $v \perp w$ .

Given a subset  $S \subset \mathbb{R}^p$ , the *orthogonal complement* of  $S$  in  $\mathbb{R}^p$  is the set  $S^\perp = \{v \in \mathbb{R}^p \mid v \perp w \text{ for all } w \in S\}$ .

Let  $V$  be a subspace of  $\mathbb{R}^p$ . Then  $(V^\perp)^\perp = V$ .

The *range* of a matrix  $A \in \mathbb{R}^{n \times p}$  is the set  $\text{range } A = \{v \in \mathbb{R}^n \mid \text{there exists } w \in \mathbb{R}^p \text{ such that } v = Aw\}$ .

Note that  $\text{range } A = \text{span } \{\text{columns of } A\}$ .

The *null space* of a matrix  $A \in \mathbb{R}^{n \times p}$  is the set  $\text{null } A = \{v \in \mathbb{R}^p \mid Av = 0\}$ .

Let  $A \in \mathbb{R}^{n \times p}$ , understood as a linear mapping from  $\mathbb{R}^p$  to  $\mathbb{R}^n$ . Then 1)  $\text{range } A = (\text{null } A^\top)^\perp$ ; 2)  $\text{null } A = (\text{range } A^\top)^\perp$ ; 3)  $\text{range } A^\top = (\text{null } A)^\perp$ ; 4)  $\text{null } A^\top = (\text{range } A)^\perp$ .

Corollary: for any matrix  $A \in \mathbb{R}^{m \times n}$ ,  $\text{range } A = \text{range } AA^\top$ .

Let  $A \in \mathbb{R}^{n \times p}$ . The *rank* of  $A$  is the maximum dimension of a square matrix, obtained from  $A$  by suppressing some rows and/or columns, with nonzero determinant.

The rank of  $A$  is equal to: 1) the dimension of the subspace of  $\mathbb{R}^n$  generated by its columns; 2) the dimension of the subspace of  $\mathbb{R}^p$  generated by its rows;

If  $A \in \mathbb{R}^{n \times p}$ , where  $n \geq p$  (“tall” matrix, i.e. more rows than columns), we say that  $A$  has *full rank* if  $\text{rank } A = p = \text{number of columns}$ . Then the columns of  $A$  are linearly independent, and the subspace of  $\mathbb{R}^n$  generated by them has dimension  $p$  (the maximum possible).

Conversely, if  $n \leq p$  (“flat” matrix, i.e. more columns than rows), we say that  $A$  has *full rank* if  $\text{rank } A = n = \text{number of rows}$ . Then the *rows* of  $A$  are linearly independent, and their span has dimension  $n$ .

In particular, if  $A \in \mathbb{R}^{p \times p}$  (square), the following statements are equivalent:

- $A$  has full rank ( $= p$ );
- the columns of  $A$  are linearly independent, and form a basis of  $\mathbb{R}^p$ ;
- the columns of  $A^\top$  (transposes of the rows of  $A$ ) are linearly independent, and form a basis of  $\mathbb{R}^p$ ;
- $A$  is invertible, i.e. non-singular, its determinant is nonzero, etc.

## 7.2 Probability

**Lemma 7.2.1 (Markov inequality).** Let  $X$  be a nonnegative random variable (that is,  $X \geq 0$  almost surely). Then for all  $\varepsilon > 0$

$$P[X \geq \varepsilon] \leq \frac{E[X]}{\varepsilon}.$$

**Proof.** Consider the function  $\mathbb{1}_{[\varepsilon, \infty)}(x)$  that takes the values 1 when  $x \geq \varepsilon$ , and 0 otherwise. Then, for all  $x \geq 0$ ,  $\mathbb{1}_{[\varepsilon, \infty)}(x) \leq x/\varepsilon$  (verify this!). Consequently,

$$P[X \geq \varepsilon] = E[\mathbb{1}_{[\varepsilon, \infty)}(X)] \leq E[X/\varepsilon] = \frac{E[X]}{\varepsilon}.$$

□

**Lemma 7.2.2 (Čebyšev inequality).** Let  $X$  be any real random variable with mean  $\mu$  and variance  $\sigma^2$ . Then

$$P[|X - \mu| \geq \varepsilon] \leq \frac{\sigma^2}{\varepsilon^2}.$$

**Proof.** Let us apply the Markov inequality to the nonnegative random variable  $(X - \mu)^2$ :

$$P[|X - \mu| \geq \varepsilon] = P[(X - \mu)^2 \geq \varepsilon^2] \leq \frac{E[(X - \mu)^2]}{\varepsilon^2} = \frac{\sigma^2}{\varepsilon^2}.$$

□

Let  $\{X_i\}_{i=1}^\infty$  be random variables.

**Definition 7.2.1**  $X_i$  is said to converge to  $X$  in probability if for all  $\varepsilon > 0$

$$\lim_{i \rightarrow \infty} P[|X_i - X| \geq \varepsilon] = 0.$$

**Definition 7.2.2**  $X_i$  is said to converge to  $X$  almost surely if the event

$$\left\{ \lim_{i \rightarrow \infty} X_i = X \right\} = \left\{ \omega \mid \lim_{i \rightarrow \infty} X_i(\omega) = X(\omega) \right\}$$

has probability 1. In other terms, the set of trajectories  $X_1(\omega), X_2(\omega), X_3(\omega), \dots$  which do not converge to  $X$  has probability 0.

Among these, almost sure convergence is the strongest form of convergence, and the most desirable; convergence in probability is an intermediate property, and convergence in distribution is the weakest one. Here, “strong” and “weak” are meant in the sense that the stronger property implies the weaker:

$X_i$  converges almost surely to  $X$

↓

$X_i$  converges in probability to  $X$

**Definition 7.2.3**  $X_i$  is said to converge to  $X$  in the mean-square if

$$\lim_{i \rightarrow \infty} E[|X_i - X|^2] = 0.$$

Mean-square convergence is also a kind-of-strong form of convergence, in the sense that

$X_i$  converges in mean-square to  $X$

↓

$X_i$  converges in probability to  $X$

However, no implication exists between almost sure convergence and convergence in the mean-square.

**Lemma 7.2.3** (*Weak law of large numbers, Čebyšev*). Let  $\{X_i\}_{i=1}^{\infty}$  be independent random variables such that for all  $i$

$$\begin{aligned}\mathbb{E}[X_i] &= \mu; \\ \text{Var}[X_i] &= \mathbb{E}[(X_i - \mu)^2] = \sigma^2;\end{aligned}$$

then

$$\frac{1}{N} \sum_{i=1}^N X_i \rightarrow \mu \quad \text{in probability.}$$

**Proof.** We have

$$\begin{aligned}\mathbb{E}\left[\frac{1}{N} \sum_{i=1}^N X_i\right] &= \frac{1}{N} \sum_{i=1}^N \mathbb{E}[X_i] = \frac{1}{N} \sum_{i=1}^N \mu = \mu; \\ \text{Var}\left[\frac{1}{N} \sum_{i=1}^N X_i\right] &= \frac{1}{N^2} \text{Var}\left[\sum_{i=1}^N X_i\right] = \frac{1}{N^2} \sum_{i=1}^N \text{Var}[X_i] = \frac{1}{N^2} N \sigma^2 = \frac{\sigma^2}{N}.\end{aligned}$$

Therefore for all  $\varepsilon > 0$ , by Čebyšev's inequality,

$$\mathbb{P}\left[\left|\frac{1}{N} \sum_{i=1}^N X_i - \mu\right| \geq \varepsilon\right] \leq \frac{\sigma^2/N}{\varepsilon^2} \rightarrow 0$$

□

**Theorem 7.2.1** (*Strong law of large numbers, Kolmogorov*). Let  $\{X_i\}_{i=1}^{\infty}$  be independent and identically distributed random variables with mean  $\mathbb{E}[X_i] = \mu$ . Then

$$\frac{1}{N} \sum_{i=1}^N X_i \rightarrow \mu \quad \text{almost surely.}$$

**Theorem 7.2.2** (*Strong law of large numbers, also by Kolmogorov*). Let  $\{X_i\}_{i=1}^{\infty}$  be independent random variables with arbitrary distributions but having the same mean  $\mathbb{E}[X_i] = \mu$ . If

$$\sum_{i=1}^{\infty} \frac{1}{i^2} \mathbb{E}[X_i^2] < \infty,$$

then

$$\frac{1}{N} \sum_{i=1}^N X_i \rightarrow \mu \quad \text{almost surely.}$$

**Definition 7.2.4**  $f(X_1, \dots, X_N)$  is called an unbiased estimator of  $\theta$  if

$$\mathbb{E}[f(X_1, \dots, X_N)] = \theta.$$

**Definition 7.2.5**  $f(X_1, \dots, X_N)$  is called a consistent estimator of  $\theta$  if

$$\lim_{N \rightarrow \infty} f(X_1, \dots, X_N) = \theta \quad \text{almost surely.}$$

*Example.* The sample average  $\bar{X} = \mathbb{M}[X] = \frac{1}{N} \sum_{i=1}^N X_i$  is a good estimator of the mean  $\mathbb{E}[X_i] = \mu$ , being

- unbiased: indeed

$$\mathbb{E}[\bar{X}] = \mathbb{E}\left[\frac{1}{N} \sum_{i=1}^N X_i\right] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[X_i] = \frac{1}{N} \sum_{i=1}^N \mu = \mu;$$

- consistent: this is precisely the statement of Theorem 7.2.1.

As another example, the sample variance  $s^2 = \frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})^2$  is a consistent estimator of the variance  $\sigma^2 = \mathbb{E}[(X_i - \mu)^2]$ , but not an unbiased one: It is not difficult to show that  $\mathbb{E}[s^2] = \frac{N-1}{N} \sigma^2$ . Hence, usually the estimator  $\tilde{s}^2 = \frac{N}{N-1} s^2 = \frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})^2$ , which enjoys both properties, is preferred. □

### 7.3 Analysis

**Definition 7.3.1** We say that  $f_n$  converges to  $f$  uniformly if

$$\lim_{n \rightarrow \infty} \sup_{x \in C} |f_n(x) - f(x)| = 0$$

This means that for all  $\varepsilon > 0$  there exists  $N \geq 0$  such that, for all  $n \geq N$ ,  $|f_n(x) - f(x)| \leq \varepsilon$  over all the domain  $C$  (in words, the  $f_n$  are all “uniformly  $\varepsilon$ -close” to  $f$  after a certain index  $N$ ).

**Definition 7.3.2** A subset  $C$  of a metric space  $X$  is called compact if from any sequence of points  $x_n \subseteq C$  it is possible to extract a sub-sequence that converges to a point belonging to  $C$ .

**Theorem 7.3.1 (Heine/Borel).** A subset  $C$  of a finite-dimensional vector space is compact if and only if it is closed and bounded.

**Definition 7.3.3** A subset  $S$  of a vector space is called convex if, whenever the points  $x, y$  belong to  $S$ , the point  $z_\lambda = \lambda x + (1 - \lambda)y$  also belongs to  $S$  for all  $\lambda \in [0, 1]$ . (Any such  $z_\lambda$  is called a convex combination of  $x$  and  $y$ .)

Examples: subspaces and their translations, hyperplanes etc.; closed and open balls (i.e. sets of the form  $\{x \in \mathbb{R}^p \mid \|x - c\| \leq r\}$  or  $\{x \in \mathbb{R}^p \mid \|x - c\| < r\}$ ).

**Lemma 7.3.1** An arbitrary intersection of convex sets is itself convex.

**Definition 7.3.4** A function  $f : S \rightarrow \mathbb{R}$  defined on a convex set  $S$  is called convex if, for all  $x, y \in S$  and all  $\lambda \in [0, 1]$ , it holds

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y).$$

Examples: linear and affine functions ( $f(x) = c^\top x + d$ ); norms.

**Definition 7.3.5** The epigraph of a function  $f : S \rightarrow \mathbb{R}$  is the subset of  $S \times \mathbb{R}$  defined as follows:

$$\text{Epi } f = \{(x, y) \in S \times \mathbb{R} \mid f(x) \leq y\}.$$

**Lemma 7.3.2** A function  $f : S \rightarrow \mathbb{R}$  is convex if and only if  $\text{Epi } f$  is a convex set.

**Definition 7.3.6** The  $k$ -sublevel set of a function  $f : S \rightarrow \mathbb{R}$ , where  $k \in \mathbb{R}$ , is the subset of  $S$  defined as follows:

$$S_k = \{x \in S \mid f(x) \leq k\}.$$

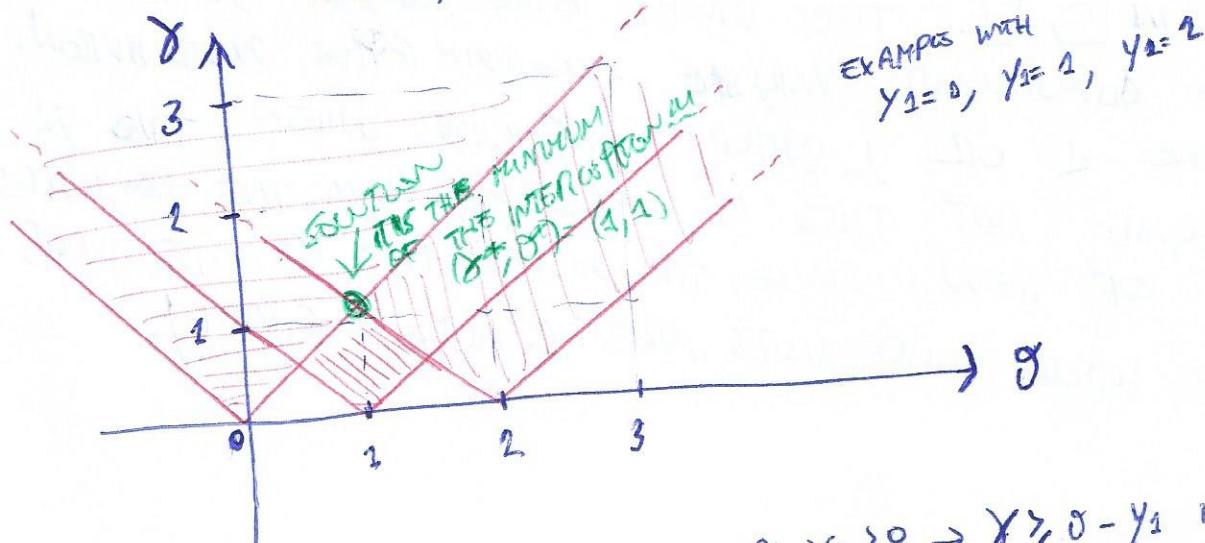
**Lemma 7.3.3** The sublevel sets of a convex function  $f : S \rightarrow \mathbb{R}$  are convex sets.

MINIMIZE  $\gamma$

SUBJECT TO  $|\delta - y_1| - \gamma \leq 0$

$$|\delta - y_2| - \gamma \leq 0$$

$\gamma \in \mathbb{R}$ ,  $\delta \in \mathbb{R}$

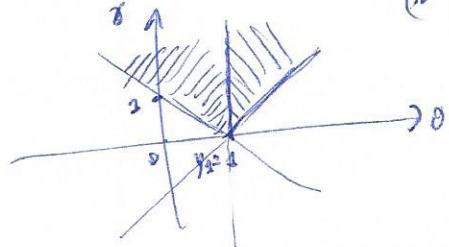


$$|\delta - y_1| \leq \gamma \rightarrow \delta - y_1 \leq \gamma$$

$$-(\delta - y_1) \leq \gamma$$

$$\text{IF } \delta - y_1 \geq 0 \rightarrow \gamma \geq \delta - y_1 \text{ IF } \delta \geq y_1$$

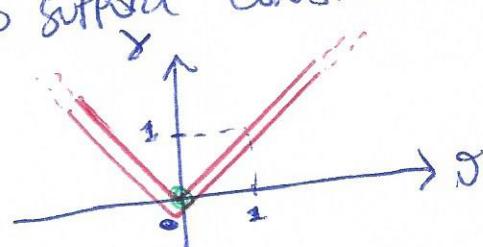
$$\text{IF } \delta - y_1 < 0 \rightarrow \gamma \geq -\delta + y_1 \text{ IF } \delta < y_1$$



A) NO SUPPORT CONSTRAINTS: YES

FOR EXAMPLE TWO  $y = 0$

IF I REMOVE ONE CONSTRAINT  
THE SOLUTION DOESN'T CHANGE



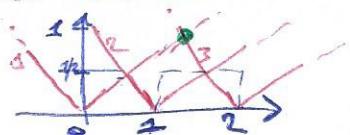
B) ONE SUPPORT CONSTRAINT: YES

IF I REMOVE 1 THE SOLUTION GETS BETTER  
( $\gamma^* = 0$  INSTEAD OF  $\gamma^* = 1$ ) SO IT'S A SUPPORT CONSTRAINT

IF I REMOVE 2 THE SOLUTION IS THE SAME  
SO IT'S NOT A SUPPORT CONSTRAINT

IF I REMOVE 3 THE SOLUTION IS THE SAME  
SO IT'S NOT A SUPPORT CONSTRAINT

C) TWO SUPPORT CONSTRAINTS: YES



IF I REMOVE 1 THE SOLUTION GETS BETTER, IT'S A SUPPORT CONSTRAINT

IF I REMOVE 2 THE SOLUTION DOESN'T CHANGE

IF I REMOVE 3 THE SOLUTION GETS BETTER, IT'S A SUPPORT CONSTRAINT

D)  $\exists$  SUPPORT CONSTRAINTS IS IMPOSSIBLE DUE TO THE  
CARPI, CALIFORNIA, GARFATU THEOREM. THE SUPPORT  
CONSTRAINTS ARE AT MOST  $d+1$  AND  $d+1 \leq 1 \leq 2$ .

IF  $y_1, \dots, y_n$  ARE RANDOM I.I.D. CONTINUOUS WITH UNIFORM  
DENSITY IN  $[0, 1]$ . THE VALUES ARE RANDOM SO THE POSITIONS  
OF THE CONSTRAINTS VARY ACCORDING TO THE DISTRIBUTION.  
WE HAVE 1 OR 2 SUPPORT CONSTRAINTS WHEN TWO  $y_i$   
ARE EQUAL BUT THIS IS IMPOSSIBLE FROM THE PROBABILISTIC  
POINT OF VIEW SINCE THE DISTRIBUTION IS CONTINUOUS.  
SO THE SUPPORT CONSTRAINTS ARE 2 ALMOST SURELY.

TRUE FUNCTION  $f(\cdot, \theta^*)$

APPROXIMATE FUNCTION  $f(\cdot, \hat{\theta})$

GENERIC FUNCTION  $f(\cdot, \theta)$   
(~~A CERTAIN FUNCTION~~)

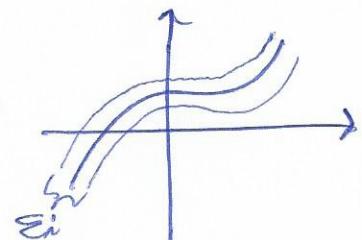
$$y_i = f(x_i, \theta^*) + \epsilon_i$$

$$y_i = f(x_i, \hat{\theta}) + \epsilon_i(\hat{\theta})$$

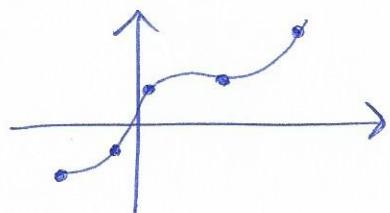
$$y_i = f(x_i, \theta) + \epsilon_i(\theta)$$

BE CAREFUL:

$\epsilon_i$  ARE DISTURBANCES DUE TO THE RANDOMNESS



$\epsilon_i(\theta)$  ARE ERRORS WE COMMIT SELECTING A FUNCTION  
EXPLOITING THE ONLY DATA AT HAND



$$\hat{\theta}_{LS} = \arg \min_{\theta \in \mathbb{R}^P} \sum_{n=1}^N (y_n - \varphi_n^\top \theta)^2$$

$$\frac{\partial}{\partial \theta} \sum_{n=1}^N (y_n - \varphi_n^\top \theta)^2 = 2 \sum_{n=1}^N (y_n - \varphi_n^\top \theta) (-\varphi_n^\top) = 0$$

$$-\sum_{n=1}^N y_n \varphi_n^\top + \sum_{n=1}^N \underbrace{\varphi_n^\top \theta \varphi_n^\top}_\text{transpose this} = 0$$

$$\sum_{n=1}^N \theta^\top \varphi_n \varphi_n^\top = \sum_{n=1}^N y_n \varphi_n^\top$$

$$\begin{matrix} y_n & \square \\ \theta & \square \\ \varphi_n & \square \end{matrix}$$

$$R = \sum_{n=1}^N \varphi_n \varphi_n^\top \in \mathbb{R}^{P \times P}$$

$$\theta^\top \left( \sum_{n=1}^N \varphi_n \varphi_n^\top \right) = \sum_{n=1}^N y_n \varphi_n^\top \quad \text{TRANSPOSING EVERYTHING}$$

$$\left( \sum_{n=1}^N \varphi_n \varphi_n^\top \right) \theta = \sum_{n=1}^N y_n \varphi_n^\top \quad \text{NORMAL EQUATION}$$

ANY  $\theta$  THAT SOLVES THE NORMAL EQUATION IS A SOLUTION FOR THE LEAST SQUARES PROBLEM.

IF  $R = \sum_{n=1}^N \varphi_n \varphi_n^\top$  IS INVERTIBLE THEN THE ONLY SOLUTION TO THE LS. PROBLEM IS:

$$\theta_{LS} = \left( \sum_{n=1}^N \varphi_n \varphi_n^\top \right)^{-1} \sum_{n=1}^N \varphi_n y_n$$

DEF: Let  $v_1, \dots, v_m \in V$

$W = \text{span} \{v_1, \dots, v_m\} = \{w \in V \mid w = a_1v_1 + \dots + a_mv_m \text{ FOR SOME COEFFICIENTS } a_1, \dots, a_m\}$

1)  $W$  IS A SUBSPACE OF  $V$

2) BY DEFINITION  $v_1, \dots, v_m$  GENERATE  $\text{span} \{v_1, \dots, v_m\}$

$L: V \rightarrow W$  IS A LINEAR MAPPING IF  $\forall x, y \in V$

$$L(\alpha x + \beta y) = \alpha L(x) + \beta L(y)$$

LINEAR MAPPINGS FROM  $\mathbb{R}^m$  TO  $\mathbb{R}^n$  ARE REPRESENTED BY MATRICES IN  $\mathbb{R}^{n \times m}$

LET  $A \in \mathbb{R}^{m \times n}$  A LINEAR MAP FROM  $\mathbb{R}^n$  TO  $\mathbb{R}^m$   
(DOMAIN) (CO-DOMAIN)

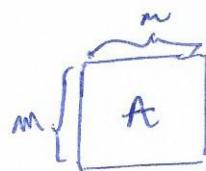
RANGE  $A = \{v \in \mathbb{R}^m \mid v = Aw \text{ for some } w \in \mathbb{R}^n\} \leftarrow$  SUBSPACE OF CO-DOMAIN  
(IMAGE OF  $A$ )

NUL  $A = \{v \in \mathbb{R}^m \mid Av = 0\} \leftarrow$  SUBSPACE OF DOMAIN  
(KERNEL OF  $A$ )

$$A \in \mathbb{R}^{m \times n}$$

$$A: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

HAVE TWO SUBSPACES:



RANGE  $A: \{v \in \mathbb{R}^m \mid \exists w \in \mathbb{R}^n, Aw = v\} \subseteq \mathbb{R}^m$  (a-DOMAIN OF  $A$ )

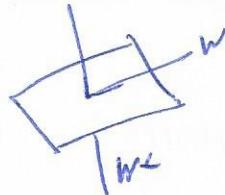
NUL  $A: \{w \in \mathbb{R}^n \mid Aw = 0\} \subseteq \mathbb{R}^n$  (a-DOMAIN OF  $A$ )

LET  $V$  BE A REAL VECTOR SPACE AND  $S \subseteq V$  THEN THE ORTHOGONAL COMPLEMENT OF  $S$  IS DEFINED AS follows:

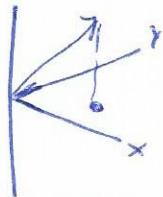
$$S^\perp := \{v \in V \mid v \perp s \text{ for all } s \in S\}$$

FACT : LET  $W$  BE A SUBSPACE OF  $V$   $\forall v \in V$  CAN BE WRITTEN  
UNIQUELY AS :  $v = w + w^\perp$   
WHERE  $w \in W$  AND  $w^\perp \in W^\perp$

FACT 1 IF  $V$  IS FINITE DIMENSIONAL ( $\mathbb{R}^n$ ) AND  $W$  IS A SUBSPACE OF  $V$   
THEN  $(W^\perp)^\perp = W$



EXAMPLE : A PROJECTION ON THE XY PLANE



$$A = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\text{RANGE } A = \text{SPAN} \left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \right\}$$

$$\dim \text{RANGE } A = 2$$

$$\text{NULL } A = \text{SPAN} \left\{ \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \right\} = \left\{ \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \right\}$$

$R$  symmetric  $\Rightarrow$  range  $R = (\text{null } R^T)^\perp = (\text{null } R)^\perp$

1)  $\forall v \in \text{null } R \rightarrow Rv = 0$

$$v^T R v = v^T \sum_{i=1}^N q_i q_i^T v = 0$$

$$\sum_{i=1}^N (q_i^T v)(v^T q_i) = 0 \rightarrow \sum_{i=1}^N (q_i^T v)^2 = 0$$

$$q_i^T v = 0 \quad \forall i = 1, \dots, N \quad v \perp q_i \text{ thi}$$

$$v \perp \alpha_1 q_1 + \alpha_2 q_2 + \dots + \alpha_N q_N \quad \forall \alpha_1, \dots, \alpha_N \in \mathbb{R}$$

$$v \perp \text{span}\{q_1, \dots, q_N\}$$

$$v \in \text{span}\{q_1, \dots, q_N\}^\perp$$

2)  $v \in \text{span}\{q_1, \dots, q_N\}^\perp$

$$v \perp \alpha_1 q_1 + \dots + \alpha_N q_N \quad \forall \alpha_1, \dots, \alpha_N \in \mathbb{R}$$

$$v \perp q_1, v \perp q_2, \dots, v \perp q_i \quad \forall i = 1, \dots, N$$

$$q_i^T v = 0 \quad \forall i \Rightarrow \left( \sum_{i=1}^N q_i q_i^T \right) v = 0 \Rightarrow Rv = 0 \Rightarrow v \in \text{null } R$$

## NORMAL EQUATIONS

$$R\theta = \sum_{i=1}^N \varphi_i y_i$$

RANGE  $R = \text{SPAN}\{\varphi_i\}$   
 $\forall \tau \in \text{SPAN}\{\varphi_i\} \rightarrow \tau \in \text{RANGE } R \rightarrow \exists \theta \in \mathbb{R}^p \text{ such that } \tau = R\theta$

COROLLARY: THE NORMAL EQUATIONS HAVE AT LEAST ONE SOLUTION

SYSTEM OF LINEAR EQUATIONS

$$Ax = b$$

$$R\theta = \sum_{i=1}^N \varphi_i y_i$$

EXACTLY ONE SOLUTION  $\Leftrightarrow R$  IS INVERTIBLE (AND THE SOLUTION IS OF COURSE  $\theta_{LS} = R^{-1}b$ )  
 INFINITELY MANY SOLUTIONS  $\Leftrightarrow R$  IS NOT INVERTIBLE

$R$  INVERTIBLE

$\Updownarrow$   
 $R$  HAS FULL RANK

$\Updownarrow$   
 RANGE  $R = \mathbb{R}^p$

$$\boxed{\text{SPAN}\{\varphi_1, \dots, \varphi_N\} = \mathbb{R}^p}$$

$$\varphi_1, \varphi_2, \dots, \varphi_{N-1}, \varphi_N \in \mathbb{R}^p$$

$p$  LINEARLY INDEPENDENT VECTORS IN  $\mathbb{R}^p$  FORM A BASIS OF  $\mathbb{R}^p$

$p$  VECTORS IN  $\mathbb{R}^p$  THAT GENERATE  $\mathbb{R}^p$  ARE LINEARLY INDEPENDENT

$\varphi_1, \dots, \varphi_N$  DO NOT GENERATE  $\mathbb{R}^p$ : WHY?

1)  $N < p$  (PATHELOGICAL)

2) THE REGRESSORS ARE FLAWED ( $\varphi_1 = x_1, \varphi_2 = x_2 + x_3, \varphi_3 = x_1 + x_2 + x_3$ )

3)  $x_i$  DO NOT CONTAIN ENOUGH INFORMATION

## MEASUREMENT MODEL

$$y_1 = \varphi_1^\top \theta + \varepsilon_1$$

$$y_2 = \varphi_2^\top \theta + \varepsilon_2$$

$$\vdots$$

$$y_N = \varphi_N^\top \theta + \varepsilon_N$$

$$\square = \square \square + \square$$

$$\square = \square \square + \square$$

}

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

$$\theta = \begin{bmatrix} \varphi_1^\top \\ \varphi_2^\top \\ \vdots \\ \varphi_N^\top \end{bmatrix}$$

$$\varepsilon \in \mathbb{R}^p$$

$$\begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_N \end{bmatrix} \in \mathbb{R}^p$$

$$Y = \theta^\top \varphi + \varepsilon$$

$$\square = \square \square + \square$$

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N (y_i - \theta)^2 \quad \hat{\theta}_{LS} = \frac{1}{N} \sum_{i=1}^N y_i = M[y]$$

ATTAINED MINIMUM:  $\frac{1}{N} \sum_{i=1}^N (y_i - M[y])^2 = \text{SAMPLE VARIANCE OF } y$   
 $S.\text{VAR.} \{y_1, \dots, y_N\}$

REGRESSORS OFTEN ARE:  $\varphi_0 = 1, \varphi_1(x), \varphi_2(x), \dots, \varphi_p(x)$

$$\varphi_i = \begin{bmatrix} 1 \\ \varphi_1(x_i) \\ \vdots \\ \varphi_p(x_i) \end{bmatrix}$$

PROBLEM: (\*)  $\min_{\theta \in \mathbb{R}^p} \frac{1}{N} \sum_{i=1}^N (y_i - \varphi_i^T \theta)^2$

IF I CHOOSE  $\theta_{\text{ML}} = \begin{bmatrix} M[y] \\ 0 \\ \vdots \\ 0 \end{bmatrix}$  THEN (\*) = S.VAR[y]

THEREFORE GET:  $RV = \text{RESIDUAL VARIANCE}$

$$= \frac{1}{N} \sum_{i=1}^N (y_i - \varphi_i^T \hat{\theta}_{LS})^2$$

AND IT FOLLOWS

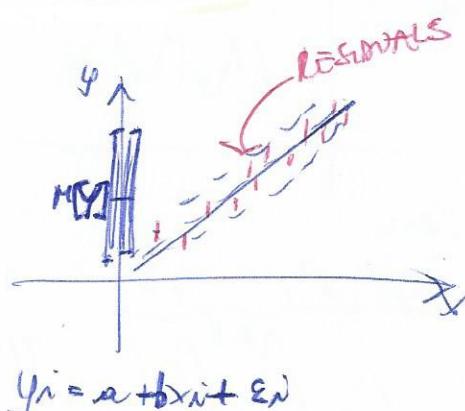
$$RV \leq S.\text{VAR}[y]$$

EV = EXPLAINED VARIANCE := S.VAR[y] - RV

$$0 \leq RV \leq S.\text{VAR}[y]$$

$$0 \leq EV \leq S.\text{VAR}[y]$$

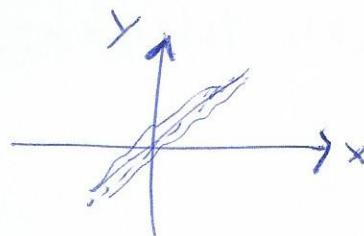
$$S.\text{VAR} = EV + RV$$



Cards I also obtain this picture



TYPICAL CLOUD  
OF TWO INDEPENDENT  
VARIABLES



TYPICAL CLOUD  
OF TWO POSITIVELY  
CORRELATED  
VARIABLES

$$\rho^2 = \frac{E[V]}{S.\text{Var}[y]}$$

$$0 \leq \rho^2 \leq 1$$

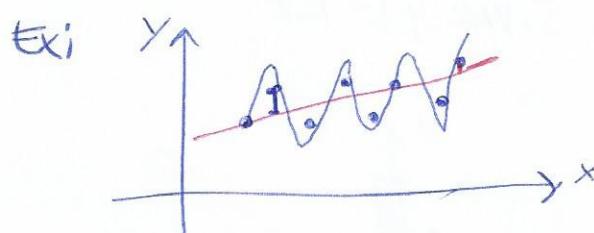
} IT'S A WAY  
TO QUANTIFY  
HOW GOOD THIS MODEL IS

In STATISTICAL MEASURES,  $\rho$  = CORRELATION COEFFICIENT

BETWEEN  $\{x_i\}$  AND  $\{y_i\}$

$$\rho = \frac{\text{cov}\{x_i, y_i\}}{\sqrt{\text{Var}\{x_i\}} \cdot \sqrt{\text{Var}\{y_i\}}} = \frac{\frac{1}{N} \sum_{i=1}^N (x_i - E[x])(y_i - E[y])}{\sqrt{\frac{1}{N} \sum (x_i - E[x])^2} \sqrt{\frac{1}{N} \sum (y_i - E[y])^2}} =$$

\* IF IT'S NEXT TO 1 DATA DESCRIBES THE MODEL IN A GOOD WAY ( $\rho^2 \approx 1$  GOOD!),  
BUT THERE IS AN EXCEPTION:  $N \approx p$



$$\begin{array}{l} N=8 \\ p=2 \\ p=8 \end{array} \quad \begin{array}{l} \text{NUMBER OF MEASUREMENTS} \\ \text{NUMBER OF REGRESSIONS} \end{array}$$

$$\underbrace{1, x_1, x_1^2, \dots, x_1^7}_{B8}$$

IF I USE LS IF FIND THE POLYNOMIAL THAT IS COMPLETELY USELESS

BECAUSE I TRY TO FOLLOW THE NOISES TOO

IF I HAD  $N=100$  AND I USE A 90 ORDER POLYNOMIAL THE SOLUTION IS COMPLETELY USELESS

$$\hat{\theta}_{LS} = \theta^* + \left( \frac{1}{N} \sum_{i=1}^N (\psi_i \psi_i^*) \right)^{-1} \cdot \frac{1}{N} \cdot \sum_{i=1}^N \psi_i \epsilon_i$$

TENTS TO  
AN APPROX  
 $\Sigma$  (SIGMA)

TENTS TO  
 $\infty$

WE EXPECT (HOPES) THAT  $\hat{\theta}_{LS} \rightarrow \theta^*$  AS  $N \rightarrow \infty$

(WE ARE DEALING WITH RANDOM VARIABLES!)

DELICATE WHEN WE  
USE PROCESSES

## PROBABILITY SPACE

$$(\Omega, \mathcal{F}, P)$$

$P: \mathcal{F} \rightarrow [0, 1]$  SUCH THAT

SET CALLED  
SAMPLE  
SPACE

FAMILY  
OF SUBSETS  $\rightarrow$  EVENTS

OF  $\Omega$   
SUCH THAT

$$1) P\left[\bigcup_{i=1}^{\infty} A_i\right] = \sum_{i=1}^{\infty} P[A_i]$$

DISJOINT

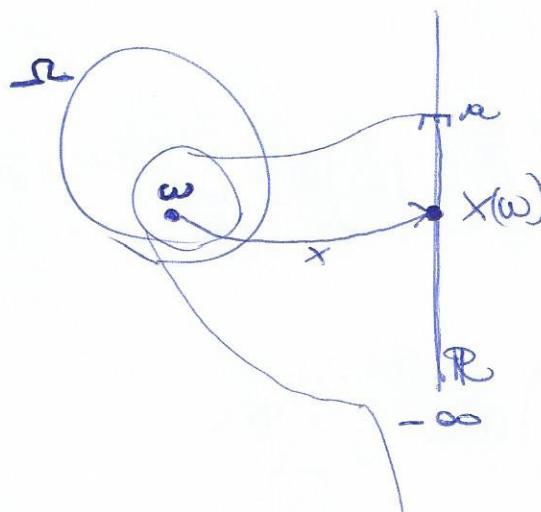
IT'S A GENERAL  
IZATION OF  
AN AREA

$$2) P[\Omega] = 1$$

$$3) A \in \mathcal{F} \Rightarrow A^c \in \mathcal{F}$$

$\left. \begin{array}{l} \text{CLOSE TO RESPECT} \\ \text{OF UNION} \end{array} \right\}$

$\left. \begin{array}{l} \text{CLOSE TO RESPECT} \\ \text{OF COMPLEMENT} \end{array} \right\}$



IF THE FUNCTION IS THAT

$$x^{-1}((-\infty, a]) = \{ \omega | x(\omega) \leq a \}$$

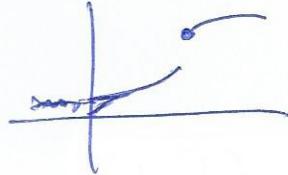
$\dots x(\omega) \in (-\infty, a] \} \in \mathcal{F}$  THEN

THEN  $x$  IS CALLED A  
RANDOM VARIABLE

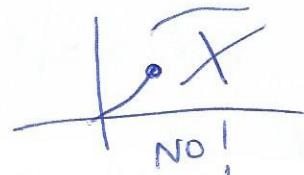
$$\underbrace{P[x^{-1}((-\infty, a])]}_{\in \mathcal{F}} = P\{ \omega | x(\omega) \leq a \} = P[x \leq a] := F(a)$$

$F_x: \mathbb{R} \rightarrow [0, 1]$  DISTRIBUTION OF  $X$

PROPOSITION:  $f(x)$  IS CONTINUOUS FROM THE RIGHT



OK



NO!

- Give  $F(\infty) = 0$   
 $\omega \rightarrow -\infty$

- Give  $F(\infty) = 1$   
 $\omega \rightarrow +\infty$

IF IN PARTICULAR  $f_x$  IS DIFFERENTIABLE (AU R), THEN

$f'_x(a) = \frac{\partial}{\partial a} F_x(a)$  IS CALLED THE DENSITY OF X AND X IS

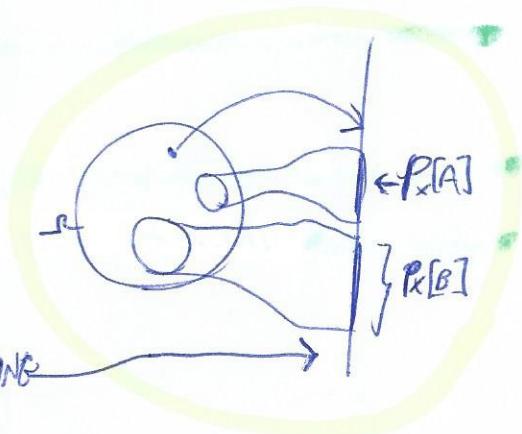
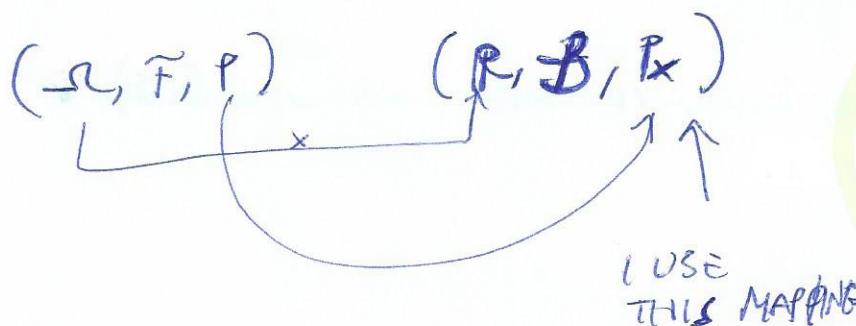
- GAUSS CONTINUOUS VARIABLE

$$F_x(a) = \int_{-\infty}^a f_x(\xi) d\xi$$

$$x^{-1}(A) = \left\{ \omega \in \Omega \mid x(\omega) \in A \right\} \quad x: \Omega \rightarrow \mathbb{R}$$

$\uparrow$   
 $\mathbb{R}$

$$\rho[x^{-1}(A)] = \rho[x \in A] := P_x[A]$$



FOR A CONTINUOUS VARIABLE

$$E[X] := \int_{-\infty}^{+\infty} x \cdot f_x(x) dx = \text{MEAN VALUE}$$

$$\text{VAR}[X] = E[(X - E[X])^2]$$

PROPOSITIONS

$$1) E[\alpha x + bY] = \alpha E[X] + b E[Y]$$

as  $a, b \in \mathbb{R}$ ,  $X, Y = R.V.$  on THE SAME  $\Omega$

$$2) |E[X]| \leq E[|X|]$$

$$3) \text{ IF } X \leq Y \text{ THEN } E[X] \leq E[Y]$$

Def: IF A PROPERTY HOLDS FOR ALL  $w \in \Omega$  EXCEPT FOR  $w \in N$  AND  $P[N] = 0$ , THEN WE SAY THAT THE PROPERTY HOLDS ALMOST SURELY

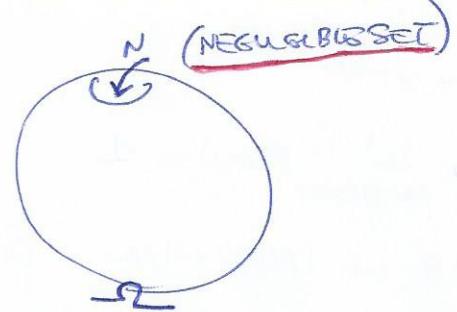
4) IF  $X = Y$  ALMOST SURELY THEN THE

$$E[X] = E[Y]$$

$$X \sim N(0, 1)$$

$$Y = \begin{cases} X & \text{if } X \neq 0 \\ 12 & \text{if } X = 0 \end{cases}$$

$$X = Y \quad \text{A.S.} \Rightarrow E[X] = E[Y]$$



LET  $X, Y$  BE R.V. [ON THE SAME  $(\Omega, \mathcal{F}, P)$ ]

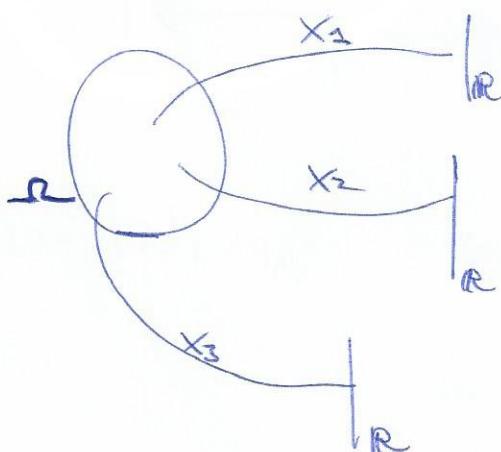
$X, Y$  ARE INDEPENDENT IF

$$P[X \in A, Y \in B] = P[X \in A] \cdot P[Y \in B]$$

JOINT DISTRIBUTION

$$\text{IF } X, Y \text{ ARE INDEPENDENT, THEN } E[XY] = E[X]E[Y]$$

A SEQUENCE  $X_1, X_2, X_3, \dots$  OF R.V. ON THE SAME  $(\Omega, \mathcal{F}, P)$  IS CALLED A RANDOM PROCESS:



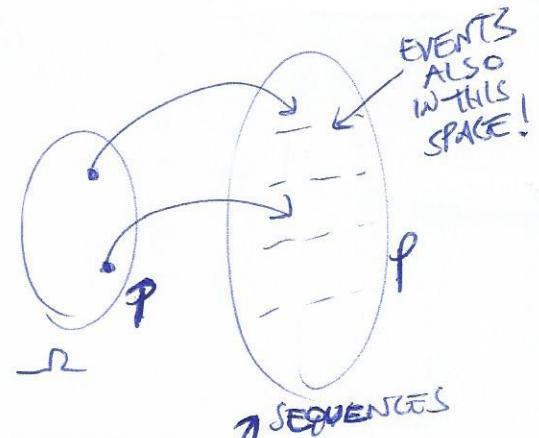
FOR FIXED  $\omega$ ,  $X_i$  IS A RANDOM VARIABLE

FOR FIXED  $\omega \in \Omega$ ,  $(X_1(\omega), X_2(\omega), \dots, X_n(\omega), \dots)$  IS A SEQUENCE OF REAL NUMBERS.  $\omega$  IS SEQUENCE

$(\Omega, \mathcal{F}, P)$

STOCH. PROCESS

((Sequences),  $A$ ,  $P_s$ )



$\{X_i \leq 5\}$  IS AN EVENT ON  $\mathbb{R}$  AXIS

$\{\lim_{i \rightarrow \infty} X_i = x\}$  IS AN EVENT ON THIS SPACE

IF I USE THIS WAY, I'M LOOKING AT THE PROCESS

- CONVERGENCES -

Def: Let  $X_1, X_2, \dots$  BE A STOCHASTIC PROCESS ON  $(\Omega, \mathcal{F}, P)$  AND  
IF  $P[\{\omega \in \Omega \mid \lim_{i \rightarrow \infty} X_i(\omega) = x(\omega)\}] = 1$   $\times$  A R.V. ON  $(\Omega, \mathcal{F}, P)$

THEN  $X_i \rightarrow x$  ALMOST SURELY AS  $i \rightarrow \infty$

Def: If  $\forall \epsilon > 0$   $\lim_{i \rightarrow \infty} P[(|X_i - x|) > \epsilon] = 0$

THEN  $X_i \rightarrow x$  ON PROBABILITY

Def: IF  $\lim_{i \rightarrow \infty} E[(X_i - x)^2] = 0$  THEN  $X_i \rightarrow x$

IN THE MEAN SQUARES

Def: IF  $\lim_{i \rightarrow \infty} F_{X_i}(x) = F_x(x)$  AT ALL THE POINTS AT WHICH  $F$  IS CONTINUOUS THEN  $X_i \rightarrow x$  IN DISTRIBUTION

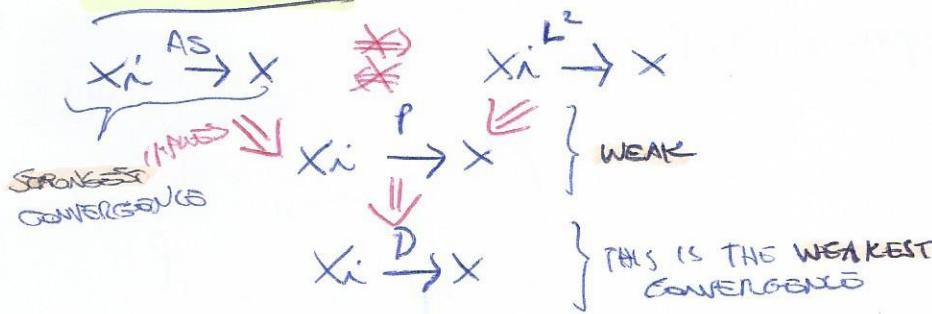
$X_n \xrightarrow{AS} X$  (AS = almost surely)

$X_n \xrightarrow{P} X$  (in probability = P)

$X_n \xrightarrow{L^2} X$  (in mean-squares =  $L^2$ )

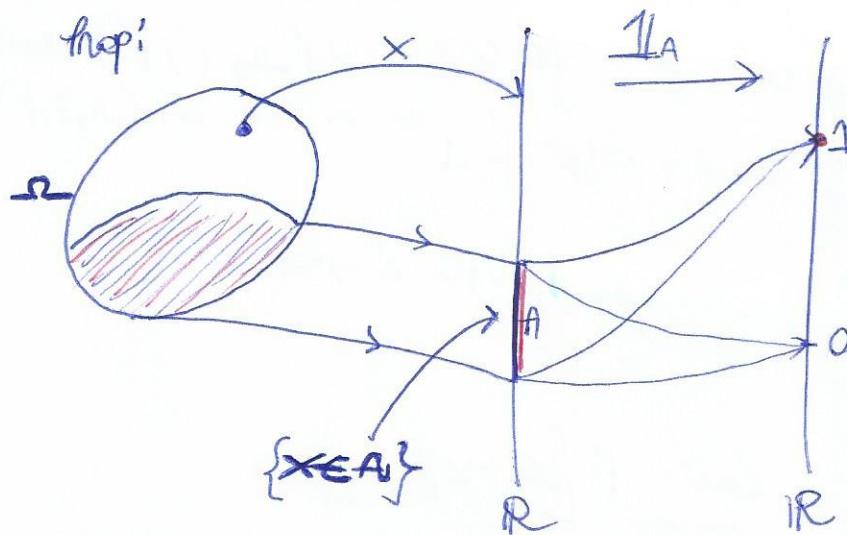
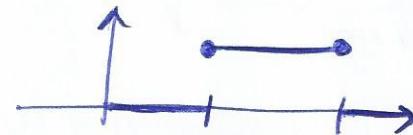
$X_n \xrightarrow{D} X$  (in distribution = D)

### THEOREMS



### INDICATOR FUNCTION of $A \subseteq \mathbb{R}$

$$\mathbb{1}(x) : \begin{cases} 1, & x \in A \\ 0, & \text{otherwise} \end{cases}$$



$$\begin{aligned} P[\mathbb{1}_A(x) = 1] &= P[x \in A] \\ P[\mathbb{1}_A(x) = 0] &= P[x \notin A] \end{aligned}$$

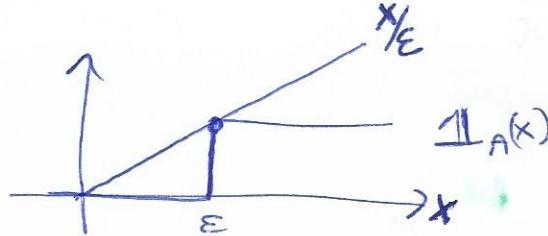
$$E[\mathbb{1}_A(x)] = 1 \cdot P[x \in A] + 0 \cdot P[x \notin A] = P[x \in A]$$

## MARKOV'S INEQUALITY

$$A = \{X \geq \varepsilon\}$$

$x \in \mathbb{R}$

$$\text{For } x \geq 0, \sqrt{\mathbf{1}_A(x)} \leq \frac{x}{\varepsilon}$$



Let  $X$  be a positive random variable  $\geq 0$ .

$$\text{Then } P[X \in A] = P[X \geq \varepsilon]$$

$$E[\mathbf{1}_A(X)] \leq E\left[\frac{X}{\varepsilon}\right] = \frac{E[X]}{\varepsilon}$$

$$\forall \varepsilon > 0, X \text{ rv. } \geq 0 \quad P[X \geq \varepsilon] \leq \frac{E[X]}{\varepsilon}$$

## CHEBYSIEV'S INEQUALITY

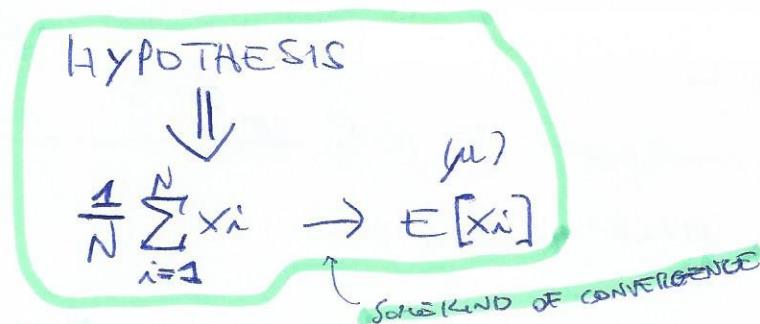
Let  $X$  be a rv. with mean  $\mu$  and variance  $\sigma^2$

$$\text{Then } P[|X - \mu| \geq \varepsilon] \leq \frac{\sigma^2}{\varepsilon^2}$$

Proof:  $(X - \mu)^2 \geq 0 \leftarrow \text{APPLY MARKOV'S INEQUALITY TO THIS}$

$$P[(X - \mu)^2 \geq \varepsilon^2] = P[|X - \mu| \geq \varepsilon] \leq E\left[\frac{(X - \mu)^2}{\varepsilon^2}\right] = \frac{\sigma^2}{\varepsilon^2} \quad \square$$

A theorem is called a LAW OF LARGE NUMBERS if it has the form:



(THE LAW CAN BE WEAK, STRONG ACCORDING TO THE KIND OF CONVERGENCE)

### 1) WEAK LAW OF LARGE NUMBERS

Let  $x_i$  independent with mean  $\mu$  and variance  $\sigma^2$

Then

$$\frac{1}{N} \sum_{i=1}^N x_i \xrightarrow{\text{P}} \mu$$

PROOF:

$$P\left[\left|\frac{1}{N} \sum_{i=1}^N x_i - \mu\right| \geq \epsilon\right] = \text{fix AN } \epsilon > 0$$

$$\text{BT } E\left[\frac{1}{N} \sum x_i\right] = \frac{1}{N} \sum E[x_i] = \mu \text{ so:}$$

$$= P\left[\left|\frac{1}{N} \sum_{i=1}^N (x_i - \mu)\right| \geq \epsilon\right] \leq E\left[\left(\frac{1}{N} \sum_{i=1}^N (x_i - \mu)\right)^2\right] = \frac{1}{N^2} E\left[\left(\sum_{i=1}^N (x_i - \mu)\right)^2\right] =$$

THIS QUANTITY HAS ANGLES ==

$$= \frac{1}{N^2} \sum_{i=1}^N \frac{E[(x_i - \mu)^2]}{\epsilon^2} = \frac{N\sigma^2}{N^2 \epsilon^2} = \frac{\sigma^2}{\epsilon^2} \cdot \frac{1}{N} \xrightarrow{N \rightarrow \infty} 0$$

(IT'S THE DEFINITION OF CONVERGENCE IN PROBABILITY)

### 2) STRONG LLN. [KOLMOGOROV]

Let  $\{x_i\}$  be r.i.d. random variables with mean

$$E[x] = \mu$$

Then  $\frac{1}{N} \sum_{i=1}^N x_i \xrightarrow{\text{A.S.}} \mu$

### 3) STRONG LLN. [KOLMOGOROV]

Let  $\{x_i\}$  be independent R.V. with arbitrary distributions but with the same  $E[x_i] = \mu$ .

If  $\sum_{n=1}^{\infty} \frac{E[x_n^2]}{n^2} < \infty$ , THEN

$$\frac{1}{N} \sum_{i=1}^N x_i \xrightarrow{\text{A.S.}} \mu$$

ADVANTAGE OF 3) IS ARBITRARY DISTRIBUTION

Ex:  $f_x(x) = \frac{c}{1+x^2}$  doesn't have a mean value!!!

$$\int_{-\infty}^{+\infty} \frac{c}{1+x^2} dx \sim [\text{out of}]^{+\infty}$$

[REDACTED]  
↓  
ERGO

### STATISTICAL JARGON

$\theta$  IS A PARAMETER TO BE ESTIMATED

$\hat{\theta}(x_1, \dots, x_N)$  IS A FUNCTION OF THE RANDOM VARIABLES  $x_1, \dots, x_N$   
IF I USE  $\hat{\theta}(\dots)$  TO "ESTIMATE"  $\theta$ , I CALL IT AN ESTIMATOR

OF  $\theta$ .  
UNBIASED ESTIMATOR (GOOD!)

$$E[\hat{\theta}(x_1, \dots, x_N)] = \theta$$

CONSISTENT ESTIMATOR (GOOD!)

IF AS  $N \rightarrow \infty$

$$\hat{\theta}(x_1, \dots, x_N) \xrightarrow{\text{A.S.}} \theta$$

$$y_i = \varphi_i^T \theta^0 + \varepsilon_i$$

random

$$\hat{\theta}_{LS} = \theta^0 + \left( \sum_{i=1}^N \varphi_i \varphi_i^T \right)^{-1} \sum_{i=1}^N \varphi_i \varepsilon_i$$

MUST BE !!! ← WE NEED THIS MATRIX INVERTIBLE  
CLAIMED...

# SINGULAR VALUE DECOMPOSITION (SVD)

For matrices  $A \in \mathbb{R}^{n \times p}$

decomposition:

$$A \geq 0$$

SYMMETRIC

WE HAVE

A BEAUTIFUL

(DECOMPOSIZIONE  
AI VALORI SINGOLARI)  
=  
PARTICOLARE ECONOMIZZANTE  
DATA SIRE USO DI AUTOVAL.  
E AUTOCARICA VILUPPATA  
PER RENDERE DELL'APPROSS.  
STIMAzione DELLA MATTRE  
ORIGINARIA CON  
PENSO RANGO

APPLICATIONS  
- FORMAS NEO PROBLEMA (RANGE)  
- KERNEL / IMMAGINE  
- ISATICA PER DEFINIRE LA  
- PSEUDODINVERSA  
- APPROSSIMA MATTRE CON  
UNA PER RANGO INFERIORE

$$A = M A M^T$$

DIAGONAL

ORTHOGONAL:  $M^T M = M M^T = I$  (IDENTITY MATRIX)

$$AM = M\lambda$$

$$\lambda = \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \lambda_n \end{bmatrix}$$

$$M = \begin{bmatrix} v_1 & v_2 & \dots \end{bmatrix}$$

$$A \sqrt{\lambda}_1 = A_1 \sqrt{\lambda}_1$$

$$A \sqrt{\lambda}_2 = A_2 \sqrt{\lambda}_2$$

eigenvalues

BE CAREFUL: I HAVE  $A \geq 0$  SO  $M^T = M^{-1}$  OR  $A = A^T \geq 0$

IF IT'S SYMMETRIC IT'S DIAGONALIZABLE

There are matrices in  $\mathbb{R}^{n \times p}$

• NORMAL BUT NOT SYMMETRIC

$$\begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$$

so can be decomposed

$$M \Lambda M^T$$

LAMBDA  $\lambda$   
COMPLEX  
COMPLEX  
KTR. CONJ

$$s^2 + 1 = 0 \quad \zeta^i$$

• NOT NORMAL BUT DIAGONALIZABLE ANYWAY

$$A = \begin{bmatrix} 1 & 4 \\ 0 & 2 \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}}_{\text{COMPLEX}} \underbrace{\begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}}_{\text{DIAGONAL}} \underbrace{\begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix}}_{\text{COMPLEX}}$$

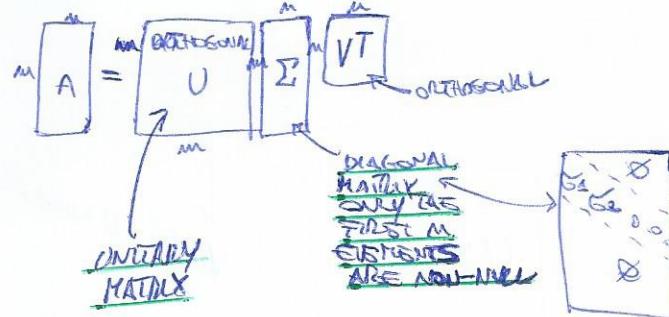
• NOT SQUARE!

• SQUARE BUT NOT DIAGONALIZABLE!

Jordan decomposition diagonal by blocks

Let  $A \in \mathbb{R}^{m \times m}$ . I can an SVD of A a decomposition like

$$A = U\Sigma V^T$$



THE ELEMENTS OF  $\Sigma$  ARE CALLED  
SINGULAR VALUES  
 $\sigma_1 > \sigma_2 > \sigma_3 - \dots > \sigma_n > 0$   
 $n = \text{RANK OF } A$   
THIS IS THE CRUCIAL MATRIX!

### THEOREM

Every matrix in  $\mathbb{R}^{m \times n}$  admits an SVD.  
Furthermore,  $\Sigma$  is uniquely determined by A.

$$A = U\Sigma V^T$$

$$\left. \begin{array}{l} A\sigma_1 = \sigma_1 u_1 \\ A\sigma_2 = \sigma_2 u_2 \\ \vdots \\ A\sigma_n = \sigma_n u_n \\ A\sigma_{n+1} = 0 \\ A\sigma_{n+2} = 0 \\ \vdots \\ A\sigma_m = 0 \end{array} \right\} \in \mathbb{R}^m$$

$\xrightarrow{\mathbb{R}^m}$   $\xrightarrow{\mathbb{R}^m}$

$\{u_1, u_2, \dots, u_n\} \subseteq \mathbb{R}^m$  IT'S A BASIS OF RANK A

$\{v_{n+1}, \dots, v_m\} \subseteq \mathbb{R}^n$  IT'S A ORTHONORMAL BASIS OF NULL A

A maps the subspace to zero

Fact: for  $A \in \mathbb{R}^{m \times n}$ ,  $B \in \mathbb{R}^{n \times m}$

$AB$  AND  $BA$  HAVE THE SAME NONZERO EIGENVALUES AND WITH THE SAME MULTIPLICITY

$AAT$  AND  $ATA$  HAVE THE SAME NON-ZERO EIGENVALUES

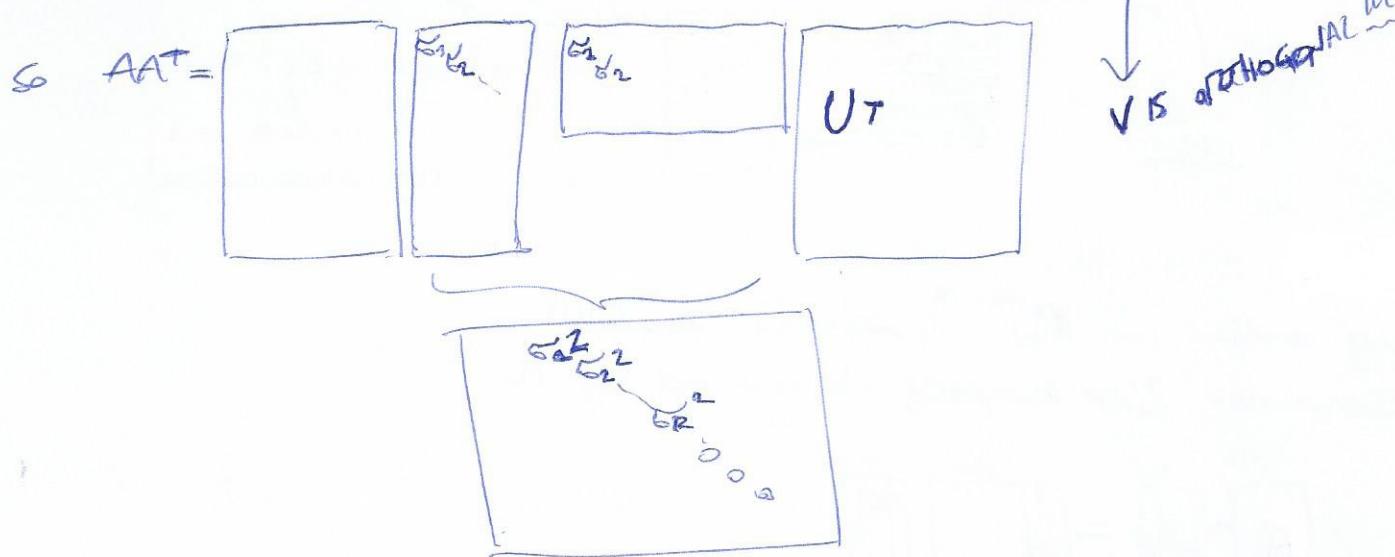
$AAT \in \mathbb{R}^{m \times m}$   
 $ATA \in \mathbb{R}^{n \times n}$

BOTH SYMMETRIC, POSITIVE DEFINITE ( $> 0$ )  
CAN BE DIAGONALIZED WITH ORTHOGONAL MATRIX

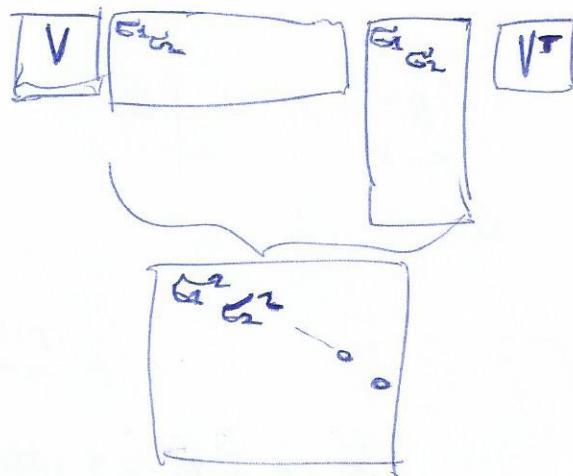
Indeed recall the definition;

$$\forall x \quad x^T (A A^T) x = (x^T A)(A^T x) = (A^T x)^T (A^T x) = \|A^T x\|^2 \geq 0$$

let  $A \in \mathbb{R}^{m \times n}$ ,  $A = U \Sigma V^T$ ,  $A A^T = U \Sigma V^T (\Sigma V^T U)$



$$A^T A = V \Sigma^T U^T U \Sigma V^T$$



### IMPORTANT THINGS:

- 1) THERE EXISTS AN SVD  $A = U \Sigma V^T$
- 2) (\*)
- 3) The columns of  $V$  form a basis of range  $A A^T$   
(orthogonal matrix in diagonalization of  $A A^T$ )
- 4) The columns of  $V$  form a basis of range  $A^T A$   
( $U$  = orthogonal matrix in diagonalization of  $A^T A$ )
- 5)  $\sigma_1 - \sigma_r$  are the square roots of the non-zero eigenvalues  
of both  $A^T A$  and  $A A^T$

# LINEAR EQUATION (System of linear equations)

$$Ax = b$$

$A \in \mathbb{R}^{m \times n}$   
 $b \in \mathbb{R}^m$   
 find  $x \in \mathbb{R}^n$

UNIQUE SOLUTION: this means  $A^{-1}$  exists  
 $(AA^T)x = A^Tb \Rightarrow x = (AA^T)^{-1}A^Tb$

INFINITELY MANY SOLUTIONS: we must choose one  
 (singular matrix)  $\text{NULL } A \neq \{0\}$

$\forall v \in \text{Null } A$ , if  $\bar{x}$  is a solution

$(A\bar{x} = b)$  then  $\bar{x} + v$  is also ( $A(\bar{x} + v) = A\bar{x} + Av = b$ )  
 ↳ many  $x \rightarrow$  CHOOSE THE ONE WITH MINIMUM NORM  $\|x\|$

NO SOLUTIONS: we gt find an approximate one

find  $x$  that attains

$$\min_x \|Ax - b\|_2$$

so, refer to the  
 CASE 2) WITH  
 INFINITELY MANY SOLUTIONS

there can be many  $x$  that minimize  
 $\|Ax - b\|_2$   
 ↳ min not unique

Def:  $S(A, b) = \{x \in \mathbb{R}^n \mid x \text{ minimizes } \|Ax - b\|_2\}$

Any  $x \in S(A, b)$  is "Almost" or really a solution.

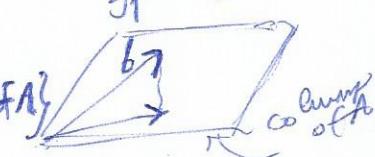
Def:  $x^* := \underset{x \in S(A, b)}{\text{arg min}} \|x\|_2$  is called the LEAST SQUARES  
 SOLUTION of the equation  $Ax = b$

PROPOSITION: A LS solution of  $Ax = b$  always exists and is

unique

We have to minimize:  $\min \|Ax - b\|_2 \Rightarrow \|Ax - b\|_2 = 0$   
 ↳ projection of  
 $b$  onto  
 span {column of  $A^T$ }

$x^*$  = projection of  
 any such  $x$   
 onto span {column of  $A^T$ }



$$E[y(t)^2] = \alpha_1^2 E[y(t-1)^2] + 5e^2$$

A process  $\{y(t)\}$  is wide-sense STATIONARY if  
 $t \in \mathbb{Z}$

1)  $E[y(t)] = \mu$  DOES NOT DEPEND ON T

2)  $E[y(t)y(t+\tau)] = \pi(\tau)$  <sub>TIME AVE</sub> CORRELATION SIGNAL, DOES NOT DEPEND ON T

$\{y(t)\}, t \in \{0, 1, 2, \dots\}$  is QUASI-STATIONARY

1)  $|E[y(t)]| \leq c$  } bounded

2)  $|E[y(t)y(t+\tau)]| \leq c$

3)  $\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=0}^{N-1} E[y(t)y(t+\tau)] = \pi(\tau) \quad (\exists \text{ A.U.M.T})$

THE AVERAGE OF THE CORRELATIONS EXISTS

**PROPERTY:** if  $y(t)$  is a quasistationary process then

$\lim_{t \rightarrow \infty} E[y(t)]$  } exists

$\lim_{t \rightarrow \infty} E[y(t)^2]$

so :  $E[y(t)^2] = \alpha_1^2 E[y(t-1)^2] + 5e^2$

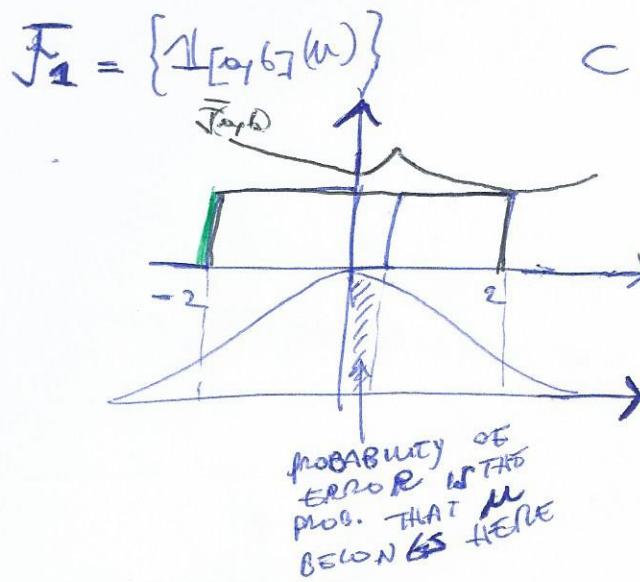
$$\Rightarrow E_y^2 = \alpha_1^2 E_y^2 + 5e^2 \rightarrow E_y^2 = \frac{5e^2}{1 - \alpha_1^2}$$

**Prop:** if  $W(z)$  is BIBO-stable 1)  $y(t)$  is quasi-stationary

2)  $\frac{1}{N} \sum_{t=0}^{N-1} y(t) \xrightarrow{\text{A.S.}} \mu = \lim_{t \rightarrow \infty} E[y(t)]$  } ERGODICITY

3)  $\frac{1}{N} \sum_{t=0}^{N-1} y(t)y(t+\tau) \xrightarrow{\text{A.S.}} \pi(\tau)$

FIRST: WE CHOOSE



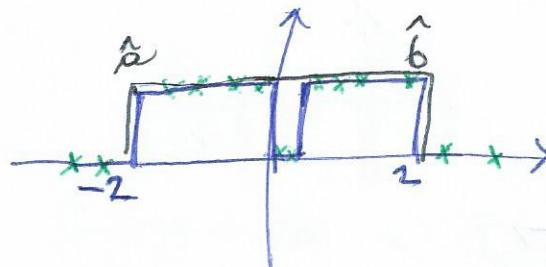
$$C = \{(\alpha, b) \in \mathbb{R}^2 \mid \alpha < b\}$$

IF FIX  $\alpha = -2$  AND ENLARGE  $b$

$$\bar{C} = (-2, 2)$$

$$\begin{aligned}\bar{J}(\bar{C}) &= \text{probability that } U \in [0, \frac{1}{2}] \\ &= \frac{1}{\sqrt{2\pi}} \int_0^{1/2} e^{-t^2/2} dt \approx 0.13\end{aligned}$$

What about  $\hat{f}_N, \hat{C}_N$ ?



$$\hat{C}_N = [-2 - \varepsilon, 2 + \varepsilon]$$

$$\text{SAY } \hat{f}_N(\hat{C}_N) = 0.125$$

IN THIS CASE THE APPROXIMATION IS GOOD  $\bar{J}(\hat{C}_N) \approx \bar{J}(C)$

SECOND: NEGLECT

$$\mathcal{F}_2 = \mathcal{F}_1 \cup \{1_{\text{finite retch}}(u)\}$$

$$C_2 = \{(\alpha, b)\} \cup \{S = \text{finite retch}\}$$

$$\mathcal{F}_2 = \{\sqsubset\} \cup \{\underbrace{\sqcap \sqcap \sqcap \sqcap}_{n \text{ (any)}}\}$$

NEEDLE FUNCTIONS  
THE PROBABILITY TO BE  
AFTER IS  $\varphi$

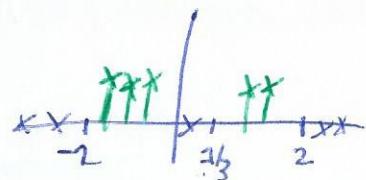
$\hat{f}_c$  minimizing  $\bar{J}$ : the same AS BEFORE

Regarding  $\bar{J}$ , the optimal  $\hat{f}_c$  in  $\bar{J}_1$  is the same or before  
 $(\bar{J}(C) = 0.13)$  But for any needle function  $\hat{f}_c$ ,  $\bar{J}(C) = \int_{-2}^{1/2} e^{-t^2} dt +$

$$+ \int_{1/2}^2 \varphi dt = 0.82$$

(IT'S BAD! WE DON'T HAVE TO CHOOSE ONE FUNCTION AMONG "NEEDLES FUNCTIONS")

What about  $\hat{f}_N, \hat{c}_N$ ?

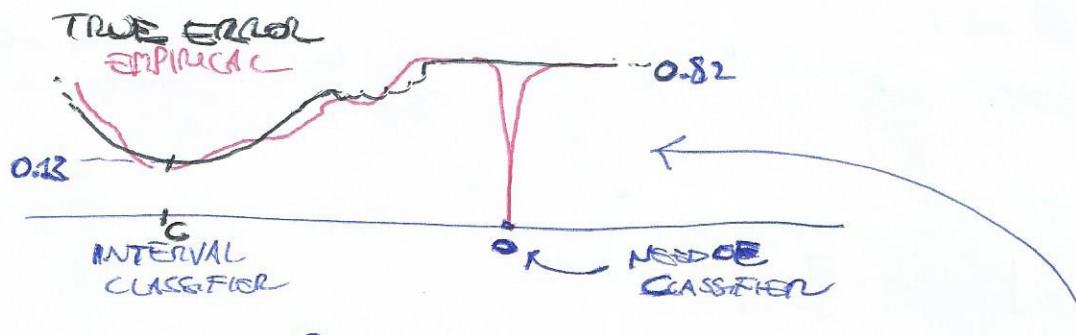


choosing  $\hat{f}_c = \{\mathbb{1}_{\{u_i | Y_i=1\}}\}$

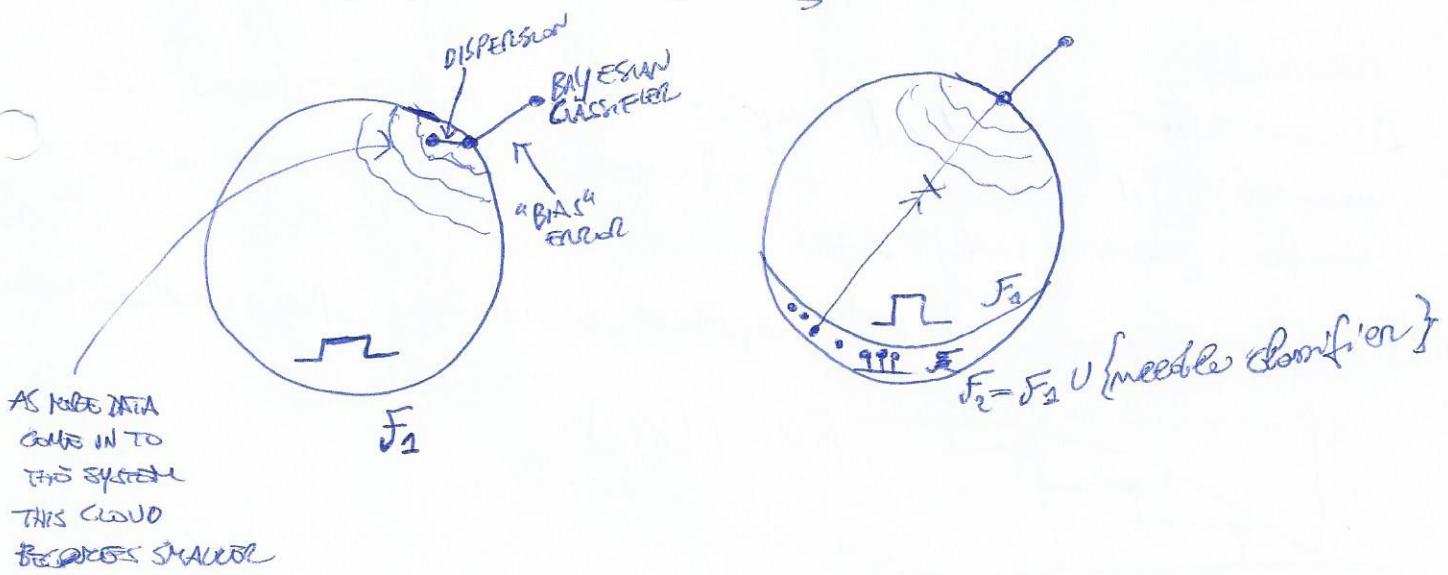
$$\sum_{i=1}^n \underbrace{\mathbb{1}(Y_i \neq \hat{f}_c(u_i))}_0 \Rightarrow \hat{f}_N(\hat{c}_N) = 0$$

$$\bar{J}(\hat{c}_N) = 0.82$$

BECAUSE THE SET  $\{\mathbb{1}_{\text{finite set } (n)}\}$  IT'S TOO COMPLEX!!!



$\hat{J}_N \xrightarrow{N \rightarrow \infty} \bar{J}$  UNIFORMLY? NO! (BECAUSE OF THIS)



$\hat{J}_N$  "escapes" to  $\bar{J}$  if  $\sup_{c \in C} |\hat{J}_N(c) - \bar{J}(c)| \leq \epsilon$