

EXERCISE : use the data CLIMATE.MAT

Analyse a multivariate regression model, having all the pollutants as response variables,

$$\underline{Y} = \begin{bmatrix} \text{PM10} \\ \text{O2} \\ \text{N} \\ \text{CO} \end{bmatrix}$$

and some climatic variable as regressor

$$\underline{X} = \begin{bmatrix} \text{meant} \\ \text{Humidity} \\ \text{meant Wind speed} \\ \text{rain, mm} \end{bmatrix}$$

Test the significance of the overall regression and then test if some regressor can be eliminated from the model.

# (UNIVARIATE) NONLINEAR REGRESSION

[Draper - Smith, Applied Regression Analysis, 1981]

Consider models which are nonlinear in the parameters

Examples:

$$1) Y = \exp(\theta_1 + \theta_2 t^2 + \varepsilon)$$

$$2) Y = \frac{\theta_1}{\theta_2 - \theta_1} [e^{-\theta_2 t} - e^{-\theta_1 t}] + \varepsilon$$

$\underline{\theta} = (\theta_1, \theta_2)$  parameters

$t = \text{response}$

$Y = \text{response} \in \mathbb{R}^1$

Note that:

- 1) Is intrinsically linear, since  $\log Y = \theta_1 + \theta_2 t^2 + \varepsilon$  is linear
- 2) Is intrinsically nonlinear, since no transformations may reduce 2) to the linear case

We will deal with models like 2) -

# LEAST SQUARES

Suppose that the postulated model is of the form

$$Y = f(\xi_1, \xi_2, \dots, \xi_k; \theta_1, \theta_2, \dots, \theta_p) + \varepsilon \quad (A)$$

with

$$\underline{\xi} = (\xi_1, \xi_2, \dots, \xi_k)^T \quad \text{REGRESSORS (before called } X\text{'s)}$$

$$\underline{\theta} = (\theta_1, \theta_2, \dots, \theta_p)^T \quad \text{PARAMETERS}$$

$\Rightarrow$  (A) is equivalent to

$$Y = f(\underline{\xi}, \underline{\theta}) + \varepsilon \quad (B1)$$

with the additional assumptions

$$\begin{cases} E(\varepsilon) = 0 \\ \text{Var}(\varepsilon) = \sigma^2 \\ \varepsilon \sim N(0, \sigma^2) \end{cases} \quad (B2)$$

If we now have a joint sample

$$(Y_u, \xi_{1u}, \xi_{2u}, \dots, \xi_{ku}) \quad u = 1, 2, \dots, n$$

We can write  $n$  equations like (B1) - (B2) for  $u = 1, \dots, n$

$$(B3) \quad \begin{cases} Y_u = f(\xi_u, \underline{\theta}) + \varepsilon_u & u = 1, \dots, n \\ \underline{\varepsilon} = N(\underline{0}, \sigma^2 I_n) \end{cases}$$

with  $\underline{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^T$ .

A least squares estimate  $\hat{\underline{\theta}}$  of  $\underline{\theta}$ , will minimise the error sum of squares:

$$S(\underline{\theta}) = \sum_{u=1}^n \varepsilon_u^2 = \sum_{u=1}^n (Y_u - f(\xi_u, \underline{\theta}))^2$$

Also in the nonlinear case the LS estimators of  $\underline{\theta}$  are also ML estimators, under the assumption of normality of errors.

In the linear case we found  $\hat{\theta}$  by solving the normal equations, which can be found by differentiating  $S(\theta)$  with respect to  $\theta_i$ :

$$\frac{\partial S(\theta)}{\partial \theta_i} = \sum_{u=1}^m (Y_u - f(\xi_u, \theta)) \left[ \frac{\partial f(\xi_u, \theta)}{\partial \theta_i} \right] = 0 \quad i=1, \dots, p \quad (N)$$

But in the linear case, where  $f$  had the form e.g. of the type

$$f(\xi_u, \theta) = \theta_1 \xi_{1u} + \theta_2 \xi_{2u} + \dots + \theta_p \xi_{pu}$$

we had  $\frac{\partial f(\xi_u, \theta)}{\partial \theta_i} = \xi_{iu}$  independent of  $\theta$

Now this is not any more true.

Example:  $f(\xi_u, \theta) = e^{-\theta \xi_u} \Rightarrow \frac{\partial f}{\partial \theta} = -\xi_u e^{-\theta \xi_u}$

i.e.  $Y_u = e^{-\theta \xi_u} + \varepsilon_u$

and the normal equations become

$$\sum_{u=1}^m [y_u - e^{-\theta \xi_u}] [-\xi_u e^{-\theta \xi_u}] = 0$$

which can not easily be solved with respect to  $\theta$ .

The minimum of  $S(\theta)$  can be found via different iterative algorithms.

Here we will introduce some of these methods.

## METHOD 1: LOCAL LINEARIZATION (GAUSS-NEWTON METHOD)

Let  $\hat{\underline{\theta}}_0$  be our initial guess for  $\hat{\underline{\theta}}$ ,  $\hat{\underline{\theta}}_0 = (\theta_{10}, \theta_{20}, \dots, \theta_{p0})^T$

We make a Taylor series expansion of  $f(\underline{x}_u, \underline{\theta})$  about  $\underline{\theta}_0$  and we approximate

$$f(\underline{x}_u, \underline{\theta}) \approx f(\underline{x}_u, \underline{\theta}_0) + \sum_{i=1}^P \left[ \frac{\partial f(\underline{x}_u, \underline{\theta})}{\partial \theta_i} \right]_{\underline{\theta} = \underline{\theta}_0} (\theta_i - \theta_{i0})$$

Now setting  $f_u^0 := f(\underline{x}_u, \underline{\theta}_0)$

$$\beta_i^0 := \theta_i - \theta_{i0}$$

$$z_{ui}^0 := \left[ \frac{\partial f(\underline{x}_u, \underline{\theta})}{\partial \theta_i} \right]_{\underline{\theta} = \underline{\theta}_0}$$

the model (B3) is approximated by

$$\left( y_u - f_u^0 \right) = \sum_{i=1}^P \beta_i^0 z_{ui}^0 + \varepsilon_u \quad u = 1, \dots, m$$

↑  
LINEAR IN THE  $\beta_i^0$ 'S !!

Now by defining

$$\underline{z}_0 = \begin{bmatrix} z_{11}^0 & z_{12}^0 & \dots & z_{1p}^0 \\ \vdots & & & \\ z_{m1}^0 & z_{m2}^0 & \dots & z_{mp}^0 \end{bmatrix} = \{z_{xi}^0\} \quad m \times p$$

$$\underline{\beta}^0 = \begin{bmatrix} \beta_1^0 \\ \vdots \\ \beta_p^0 \end{bmatrix} = \begin{bmatrix} \theta_1 - \theta_{10} \\ \vdots \\ \theta_p - \theta_{p0} \end{bmatrix} \quad \text{and} \quad \underline{y}_0 = \begin{bmatrix} y_1 - \mu_1^0 \\ \vdots \\ y_m - \mu_m^0 \end{bmatrix} \quad \underline{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_m \end{bmatrix}$$

the model becomes, with matrix notation

$$\begin{cases} \underline{y}_0 = \underline{z}_0 \underline{\beta}^0 + \underline{\varepsilon} \\ \underline{\varepsilon} \sim N(\underline{0}, \sigma^2 \underline{I}_m) \end{cases} \quad \text{LINEAR MODEL}$$

and the LS-estimator of  $\underline{\beta}^0$  is  $\hat{\underline{\beta}}^0 = \underline{S}_0^{-1} \underline{z}_0^T \underline{y}_0$   
with  $\underline{S}_0 = \underline{z}_0^T \underline{z}_0$ , assuming it is invertible.



From  $\hat{\beta}^0$  we can recover one update of the estimate of  $\theta$ :

$$\hat{\theta}_1 = \theta_0 + \hat{\beta}_0$$

Then we use  $\hat{\theta}_1$  as initial guess and iterate the procedure:

$$\Rightarrow \hat{\theta}_{j+1} = \hat{\theta}_j + \hat{\beta}_j = \hat{\theta}_j + S_j^{-1} z_j (y_j - f_j) = \hat{\theta}_j + S_j^{-1} z_j (y_j - f(\xi_j, \hat{\theta}_j^j))$$

$\uparrow$   
 matrix of first derivatives of  $f$  computed in  $\hat{\theta}_j^j$

up to when

$$|\{\theta_{i(j+1)} - \theta_{ij}\} / \theta_{ij}| < \delta \quad \forall i=1, \dots, p$$

$\uparrow$   
 fixed tolerance

At each iteration,  $S(\hat{\theta}_j)$  can be computed to control if a reduction in its value has been achieved.

## Problems of this method:

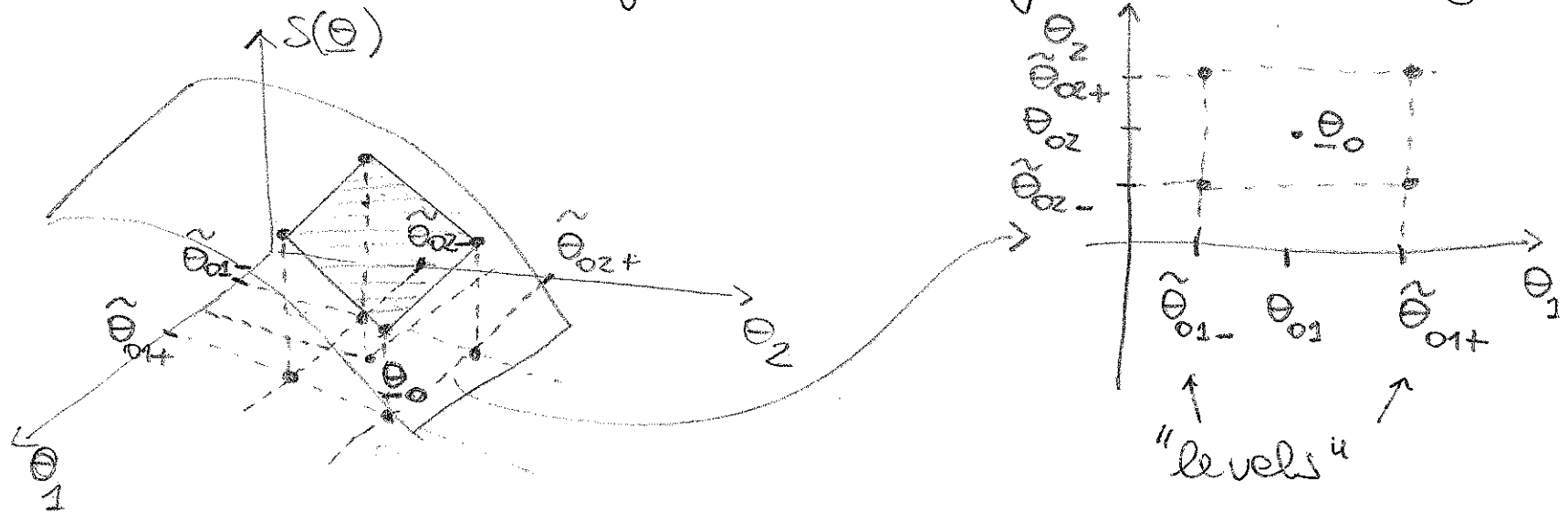
- 1) It may converge very slowly
- 2)  $S(\hat{\theta}_j)$  may oscillate very widely, even if at some point it stabilizes
- 3) It may not converge at all, and even diverge  $\Rightarrow$   
In this case, reparameterize the model or use the Marquardt's Compromise
- 4) It may converge to a local minimum of  $S(\theta)$  or to a maximum  $\Rightarrow$  test many and widely spread initial guesses  $\theta_0$  !!

## METHOD 2: OF STEEPEST DESCENT

Concentrate on  $S(\underline{\theta})$ .

Start from an initial point (guess)  $\underline{\theta}_0$  then

- fix an experimental factorial design around  $\underline{\theta}_0$



- Compute  $S(\underline{\theta})$  in the points of the experimental design
- Approximate  $S(\underline{\theta}_0)$  locally with the hyperplane passing through the points of the design:

$$S(\underline{\theta}) \approx \beta_0 + \sum_{i=1}^P \frac{\beta_i (\theta_i - \bar{\theta}_i)}{d_i} + \varepsilon \quad (*)$$

↑ scaling factor such that  
 $\sum_{u=1}^m (\theta_{iu} - \bar{\theta}_i)^2 / d_i^2 = \text{const.}$

$\bar{\theta}_i$  = means of the levels  $\theta_{iu}$ ,  $u=1, 2, \dots, m$   
of the levels of the design.

The estimated coefficients of (\*)

$$\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_p$$

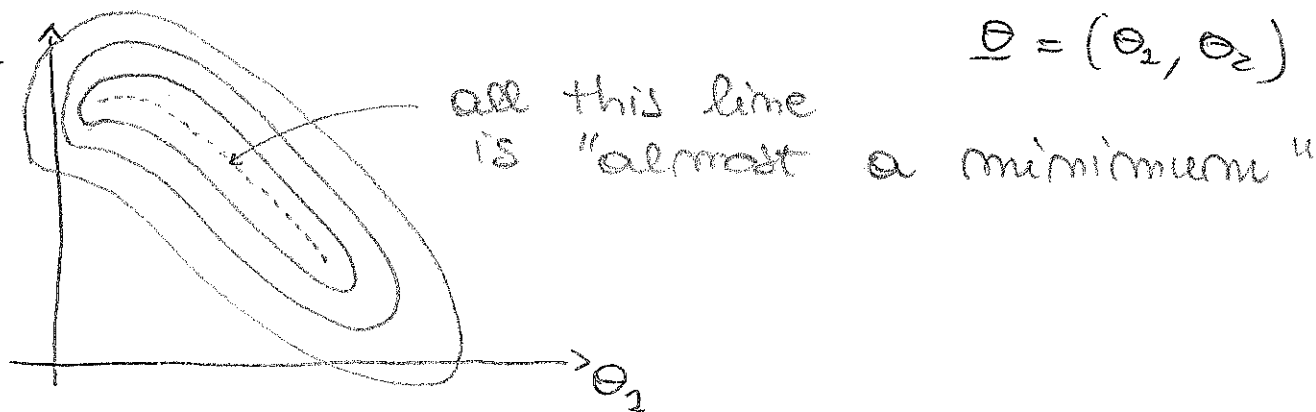
indicate the direction of the steepest ascent, thus

$$-\hat{\beta}_1, -\hat{\beta}_2, \dots, -\hat{\beta}_p$$

indicate the direction of the steepest descent  $\Rightarrow$  we  
will move on  $S(\underline{\theta})$  to a  $\underline{\theta}_1$  following this direction,  
then iterate.

## Problems of this method:

- 1) It may converge very slowly, in particular when the contour levels of  $S(\underline{\theta})$  are "banana shaped", leading to "zigzagging".



- 2) The method is not scale invariant, i.e. is sensitive to the choice of  $\alpha_i$ , if they are not changed by the same factor.

An advantage of this method is that it works well when  $\underline{\theta}_0$  is far from the point of global minimum  $\hat{\underline{\theta}}$ , which is often the case.

### METHOD 3: MARQUARDT'S COMPROMISE

It is a compromise between linearization (Gauss-Newton) and the steepest descent, and works better of both of them in most practical cases.

IDEA: Start from an initial guess  $\underline{\theta}_0$ .

- Compute the direction  $\underline{\delta}_g$  where to move on  $S(\underline{\theta})$  to find  $\underline{\theta}_1$ , using the steepest descent method;
- Compute the direction  $\underline{\delta}$  where to move on  $S(\underline{\theta})$  to find  $\underline{\theta}_2$ , using the linearization method;
- Interpolate suitably  $\underline{\delta}_g$  and  $\underline{\delta}$  for obtaining a new "best direction" and a suitable step size

Remark: on many examples Marquardt found that the angle  $\phi$  between  $\underline{\delta}_g$  and  $\underline{\delta}$  is  $80^\circ < \phi < 90^\circ$  !!

## CONFIDENCE REGIONS

① Use an ellipsoidal confidence region by assuming that in a neighborhood of  $\hat{\underline{\theta}}_{LS}$ , the model can be approximated by its linearized form. The ellipsoid is

$$(\underline{\theta} - \hat{\underline{\theta}}_{LS})^T \underset{\sim}{\hat{\underline{z}}} \underset{\sim}{\hat{\underline{z}}}^T (\underline{\theta} - \hat{\underline{\theta}}_{LS}) \leq p \sigma^2 F_{1-\alpha; p, m-p} \quad (**)$$

where  $\underset{\sim}{\hat{\underline{z}}}$  is the matrix of first derivatives of  $f$ , computed in  $\hat{\underline{\theta}}_{LS}$  and

$$\sigma^2 = \frac{S(\hat{\underline{\theta}}_{LS})}{m-p}$$

$p = m - \text{of parameters to be estimated}$

Note that (\*\*) is only an approximate  $1-\alpha$  confidence region.

② An exact confidence contour is defined by taking

$$S(\underline{\theta}) = \text{const.} = c$$

but since we don't know the distribution of  $S(\underline{\theta})$ , we can't determine  $c$  such that the confidence level of the region is exactly  $1-\alpha$ .

We can anyway use the contour such that

$$S(\underline{\theta}) = S(\hat{\underline{\theta}}_{LS}) \left\{ 1 + \frac{p}{m-p} F_{1-\alpha; p, m-p} \right\}$$

which is exactly of level  $1-\alpha$  if the model is linear, and is approximately of level  $1-\alpha$  in the non linear case.

Note that this is a confidence contour, is its confidence level which is approximate!



Confidence intervals for the single parameters can be obtained by sectioning the confidence regions.

## Exercise

Consider the climate-pollutants data.

Find a suitable (non linear) model to relate  $O_3$  and mean temperature (meant)

$$O_3 = f(\beta, \text{meant}) + \varepsilon$$

↑  
?  
,