

Modeling and design of experiments

Chapitre 2: Modeling

Dr Jean-Marie Fürbringer

École Polytechnique Fédérale de Lausanne

Fall 2025

Modeling

Empirical Modeling

Linear system

Geometric interpretation

Model coefficient variance

Example of elasticity

2.1.1 The school of Athens - 1510 - Raphael



2.1.2 Plato & Aristotle

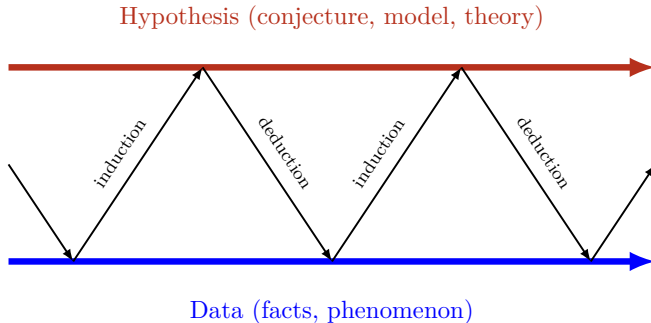


2.1.3 Empiricism vs theory

The men of experiment are like the ant, they only collect and use; the reasoners resemble spiders, who make cobwebs out of their own substance. But the bee takes a middle course : it gathers its material from the flowers of the garden and of the field, but transforms and digests it by a power of its own.

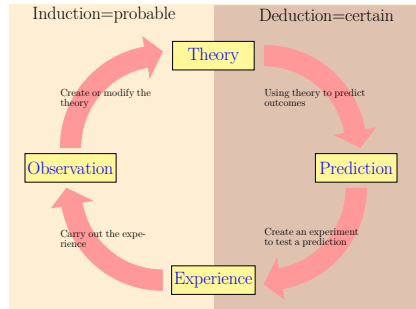
"The new organon", Francis Bacon, 1620.

2.1.4 Scientific process of modelling

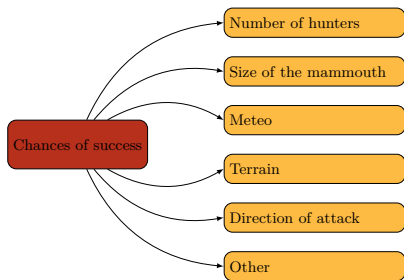


2.1.5 The scientific method

- ▶ **Deductive argument**
 - ▶ The validity of the argument is based on its structure (premises \rightarrow conclusion)
 - ▶ Application : predictions
- ▶ **Inductive argument**
 - ▶ Relevance is based on the representativeness of the sample (quality)
 - ▶ Sufficiency is based on sample size(quantity)
 - ▶ Application : conceptualization, learning



2.1.6 Causation and mental model

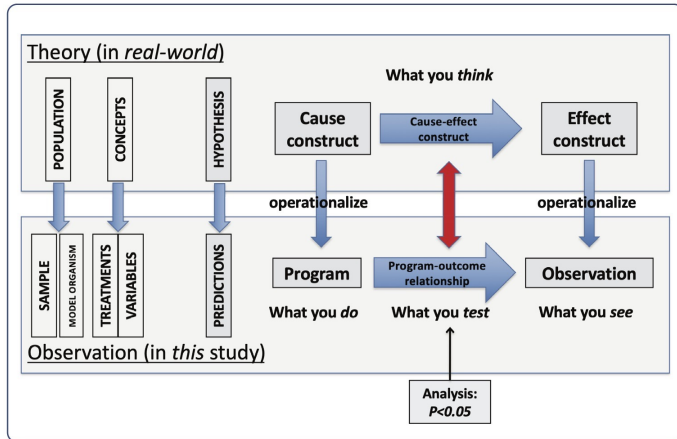


Knowledge revolution $\approx 50'000$ years
"The book of why", Judea Pearl, 2018

2.1.7 What is a good hypothesis ?

1. It can be **tested experimentally**
 - ▶ Question to which it can be answered by yes or not (truth proposition)
 - ▶ Excludes supernatural explanations
2. It is **sufficient** : it explains all the relevant facts
 - ▶ What is relevant ? (circular approach)
3. It is based on the **minimum of premises** (does not generate puzzles)
 - ▶ Ockam's razor (Parsimony principle)
4. It takes into account **established knowledge**
 - ▶ Problem : established knowledge can be erroneous ! to be verified by experience
5. It is more plausible than the **alternative hypotheses**
 - ▶ It investigates a maximum of hypothesis : imagination and knowledge

2.1.8 Operationalisation : from concept to reality

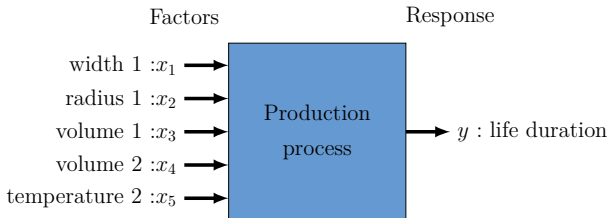


2.1.9 Result interpretation

Type of validity	Question	Threat
4. External validity (generalisation)	Is it possible to generalize the conclusions of the experiments to the real world?	<ul style="list-style-type: none"> ▶ weak model ▶ weak sampling
3. Construct validity (translation of concepts in variables)	Treatments and the measurements are correctly connected to the definition of causes and effects?	<ul style="list-style-type: none"> ▶ weak model ▶ weak indicators ▶ weak signals ▶ Bias
2. Internal validity (verify that treatments refute alternative explanations)	Did the treatment cause the results?	<ul style="list-style-type: none"> ▶ weak control ▶ weak blocking ▶ weak randomization ▶ attrition
1. Conclusion validity (statistical analysis)	Does it exist a relation between the treatment and the results	<ul style="list-style-type: none"> ▶ P-value hacking

P-hacking is manipulating analyses or selectively reporting results until a "significant" p-value appears, which inflates false positives.

2.1.10 Bloc diagram



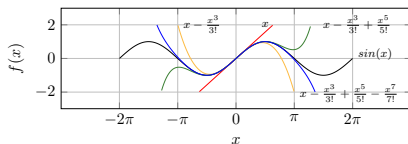
2.1.11 Taylor series

Theorem : If $U \subset \mathbb{R}^n$ and f is an application of U in \mathbb{R}^p of class C^q . $\forall \mathbf{x}_o \in U$ and $\mathbf{h} \in U$ so that $\mathbf{x}_o + \mathbf{h} \in U$, we have

$$f(\mathbf{x}_o + \mathbf{h}) = f(\mathbf{x}_o) + df_a(\mathbf{h}) + \frac{1}{2!}d^2f_a(\mathbf{h}, \mathbf{h}) + \dots + R_n$$

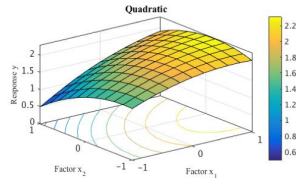
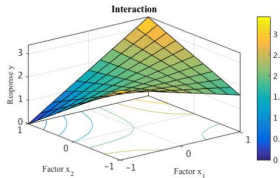
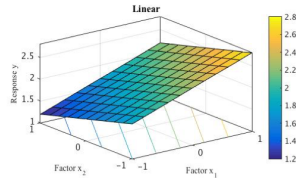
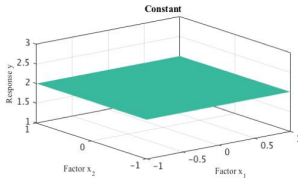
Taylor polynomial :

$$f(\mathbf{x}) = a_o + a_1x_1 + a_2x_2 + a_{12}x_1x_2 + a_{11}x_1^2 + a_{22}x_2^2 + \epsilon$$



2.1.12 Empirical model

$$y(x) = a_0 + a_1x_1 + a_2x_2 + a_{12}x_1x_2 + a_{11}x_1^2 + a_{22}x_2^2 + \epsilon$$



2.2.1 Matrix of experiments

$N_{exp} \times N_{fact}$ matrix whose components are the levels x_{ij} of the factors for each experiment :

$$E = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1N_{fact}} \\ x_{21} & x_{22} & \dots & \vdots \\ \vdots & & & \vdots \\ x_{N_{exp}1} & x_{N_{exp}2} & \dots & x_{N_{exp}N_{fact}} \end{bmatrix}$$

The component can be standardized or not

2.2.2 Model matrix

$N_{exp} \times N_{coef}$ matrix whose components are the product of the levels of the factors corresponding to each coefficients of the model for each experiment, and corresponding to a given matrix of experiment E .

To the model $y = a_o + a_1x_1 + a_2x_2 + a_{12}x_1x_2 + a_{11}x_1^2$, linked to the matrix of experiments E corresponds the model matrix :

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{11}x_{12} & x_{11}^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{N_{exp}1} & x_{N_{exp}2} & x_{N_{exp}1}x_{N_{exp}2} & x_{N_{exp}1}^2 \end{bmatrix} \text{ so that}$$

$$Y = X\eta = X \begin{bmatrix} a_o \\ a_1 \\ a_2 \\ a_{12} \\ a_{11} \end{bmatrix}$$

is a linear system. η is the vector of the model coefficients.

2.2.3 Least square fit

The ideal system $Y = X\eta$ is replaced by $Y = X\alpha + \epsilon$

The least square fit method makes possible the determination of α and ϵ based on standard hypotheses over ϵ :

$$Y = X\alpha + \epsilon \quad \text{with } \epsilon' \epsilon \rightarrow \min \quad (1)$$

Because optimal ϵ should be orthogonal to $X\alpha$, then

$$X' \epsilon = \vec{0} \quad (2)$$

$$X'(Y - X\alpha) = \vec{0} \quad (3)$$

$$X'X\alpha = X'Y \quad (4)$$

$$\hat{\alpha} = (X'X)^{-1}X'Y \quad (5)$$

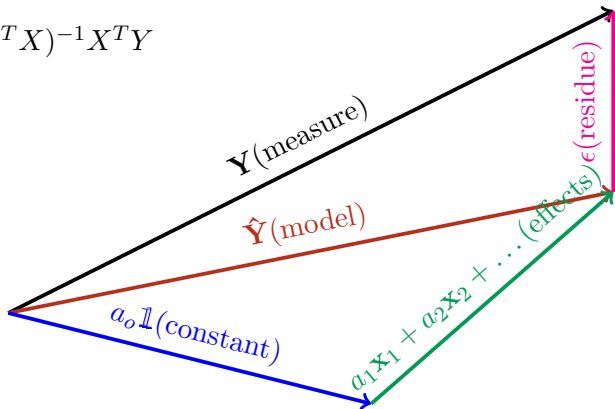
2.2.4 Dispersion matrix

$$\hat{a} = (X^T X)^{-1} X^T Y$$

- ▶ Also called variance-covariance matrix
- ▶ It represents the transfer of the experimental error to the model coefficients
- ▶ When possible a diagonal dispersion matrix is preferred
- ▶ At least with diagonal elements close to $1/N_{exp}$ and non-diagonal elements smaller (in absolute value) than the diagonal elements

2.3.1 Geometric point of view

$$\hat{\alpha} = (X^T X)^{-1} X^T Y$$



2.3.1 Correlation matrix

- ▶ Correlation quantifies the collinearity between two regressors
- ▶ The correlation coefficients are defined as :

$$\text{cor}(x_i, x_j) = \frac{\text{cov}(x_i, x_j)}{\sigma_{X_i} \sigma_{X_j}}$$

- ▶ The components of the matrix of correlation C can be computed with : $c_{ij} = \frac{D_{ij}}{\sqrt{D_{ii} D_{jj}}}$
with D_{ij} being the dispersion matrix components

MATLAB

$C = \text{corrcoef}(D)$

2.4.1 Variance inflation factor (VIF)

- ▶ Quantification of the severity of the collinearity of the regressors
- ▶ It is an index that measures how the variance of the coefficients is multiplied because of the collinearity
- ▶ The square root of VIF indicates how many times the standard error is multiplied
- ▶
$$\text{var}(\hat{a}_j) = \frac{s^2}{(N-1) \text{var}(x_j)} \frac{1}{1-R_j^2} = \frac{s^2}{(N-1) \text{var}(x_j)} \text{VIF}(\hat{a}_j)$$

MATLAB

$VIF = \text{diag}(\text{inv}(C))$

2.4.2 Variance of the model coefficient

▶ $\hat{a} = (X^T X)^{-1} X^T Y = UY$ with $U = (X^T X)^{-1} X^T$

▶ $var(\hat{a}) = var(UY) = U U^T var(y) = (X^T X)^{-1} var(y)$
because $U U^T = (X^T X)^{-1}$

▶ $var(y) \approx s^2 = \frac{\epsilon^T \epsilon}{N-P}$

▶ N is the number of experiments

▶ P is the number of coefficients

▶ $var(\hat{a}) = (X^T X)^{-1} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} s^2$

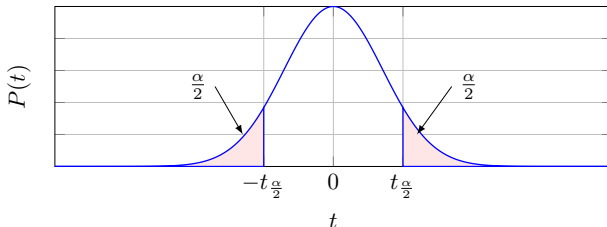
2.4.3 Confidence Intervals

For computing the level of confidence at $\beta * 100\%$:

- ▶ If $var(\hat{a}) = (X^T X)^{-1} s^2 = D s^2$
- ▶ $\alpha = 1 - \beta$
- ▶ If the degree of freedom is $\nu = (N - P)$
- ▶ If the quantile of the Student distribution is $t_{\alpha/2, \nu}$

Then

$$CI_{\beta}(a_i) = t_{\frac{\alpha}{2}, \nu} \sqrt{D_{ii} s^2}$$



2.4.3 Model coordinate $\vec{f}(\vec{x})$

- ▶ Vector of the model space parametrised with the coordinates x_i of the experimental space
- ▶ The model can then be written $y(\vec{x}) = \vec{f}(\vec{x}) \cdot \vec{\alpha}$ with $\vec{\alpha}$ being the model coefficients

$$\vec{f}(\vec{x}) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_1 x_2 \\ x_1 x_3 \\ \vdots \\ x_1^2 \\ x_2^2 \\ x_3^2 \end{bmatrix}$$

2.4.4 A new definition of the model Matrix X

The model matrix X can be written as the model coordinates $\vec{f}'(\vec{x}_i)$ computed at the points of measurement \vec{x}_i

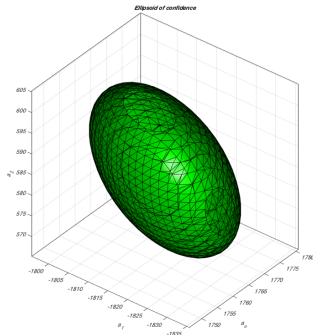
$$X = \begin{bmatrix} \vec{f}'(\vec{x}_1) \\ \vec{f}'(\vec{x}_2) \\ \vdots \\ \vec{f}'(\vec{x}_{N_{exp}}) \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{11}x_{12} & \dots & x_{11}^2 \\ \vdots & \vdots & \vdots & & \vdots & & \vdots \\ 1 & x_{N_{exp}1} & x_{N_{exp}2} & \dots & x_{N_{exp}1}x_{N_{exp}2} & \dots & x_{N_{exp}1}^2 \end{bmatrix}$$

2.4.5 Ellipsoid of confidence

- ▶ If the coefficients α of a model $g(x_i, \alpha_i)$ are considered as Normal random function :

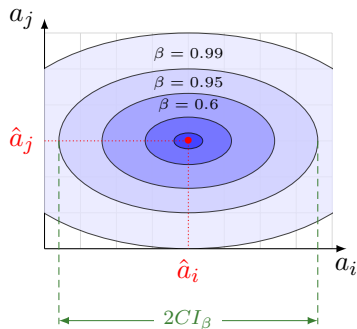
$$\hat{\alpha} \sim N(\alpha, (X^T X)^{-1} \sigma^2),$$
- ▶ Then, if η represents the true model, it can be written $\eta - \hat{Y} = X(\alpha - \hat{\alpha}),$
- ▶ And there is a probability $1 - \beta$ for finding the true value of the coefficients α within the ellipsoid defined by

$$(\alpha - \hat{\alpha})^T (X^T X) (\alpha - \hat{\alpha}) / ps^2 = F_\beta(p, \nu)$$
- ▶ If the regressors are not orthogonal, the confidence intervals of the coefficients are not independent

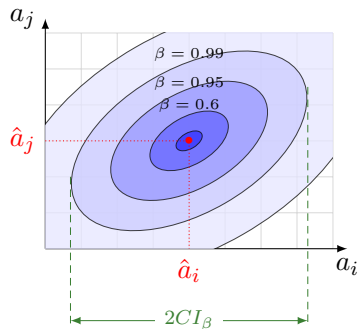


2.4.6 Interdependent confidence intervals

Orthogonal regressors



Non-orthogonal regressors



2.4.7 The relative variance function

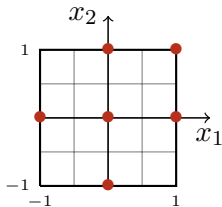
- ▶ Function of the experimental space
- ▶ Gives a prediction before the experiments of the ratio of the experimental uncertainty that will be transferred to the model
- ▶ We usually look for a low and uniform variance over the experimental space

$$\frac{\text{var}_Y(x)}{\sigma^2} = f^T(x) (X^T X)^{-1} f(x)$$

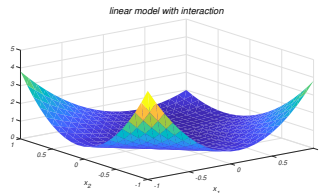
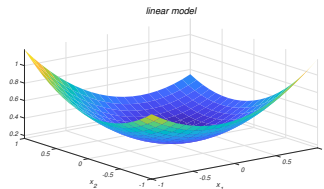
$\vec{f}(\vec{x}_i)$ is the model coordinate at the point \vec{x}_i

2.4.8 Example of the relative variance function

$$\frac{\text{var}_Y(x)}{\sigma^2} = f^T(x) (X^T X)^{-1} f(x)$$



$$E = \begin{pmatrix} 0 & 0 \\ 1 & 1 \\ -1 & 0 \\ 1 & 0 \\ 0 & -1 \\ 0 & 1 \end{pmatrix}$$



2.4.9 Computation of the variance function

If the model is $y = a_0 + a_1x_1 + a_2x_2$:

```
%matrix of essays
E=[-1 -1;-.6 0.17;0.2 -.5;1 1]
% matrix of the model
X=x2fx(E)
% matrix of dispersion
D=inv(X'*X)
% declaration of symbolic variable
syms x y
% built the model coordinate
f=[1;x;y]
% compute the variance function
v(x,y)=f'*D*f
% surface plot
h=fsurf(x,y,v,[-1 1]);
```

2.4.10 Information function

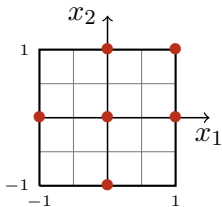
- ▶ Function of the experimental space
- ▶ Inverse of the variance function
- ▶ We usually look for high and uniform information in the experimental space

$$I_Y(x) = \frac{\sigma^2}{\text{var}_Y(x)} = \left[f^T(x) (X^T X)^{-1} f(x) \right]^{-1}$$

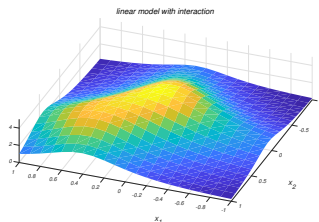
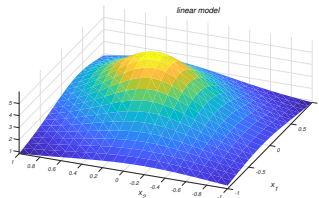
$\vec{f}(\vec{x}_i)$ is the model coordinate at the point \vec{x}_i

2.4.11 Example of the information function

$$I_Y(x) = \left(f^T(x) (X^T X)^{-1} f(x) \right)^{-1}$$



$$E = \begin{pmatrix} 0 & 0 \\ 1 & 1 \\ -1 & 0 \\ 1 & 0 \\ 0 & -1 \\ 0 & 1 \end{pmatrix}$$



2.4.12 Main concepts

Design

- ▶ Matrix of experiments
- ▶ Model coordinate
- ▶ Model matrix
- ▶ Dispersion matrix
- ▶ Correlation matrix
- ▶ Variance inflation factor
- ▶ Variance function
- ▶ Information function
- ▶ Alias matrix

Analysis

- ▶ Model coefficients
- ▶ Ellipsoid de confidence
- ▶ Interval de confidence
- ▶ ANOVA

2.5.1 Young modulus of steel

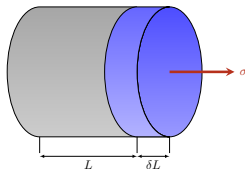


A company is in charge of the design of a bridge and then has to choose the steel that will be used for the beams. An important characteristic of a material when analysing its elasticity is the Young modulus E . In this case, the engineer want to test the Young modulus in function of the temperature and the concentration of carbon and sulfur characterizing the steel available in the market.

2.5.2 Example of multilinear regression

Hooke law : $\epsilon = \frac{\delta L}{L} = \frac{\sigma}{E}$

1. ϵ : relative strain [-]
2. σ : stress [kpa]
3. L : length of the sample [m]



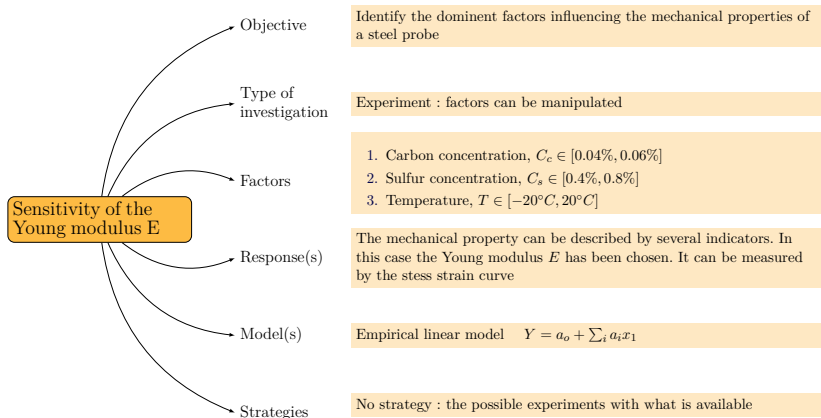
Sensitivity of the Young modulus : $E = f(T, C, S)$

T : temperature [$^{\circ}C$]

S : concentration of sulfur [%]

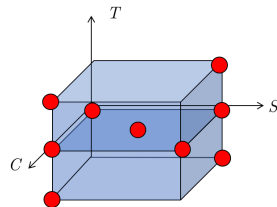
C : concentration of carbon [%]

2.5.3 Mindmap



2.5.4 Data Young modulus

	Carbon [%]	Sulfur [%]	Temperature [°C]	E [kpa]
1	0.04	0.4	-20	210.31
2	0.04	0.4	0	210.37
3	0.04	0.4	20	210.28
4	0.04	0.8	0	209.18
5	0.05	0.5	0	210.31
6	0.06	0.4	0	210.81
7	0.06	0.8	-20	209.70
8	0.06	0.8	0	209.58
9	0.06	0.8	20	209.67



2.5.5 Linear system

Model : $E = \alpha_o + \alpha_c C + \alpha_s S + \alpha_T T + \epsilon$

Matrix equation : $Y = X\alpha + \epsilon$

$$\begin{pmatrix} 210.31 \\ 210.37 \\ 210.28 \\ 209.18 \\ 210.31 \\ 210.81 \\ 209.70 \\ 209.58 \\ 209.67 \end{pmatrix} = \begin{pmatrix} 1 & 0.04 & 0.4 & -20 \\ 1 & 0.04 & 0.4 & 0 \\ 1 & 0.04 & 0.4 & 20 \\ 1 & 0.04 & 0.8 & 0 \\ 1 & 0.05 & 0.5 & 0 \\ 1 & 0.06 & 0.4 & 0 \\ 1 & 0.06 & 0.8 & 0 \\ 1 & 0.06 & 0.8 & 0 \\ 1 & 0.06 & 0.8 & 20 \end{pmatrix} \times \begin{pmatrix} \alpha_o \\ \alpha_C \\ \alpha_S \\ \alpha_T \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \epsilon_4 \\ \epsilon_5 \\ \epsilon_6 \\ \epsilon_7 \\ \epsilon_8 \\ \epsilon_9 \end{pmatrix}$$

2.5.6 Computing the effects

Estimator : $\hat{\alpha} = (X^T X)^{-1} X^T Y$

Matrix of dispersion :

$$(X^T X)^{-1} = \begin{pmatrix} \mathbf{3.3} & -59 & -0.36 & 0 \\ -59 & \mathbf{1652} & -40.18 & 0 \\ -0.36 & -40.18 & \mathbf{4.02} & 0 \\ 0 & 0 & 0 & \mathbf{0.0006} \end{pmatrix}$$

Coefficients :

$$\hat{\alpha}^T = (210.5[kpa], 24.1[kpa/\%C], -2.9[kpa/\%S], -8 \cdot 10^{-4}[kpa/K])$$

Differential : $dE = 24.1dC - 2.9dS - 0.0008dT$

Effects (variation of the output between the extremes of the experimental range) :

$$effect_C = 0.48[kpa], effect_S = -1.15[kpa], effect_T = -0.03[kpa]$$

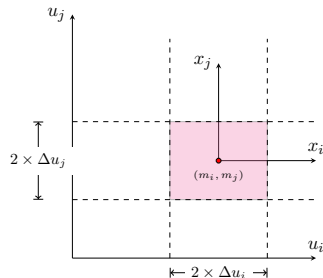
2.5.7 Let's take stock

- ▶ We have started with a matrix of experiments : A
- ▶ We have chosen a model with a linear response :
$$E = \alpha_o + \alpha_c C + \alpha_s S + \alpha_T T + \epsilon$$
- ▶ A model matrix is built :
- ▶ The least square fit method is used to compute the model coefficients
- ▶ The quality of the fit is evaluated with the statistic
$$R^2 = \frac{\hat{Y}^T \hat{Y}}{Y^T Y}$$
- ▶ The theorem of the confidence interval is used to verify the quality of the coefficients

2.5.8 Normalisation of the experimental domain

- ▶ Original variables : u_i, u_j
- ▶ Normalized variables x_i, x_j
- ▶ Median : $m_i = \frac{\max(u_i) + \min(u_i)}{2}$
- ▶ Half ranges : $\Delta u_i = \frac{\max(u_i) - \min(u_i)}{2}$
- ▶ Linear relations : $x_i = \frac{u_i - m_i}{\Delta u_i}$

$$u_i = m_i + x_i \Delta u_i$$



2.5.9 Normalizing the matrix of essay

$$E = \begin{pmatrix} 0.04 & 0.4 & -20 \\ 0.04 & 0.4 & 0 \\ 0.04 & 0.4 & 20 \\ 0.04 & 0.8 & 0 \\ 0.05 & 0.5 & 0 \\ 0.06 & 0.4 & 0 \\ 0.06 & 0.8 & -20 \\ 0.06 & 0.8 & 0 \\ 0.06 & 0.8 & 20 \end{pmatrix} \rightarrow E_{st} = \begin{pmatrix} -1 & -1 & -1 \\ -1 & -1 & 0 \\ -1 & -1 & 1 \\ -1 & 1 & 0 \\ 0 & -0.5 & 0 \\ 1 & -1 & 0 \\ 1 & 1 & -1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}$$

MATLAB

```
Est = rescale(E,-1,1,'InputMin',min(E),'InputMax',max(E))
```

2.5.10 Computing the normalized effects

Estimator : $\hat{\alpha} = (X^T X)^{-1} X^T Y$

Dispersion matrix : $(X^T X)^{-1} = \begin{pmatrix} \mathbf{0.12} & 0.05 & -0.03 & 0 \\ 0.05 & \mathbf{0.30} & -0.21 & 0 \\ -0.03 & -0.21 & \mathbf{0.28} & 0 \\ 0 & 0 & 0 & \mathbf{0.25} \end{pmatrix}$

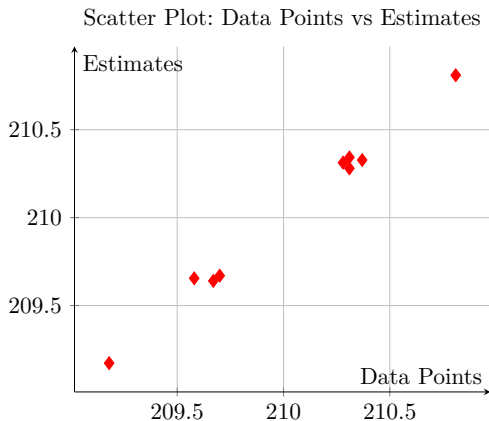
Half effects : $\hat{\alpha}^T = (210[kpa], 0.23[kpa], -0.62[kpa], -0.015[kpa])$

Model : $E = 210 + 0.23x_C - 0.62x_S - 0.015x_T$ avec $x_i \in [-1, 1]$

Relative effects : $ER_C = 0.11[\%]$, $ER_S = -0.3[\%]$, $ER_T = -0.01[\%]$
 ER= half effect divided by the constant a_o

2.5.11 Compare estimates and data

- ▶ Mandatory to loop the loop and check that prediction are sufficiently close to the data
- ▶ $\hat{y} = X\hat{\alpha}$



2.5.12 Coefficients in the original units

To have the coefficients in the units of the laboratory it is necessary to go back to the non normalized system. It is different for each model. Here is the case of a linear model.

$$y = a_o + \sum a_i x_i$$
$$x_i = \frac{u_i - \bar{u}_i}{\Delta u_i}$$

$$\begin{aligned} y &= a_o + \sum a_i \frac{u_i - \bar{u}_i}{\Delta u_i} \\ &= \left(a_o - \sum a_i \frac{\bar{u}_i}{\Delta u_i} \right) + \sum \frac{a_i}{\Delta u_i} u_i \\ &= b_o + \sum b_i u_i \end{aligned}$$