



Les ajustements linéaires multiples :

Règles de bonne pratique

par Bernard Beauzamy, avec la collaboration d'Alisson Stochetti.

Octobre 2018

Cet article est rédigé à l'attention des Autorités de Contrôle, des Autorités de Sûreté, des Bureaux de Vérification, etc., pour leur permettre de juger de la qualité des arguments mathématiques invoqués lors d'une démonstration de sûreté.

Avertissement :

Les mathématiques sont là pour décrire les lois de la Nature, et en aucune manière pour apporter des outils factices destinés à des démonstrations factices, qui seront rejetées avec violence et dont les auteurs seront châtiés avec sévérité.

I. Présentation du problème

Etant donnée une "variable d'intérêt", généralement notée Y , et des "paramètres explicatifs", généralement notés X_1, \dots, X_n , on appelle ajustement linéaire de Y par X_1, \dots, X_n une représentation de la forme :

$$Y = a_1 X_1 + \dots + a_n X_n + \text{reste} \quad (1)$$

Il s'agit d'un ajustement "multiple" si $n \geq 2$. On parle aussi quelquefois de "régression", ce mot ayant été utilisé dans le cadre d'études très anciennes.

Le besoin de réaliser un ajustement linéaire se rencontre souvent dans les applications :

- Industrielles : par exemple Y désigne les propriétés d'un matériau (comme une résistance élastique) et les X_j représentent des paramètres mesurés lors du process de fabrication (des températures, des pressions, etc.) ;
- Environnementales : Y désigne la concentration en un certain polluant et les X_j représentent différents paramètres mesurés : la température, le débit du fleuve, la vitesse du vent, le nombre de véhicules qui passent, la densité de population, etc. ;
- Epidémiologiques : Y peut désigner le nombre de personnes atteintes d'une maladie, et les X_j seront divers facteurs explicatifs (âge, expositions diverses, etc.) ;
- Financières : Y peut désigner la performance d'un indice ou le rendement d'un portefeuille, et les X_j seront divers facteurs explicatifs associés (taux de change, niveau de la dette, etc.).

De manière générale, il s'agit de ramener une variable que l'on ne connaît pas, ou que l'on connaît mal, à savoir Y , à une combinaison linéaire de variables que l'on connaît mieux (les X_j) : c'est donc une idée complètement naturelle. On peut ensuite chercher à "hiérarchiser" les paramètres : parmi les X_j , lequel a la plus grande influence sur Y ? Cette question est beaucoup plus floue.

En théorie, il pourrait se faire que Y et les X_j soient données sous forme "abstraite" (sous la forme d'une fonction parfaitement définie), mais ce n'est jamais le cas en pratique. En pratique, il s'agit de relevés faits sur diverses expériences. On a procédé à N expériences, et, pour chacune d'elles, on a relevé les valeurs de Y , notées y_i et des X_j , notées $x_{i,j}$, $i = 1, \dots, N$. On dispose donc d'un tableau à double entrée, du type suivant :

Y	X_1	\dots	X_n
y_1	$x_{1,1}$	\dots	$x_{1,n}$
\vdots	\vdots	\dots	\vdots
y_N	$x_{N,1}$	\dots	$x_{N,n}$

Donc – et il est important de bien comprendre ceci – le problème de l'ajustement linéaire ne concerne pas des variables, plus ou moins bien définies, mais simplement la première colonne d'un tableau, que l'on cherche à représenter comme combinaison linéaire des n autres (avec cette représentation, le tableau a N lignes et $n+1$ colonnes). Un tel problème n'a rien de probabiliste : il sera traité avec des outils d'analyse. Nous dirons que Y et les X_j sont des "listes", ceci afin d'éviter toute confusion avec le vocabulaire probabiliste.

II. Résumé des recommandations

1. Un peu de bon sens

Vérifier que les données sont en nombre suffisant pour ne pas dire de sottises. Si l'on dispose de trois listes de paramètres, trois expériences suffiront pour déterminer (a_1, a_2, a_3) ; prenons le tableau suivant (âge, taille, poids) :

âge	taille	poids
15	1,6	53
20	1,8	70
60	1,75	85

Trois expériences ont été faites et on pourrait ajuster n'importe quelle variable Y (par exemple la sensibilité à un vaccin) aux trois paramètres âge, taille, poids, et en déduire, grâce à cet ajustement, la valeur de Y pour n'importe quelle valeur de (X_1, X_2, X_3) . Or ce serait absurde : on ne sait même pas si ces trois expériences concernent des hommes ou des femmes, et il n'y a aucun test pour la tranche d'âge entre 20 et 60 ans !

2. Avant d'utiliser Excel

- Vérifier que les listes X_1, \dots, X_n sont linéairement indépendantes ;
- Remplacer Y et les X_j par des listes centrées réduites : Y est remplacé par $Y = (Y - m_Y) / \sigma_Y$, où m_Y , σ_Y désignent respectivement la moyenne et l'écart-type de la liste. Faire de même avec les X_j , remplacés par les variables réduites X_j .

3. Utilisation d'Excel

Une fois que Excel fournit les coefficients a_1, \dots, a_n de l'ajustement et une estimation de $\varepsilon = \|reste\|_2$, bien se souvenir que :

- On n'a pas $y_i = a_1 x_{i,1} + \dots + a_n x_{i,n}$ pour tout i , mais seulement $\left| y_i - (a_1 x_{i,1} + \dots + a_n x_{i,n}) \right| < \varepsilon \sqrt{N}$, estimation qui peut être très mauvaise si N est grand ;
- Cette estimation ne porte que sur la liste normalisée Y ; si l'on veut revenir à la liste d'origine Y , on aura, avec d'autres coefficients b_1, \dots, b_n , $\left| y_i - (b_1 x_{i,1} + \dots + b_n x_{i,n}) \right| < \varepsilon \sqrt{N} \sigma_Y$.

4. Extrapolation des résultats

Bien entendu, l'écart mesuré entre Y et la combinaison $a_1 X_1 + \dots + a_n X_n$ ne porte que sur les listes elles-mêmes. Il est incorrect de prétendre que cet écart sera maintenu à l'identique si on attribue à X_1, \dots, X_n des valeurs qui n'ont rien à voir avec celles de l'expérience.

5. En conclusion de nos recommandations

Vouloir remplacer Y par une combinaison linéaire $a_1X_1 + \dots + a_nX_n$ peut sembler raisonnable a priori, mais le choix d'une dépendance linéaire est artificiel. Les résultats obtenus seront faux dans 80% des cas et les démonstrations fausses dans 99% des cas. Il est bien préférable de rechercher (au moyen de lois de probabilités conditionnelles) la dépendance de Y par rapport à chacune des variables X_j .

Pour bien faire comprendre ceci, prenons l'exemple de données environnementales. Une pollution Y , mesurée en grammes par litre, peut parfaitement dépendre linéairement de la population (mesurée en nombre d'habitants), du moins dans une certaine plage, et peut parfaitement dépendre linéairement du débit d'un fleuve (mesuré en mètres cubes par seconde), mais ne peut en aucune manière s'écrire comme la somme $a_1X_1 + a_2X_2$, parce que l'on ne peut pas ajouter un nombre d'habitants à des m^3/s .

Les mathématiques sont là pour décrire les lois de la Nature, et en aucune manière pour apporter des outils factices destinés à des démonstrations factices, qui seront rejetées avec violence et dont les auteurs seront châtiés avec sévérité.

III. Approche mathématique

On peut évidemment, pour chaque expérience, remplacer y_i par une combinaison linéaire $a_1x_{i,1} + \dots + a_nx_{i,n}$ et il y a une infinité de manières de le faire, puisque l'on a une seule équation et n coefficients à déterminer (les a_1, \dots, a_n), mais, en procédant ainsi, les coefficients ainsi choisis seront différents d'une expérience à l'autre. Ce n'est pas ce que l'on veut : on voudrait qu'ils soient les mêmes pour toutes les expériences. Mais alors on ne peut espérer "tomber juste" et avoir pour tout $i = 1, \dots, N$: $y_i = a_1x_{i,1} + \dots + a_nx_{i,n}$. Tout au plus aurons-nous :

$$y_i = a_1x_{i,1} + \dots + a_nx_{i,n} + r_i \quad (2)$$

où les r_i seront "les plus petits possible", en un sens à déterminer. En pratique, trois significations sont possibles :

1. On exige que tous les restes soient petits. Cela signifie que les coefficients a_1, \dots, a_n seront choisis pour que :

$$\max_i \left| y_i - (a_1x_{i,1} + \dots + a_nx_{i,n}) \right|$$

soit le plus petit possible. Mais ceci est peu réaliste, parce que, en pratique, il peut se faire que plusieurs expériences soient ratées : ou bien y_i , ou bien l'un des $x_{i,j}$ ont été mal mesurés. Lorsqu'on le sait, on élimine cette expérience-là, ou ces expériences-là, s'il y en a plusieurs.

Une variante, que l'on rencontre quelquefois dans les questions liées à l'environnement, est que l'on autorise l'industriel à éliminer quelques expériences.

2. On exige que la moyenne des restes soit petite. Cela signifie que les coefficients a_1, \dots, a_n seront choisis pour que :

$$\frac{1}{N} \sum_{i=1}^N |y_i - (a_1 x_{i,1} + \dots + a_n x_{i,n})|$$

soit le plus petit possible.

3. On exige que la moyenne quadratique des restes soit petite. Cela signifie que les coefficients a_1, \dots, a_n seront choisis pour que :

$$\frac{1}{N} \sum_{i=1}^N (y_i - (a_1 x_{i,1} + \dots + a_n x_{i,n}))^2$$

soit le plus petit possible.

Rappelons (voir par exemple [BB_Banach]) que pour des nombres réels $\alpha_1, \dots, \alpha_N$ on a toujours :

$$\frac{1}{N} \sum_{i=1}^N |\alpha_i| \leq \left(\frac{1}{N} \sum_{i=1}^N \alpha_i^2 \right)^{1/2} \leq \max_i |\alpha_i| \quad (3)$$

Par conséquent, contrôler le troisième est plus fort que contrôler le second, lui-même plus fort que contrôler le premier. Comme le troisième est exclu pour des raisons pratiques, on préfère le second, d'autant que les outils disponibles sont plus faciles à mettre en œuvre que pour le premier.

Mais ce choix aura des conséquences négatives inattendues, comme nous le verrons par la suite.

IV. Résolution mathématique

Concrètement, le problème de l'ajustement linéaire revient donc à chercher les coefficients a_1, \dots, a_n qui minimisent :

$$Q(a_1, \dots, a_n) = \frac{1}{N} \sum_{i=1}^N (y_i - (a_1 x_{i,1} + \dots + a_n x_{i,n}))^2 \quad (4)$$

La fonction Q est convexe (voir [BB_Banach]), ce qui signifie qu'elle ressemble à une parabole pour chacune des variables. Il en résulte que l'on est sûr qu'il existe un minimum unique.

Il existe trois méthodes pratiques pour résoudre le problème (4) :

6. Algorithme de minimisation

On commence avec une valeur arbitraire $(a_1^{(0)}, \dots, a_n^{(0)})$ et on la déplace de manière à diminuer Q à chaque étape. Un tel algorithme est facile à mettre en œuvre, mais on n'est jamais sûr qu'il s'arrête exactement à la bonne valeur : tout dépend en particulier des erreurs d'arrondi.

7. Calcul matriciel

Pour résoudre le problème (4) :

$$Q(a_1, \dots, a_n) = \sum_{i=1}^N \left(y_i - \sum_{j=1}^n a_j x_{i,j} \right)^2 \rightarrow \min$$

on peut procéder comme suit. Calculons les dérivées partielles :

$$\frac{\partial Q}{\partial a_k}(a_1, \dots, a_n) = -2 \sum_{i=1}^N x_{i,k} \left(y_i - \sum_{j=1}^n a_j x_{i,j} \right)$$

Par conséquent, $\frac{\partial Q}{\partial a_j}(a_1, \dots, a_n) = 0$ donne :

$$\sum_{i=1}^N x_{i,k} \sum_{j=1}^n a_j x_{i,j} = \sum_{i=1}^N x_{i,k} y_i \quad (5)$$

Ceci a une interprétation analytique. Considérons les listes Y, X_j comme des points dans l'espace vectoriel euclidien \mathbb{R}^N et introduisons les produits scalaires :

$$\langle X_j, X_k \rangle = \sum_{i=1}^N x_{i,j} x_{i,k} \quad (6)$$

et :

$$\langle X_j, Y \rangle = \sum_{i=1}^N x_{i,j} y_i \quad (7)$$

Les équations (5) deviennent, pour $k = 1, \dots, n$:

$$\sum_{j=1}^n a_j \langle X_j, X_k \rangle = \langle X_k, Y \rangle \quad (8)$$

Soit A le vecteur colonne des (a_1, \dots, a_n) , soit M la matrice $n \times n$ faite des $\langle X_i, X_j \rangle$, et soit B le vecteur colonne des $\langle X_j, Y \rangle$; le système d'équations peut se mettre sous forme matricielle :

$$MA = B \tag{9}$$

et la solution est $A = M^{-1}B$. Cette approche est très satisfaisante et elle retourne une valeur exacte, mais elle demande l'inversion d'une matrice. Comme celle-ci est symétrique et que, en pratique, la dimension satisfait presque toujours $n \leq 10$, l'inversion se fait facilement. L'inconvénient de la méthode est qu'il s'agit d'une "boîte noire" : elle permet d'obtenir la solution, sans que l'on comprenne bien comment elle est apparue. C'est pourquoi la présentation qui suit est indispensable.

8. Approche géométrique

Une approche géométrique permet d'avoir une bien meilleure compréhension du problème. Là encore, nous considérons que les listes Y et X_j représentent des points dans l'espace euclidien \mathbb{R}^N de dimension N (nombre total d'expériences). Le produit scalaire $\langle X, Y \rangle = \sum_{i=1}^N x_i y_i$ correspond à la norme euclidienne :

$$\|X\|_2 = \sqrt{\langle X, X \rangle} = \sqrt{\sum_{i=1}^N x_i^2}.$$

Voir Annexe pour les propriétés de cette norme. Soit $F = \text{span}\{X_1, \dots, X_n\}$ le sous-espace vectoriel de \mathbb{R}^N engendré par les vecteurs X_1, \dots, X_n ; par définition, F est l'ensemble des combinaisons $\sum_{i=1}^n \alpha_i X_i$, pour tous les réels $(\alpha_1, \dots, \alpha_n)$; la dimension de F est $\leq n$ (exactement n si les vecteurs X_1, \dots, X_n sont linéairement indépendants, ce qui est le cas en général ; nous en reparlons plus bas).

Soit P la projection orthogonale de Y sur F ; le point P est la solution du problème (4), et sa décomposition $P = \sum_{i=1}^n a_i X_i$ donne la solution cherchée : par définition, le point P , projection linéaire, est le plus proche de Y , parmi tous les points de F .

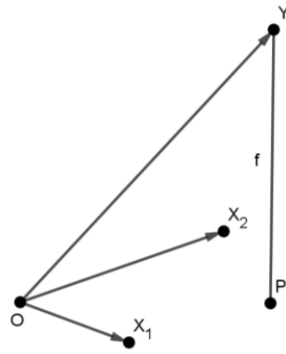


Fig. 1 : projection sur un plan

Sur la figure 1, on voit un cas particulier : il n'y a que deux vecteurs X_1, X_2 et P est la projection linéaire de Y sur le plan qu'ils engendrent.

Cette façon de voir les choses ne résout pas le problème numériquement (reste à calculer les coordonnées de P), mais elle est essentielle pour bien comprendre le problème. On constate en effet que le problème de minimisation (4) se réduit à un problème de projection orthogonale dans un espace euclidien. On a donc deux problèmes distincts :

- Trouver la position de P dans le sous-espace F ;
- Décomposer P en se servant des X_1, \dots, X_n .

Pour le premier, nous notons que la position du point P ne dépend que de celle de Y et du s.e.v F ; elle ne sera pas modifiée si, par exemple, les vecteurs X_1, \dots, X_n sont permutés, ou remplacés par des combinaisons linéaires (par exemple remplacer X_2 par $\frac{X_1 + X_2}{2}$) ; nous notons aussi que la position de P est unique : il n'existe qu'une seule projection orthogonale d'un point sur un sous-espace vectoriel. Soit B la boule de centre Y et de rayon $r = \text{dist}(Y, P) = \|Y - P\|_2$; elle est tangente en P au sous-espace F ; voir figure :

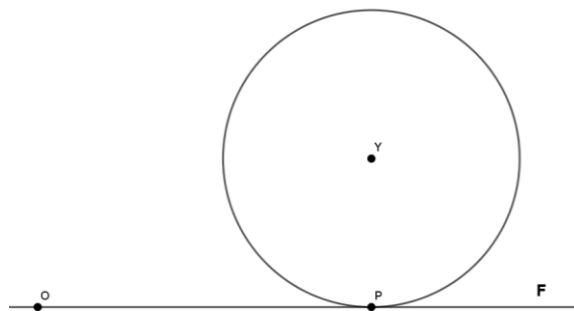


Fig. 2 : la boule est tangente au sous-espace vectoriel

Pour le second problème, c'est plus compliqué. Il faut d'abord vérifier que les vecteurs X_1, \dots, X_n sont linéairement indépendants. Formellement, cela signifie que la seule combinaison linéaire $\lambda_1 X_1 + \dots + \lambda_n X_n = 0$ est celle dont tous les coefficients sont nuls. Il existe des routines sous Excel pour calculer le rang d'un ensemble de vecteurs : ils sont indépendants si le

rang est égal au nombre de vecteurs (noté ici n). Si le rang est $n' < n$, cela signifie que n' vecteurs suffisent à reconstituer l'ensemble.

Attention : il faut bien distinguer ici entre "linéairement indépendants" et "stochastiquement indépendants" : la confusion est possible parce que nous avons parlé de variables aléatoires, mais aussi de vecteurs. Linéairement indépendants est une propriété très faible, qui signifie simplement qu'aucun vecteur ne peut être reconstitué à partir d'une combinaison linéaire des autres. Stochastiquement indépendants est une propriété très forte, qui signifie que toute connaissance sur $n-1$ d'entre eux ne donne rien sur le dernier. Par exemple, X et X^2 ne sont certainement pas stochastiquement indépendants, parce que si l'on connaît X on connaît X^2 , par contre, ils sont en général linéairement indépendants. Dans ce paragraphe, nous parlons d'indépendance linéaire.

En général, dans la pratique, les vecteurs X_1, \dots, X_n sont linéairement indépendants, parce que le nombre d'observations est grand devant le nombre de vecteurs. Mais il peut arriver (artificiellement) que par exemple $X_2 = \frac{X_1 + X_3}{2}$: la seconde mesure est toujours la moyenne entre la première et la troisième, auquel cas on l'élimine : elle ne sert à rien.

Nous supposons dans la suite que les vecteurs X_1, \dots, X_n sont linéairement indépendants : les éventuels vecteurs dépendants ont été éliminés.

Dans ces conditions, ces vecteurs forment une base de F et le vecteur P se décompose de manière unique en une combinaison linéaire :

$$P = a_1 X_1 + \dots + a_n X_n \quad (10)$$

et les coefficients a_1, \dots, a_n sont les solutions du problème (4).

Les vecteurs qui sont donnés, X_1, \dots, X_n , correspondent à des mesures physiques (par exemple des températures en tel point du processus de fabrication). Il est possible de les normer (voir plus loin), mais on ne peut pas les remplacer par des combinaisons linéaires, qui n'auraient plus de sens physique. On ne peut donc pas, en général, engendrer F au moyen d'une base orthogonale (deux vecteurs quelconques ayant un produit scalaire nul).

V. Précautions préliminaires

A. Il faut normaliser les vecteurs

Si on opère directement sur les vecteurs X_j , l'espace doit avoir une structure vectorielle : il faut pouvoir additionner deux vecteurs (puisqu'on écrit $a_1 X_1 + a_2 X_2 + \dots$). Les colonnes doivent donc être d'une nature qui permet l'addition. Ceci a un sens si, par exemple, il s'agit constamment de longueurs, de poids, de volumes, etc., mais n'en a aucun s'il s'agit de concepts de nature différente. Par exemple, si X_1 est une concentration exprimée en grammes par litre et

X_2 une température, la somme $X_1 + X_2$ n'a aucun sens. Or, si on se contente de remplir des colonnes de chiffres (par exemple $X_1 = 0.1$ et $X_2 = 56$), la somme a toujours un sens.

Lorsqu'on utilise "à l'aveugle" un tableau Excel, on le remplit avec des nombres, sans se soucier de leur signification. Mais dans ces conditions toutes les opérations sont possibles : on peut toujours faire la somme, le produit, le quotient, de deux nombres, même si ces opérations n'ont aucun sens sur les données que ces nombres représentent. Excel n'émet aucun "warning" : c'est à l'utilisateur de le savoir et de prendre ses précautions.

Si on a au préalable normalisé les vecteurs, cette objection disparaît : les $\tilde{X}_j = \frac{X_j}{\|X_j\|_2}$, avec

$\|X_j\|_2 = \sqrt{\sum_{i=1}^N x_{i,j}^2}$, sont en effet des nombres sans dimension, invariants par changement

d'échelle : il est indifférent, par exemple, qu'ils soient mesurés en grammes par litre ou en kg par m³. On peut accepter l'idée de faire des combinaisons linéaires $a_1\tilde{X}_1 + \dots + a_n\tilde{X}_n$ et de chercher si elles sont proches de Y , qui aura lui aussi été normé. Bien entendu, ces combinaisons linéaires n'ont pas de sens physique.

L'opération de normalisation, que les vecteurs aient une signification commune ou non, est de toute façon nécessaire si l'on cherche à comparer entre eux les coefficients de la décomposition (en particulier pour voir lequel est le plus grand).

En effet, dans (10), si X_1 est multiplié par 2, le coefficient correspondant a_1 sera divisé par 2. On ne peut donc comparer les coefficients que si l'on s'est au préalable assuré que les X_j avaient tous la même taille, ici une longueur égale à 1. Dans ces conditions, les $\tilde{X}_1, \dots, \tilde{X}_n$ forment une base normée (mais non orthonormée).

Une fois que les vecteurs ont été normés, la décomposition :

$$P = a_1\tilde{X}_1 + \dots + a_n\tilde{X}_n \quad (11)$$

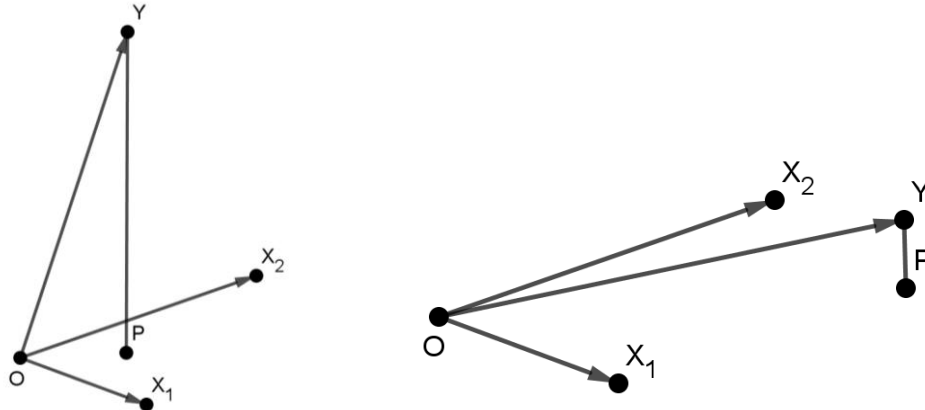
présente un intérêt particulier : le plus grand des coefficients a_j indique la plus grande dépendance par rapport aux paramètres. Si par exemple $a_1 > \dots > a_n$, cela signifie que P "dépend" plus de X_1 que de X_2, \dots, X_n .

Si les vecteurs X_1, \dots, X_n ne sont pas indépendants, la décomposition (10) n'est pas unique et on s'expose à de graves ennuis. L'application aveugle des routines de Excel ne permettra pas de s'en apercevoir. On aura par exemple l'impression d'avoir deux variables prépondérantes, alors que c'est la même.

Si on a normalisé les X_j , il faut également normaliser Y . On ne peut pas avoir une égalité du type $Y = a_1\tilde{X}_1 + \dots + a_n\tilde{X}_n + \text{reste}$ si les \tilde{X}_j sont des nombres sans dimension et Y est exprimé dans une unité quelconque, par exemple des grammes par litre. La formule doit être cohérente du point de vue de la physique.

B. Analyse de la qualité du résultat

La projection de Y sur P représente nécessairement une perte d'information, mais cette perte est d'autant plus faible que Y est proche de P (et la perte est nulle si $Y = P$).



Dans le premier cas, les points Y et P sont éloignés ; dans le second ils sont proches.

Ceci peut être mesuré par le quotient $\rho = \frac{\|P\|_2}{\|Y\|_2}$, qui est une mesure de la qualité de l'ajustement. Le triangle OPY est rectangle en P et OY est l'hypoténuse : la longueur OP est inférieure à la longueur OY et on a toujours $0 \leq \rho \leq 1$. On a $\rho = 0$ si le point P est à l'origine (la projection de Y est en O) et on a $\rho = 1$ si $Y = P$.

La liste Y représente les enregistrements d'une expérience qui a été répétée : la valeur y_i , $i = 1, \dots, N$ est le résultat de la $i^{\text{ème}}$ expérience. Il y a donc habituellement une "dispersion" des résultats : toutes les expériences ne donnent pas le même résultat. On notera :

$$m_Y = \frac{1}{N} \sum_{i=1}^N y_i \quad : \text{moyenne des résultats}$$

$$v_Y = \frac{1}{N} \sum_{i=1}^N y_i^2 - \left(\frac{1}{N} \sum_{i=1}^N y_i \right)^2 \quad : \text{variance des résultats}$$

$$\sigma_Y = \sqrt{v_Y} \quad : \text{écart-type des résultats.}$$

Par projection linéaire sur le s.e.v. F , la moyenne de Y va se projeter en la moyenne de P : ceci est évident parce que la projection est linéaire. On a $m_Y = m_P + m_R$.

L'écart-type de la variable P est toujours inférieur (ou égal) à celui de la variable Y : ceci résulte du fait qu'une projection ne peut que "contracter" les données.

Voici une démonstration abrégée de ce fait. Nous voulons montrer que

$$m(P - m_p)^2 \leq m(Y - m_Y)^2$$

Il suffit de montrer que, pour tout Y , en notant $P = \pi(Y)$ la projection :

$$m(\pi(Y))^2 \leq m(Y)^2 : \text{on appliquera ce résultat à } Y - m_Y.$$

Ecrivons $Y = P + R$ où R est le reste, orthogonal à P . On a (voir Annexe) :

$$m(Y)^2 = \frac{1}{N} \|Y\|_2^2 = \frac{1}{N} \|P\|_2^2 + \frac{1}{N} \|R\|_2^2 \geq \frac{1}{N} \|P\|_2^2 = m(P)^2, \text{ ce qui prouve notre assertion.}$$

On peut donc se poser la question : la dispersion des résultats de P est-elle proche ou non de celle de Y ? Pour répondre à cette question, calculons la variance de Y en fonction de celles de P et R , supposés orthogonaux. On a :

$$v(Y) = m(Y - m(Y))^2 = m(Y^2) - m(Y)^2$$

Or, puisque P et R sont orthogonaux :

$$m(Y^2) = m(P + R)^2 = m(P^2) + m(R^2)$$

$$m_Y^2 = (m_p + m_r)^2 = m_p^2 + m_r^2 + 2m_p m_r$$

et par conséquent :

$$v(Y) = m(Y^2) - m_Y^2 = m(P^2) + m(R^2) - (m_p^2 + m_r^2 + 2m_p m_r)$$

D'où la formule générale, valable dès que les listes P et R sont orthogonales dans l'espace euclidien :

$$v(Y) = v(P) + v(R) - 2m_p m_r \tag{1}$$

Attention. – On n'a pas nécessairement $v(Y) = v(P) + v(R)$ si les listes sont orthogonales.

Voici un exemple en dimension 3.

$Y = (1, 3, -4)$, $m_Y = 0$, $V(Y) = \frac{26}{3}$. On projette sur le plan $z=0$; on a donc $P = (1, 3, 0)$ et $R = (0, 0, -4)$, d'où $m_P = \frac{4}{3}$, $V(P) = \frac{14}{9}$, $m_R = -\frac{4}{3}$, $V(R) = \frac{32}{9}$ et $\frac{26}{3} = \frac{78}{9} \neq \frac{14}{9} + \frac{32}{9} = \frac{46}{9}$

Corollaire. – Soient deux listes P et R orthogonales et $Y = P + R$ leur somme ; on a :

$$v(Y) = v(P) + v(R)$$

si et seulement si l'une des listes P ou R est centrée (c'est-à-dire $m_P = 0$ ou $m_R = 0$).

Revenons au problème de l'ajustement : nous disposons d'une liste Y , projetée sur le s.e.v. F de \mathbb{R}^N engendré par les listes X_1, \dots, X_n . Donc, par définition de la projection, on a, pour $j = 1, \dots, n$:

$$\langle Y, X_j \rangle = \langle P, X_j \rangle$$

Si la liste X_1 est faite constamment de 1, c'est-à-dire $X_1 = (1, \dots, 1)$, alors

$$\langle Y, X_1 \rangle = \sum_{i=1}^N y_i = Nm_Y$$

On en déduit le résultat suivant :

Théorème. – Si la liste Y a été centrée (avant tout ajustement) et si la liste X_1 est constante, alors, pour tout ajustement :

$$v(Y) = v(P) + v(R)$$

En effet, puisque $m_Y = 0$, on a $\langle Y, X_1 \rangle = 0$ et donc $\langle P, X_1 \rangle = 0$ et $m_P = 0$; le théorème résulte alors du Corollaire.

VI. Que déduire de l'ajustement ?

L'idée est évidemment de reconstituer Y à partir des X_j en un endroit inconnu, c'est-à-dire pour lequel aucune expérience n'a été faite. Si on écrit

$$Y \approx a_1 X_1 + \dots + a_n X_n$$

on peut être tenté d'utiliser cette égalité approximative pour calculer Y à partir de n'importe quelle valeur donnée aux X_j : c'est ce que font beaucoup d'études. Mais cette approche est incorrecte ; voyons ceci en détail.

Admettons que l'on ait obtenu par ajustement :

$$\|Y - P\|_2 < \varepsilon$$

cela signifie par définition :

$$\frac{1}{N} \sum_{i=1}^N (y_i - p_i)^2 < \varepsilon^2$$

et donc :

$$\sum_{i=1}^N (y_i - p_i)^2 < N\varepsilon^2$$

D'où il résulte que :

$$\max_i |y_i - p_i| < \sqrt{N}\varepsilon$$

C'est ici que l'on voit une conséquence fâcheuse du choix de la norme $\| \cdot \|_2$, bien pratique par ailleurs : elle contrôle la somme des carrés, mais donne une mauvaise estimation sur chacun des termes (du fait du facteur \sqrt{N}).

Voyons un exemple pour illustrer ceci. Générons pour X_1 et X_2 $N=1000$ nombres aléatoires entre 0 et 1, normalisons-les, et posons $Y = 3X_1 - 5X_2$; Y est évidemment une combinaison linéaire parfaite de X_1, X_2 et l'ajustement est idéal. Maintenant, modifions simplement les deux premières lignes de Y : on ajoute 50 pour la première et on retranche 50 pour la seconde, ce qui fait que la moyenne reste nulle :

Y modifié	X1 normalisé	X2 normalisé
51,1833677	0,45853685	0,03844857
-58,0565253	-0,60694961	1,24713531
-10,7243127	-1,5176585	1,23426744
-0,63298703	0,20090345	0,24713947

L'ajustement reste très bon (coefficient de détermination multiple égal à 0.93) parce qu'il y a 1000 lignes et que la modification de deux d'entre elles intervient peu. Pourtant, le reste sur la première ligne est maintenant de 49.98 : l'approximation est très mauvaise.

Il y a un second élément à prendre en compte : nous n'avons pas travaillé sur Y et les X_j , mais sur leur normalisation. Si nous avons :

$$\left| \frac{y}{\sigma(Y)} - \frac{p}{\sigma(P)} \right| < \varepsilon$$

il en résulte :

$$\frac{p}{\sigma(P)} - \varepsilon < \frac{y}{\sigma(Y)} < \frac{p}{\sigma(P)} + \varepsilon$$

ou encore :

$$\frac{\sigma(Y)}{\sigma(P)} p - \varepsilon \sigma(Y) < y < \frac{\sigma(Y)}{\sigma(P)} p + \varepsilon \sigma(Y)$$

Autrement dit, l'écart-type $\sigma(Y)$ va intervenir dans l'intervalle final pour l'approximation. Si la taille du reste était 0.1 lorsqu'on a fait l'ajustement sur liste normalisée, et si $\sigma(Y) = 50$, la taille du reste sera 5 pour la liste elle-même.

Bien entendu, l'écart mesuré entre Y et la combinaison $a_1 X_1 + \dots + a_n X_n$ ne porte que sur les listes elles-mêmes. Il est incorrect de prétendre que cet écart sera maintenu à l'identique si on attribue à X_1, \dots, X_n des valeurs qui n'ont rien à voir avec celles de l'expérience.

VII. Mise en œuvre pratique

Excel permet de calculer un ajustement linéaire multiple de deux manières :

- Tout d'abord en utilisant directement une formule dans le tableau Excel, sous la forme :
=DROITEREG(E2:E39;F2:G39;VRAI;FAUX)

Dans cet exemple, les données de Y sont dans la plage $E2 - E39$ et celles des X_j dans la plage $F2 - G39$; Excel fait l'hypothèse que la première variable explicative est une constante ($x_{i,1} = 1$ pour tout i).

- En utilisant un complément Excel appelé "utilitaire d'analyse", "régression linéaire", qui laisse choisir directement les plages pour les différentes variables (ce qui est plus simple à l'usage). Voici le résultat que donne cet utilitaire (process de fabrication, 38 expériences, 9 paramètres explicatifs) :

RAPPORT DÉTAILLÉ								
Statistiques de la régression								
Coefficient de détermination multiple	0,967909946							
Coefficient de détermination R^2	0,936849663							
Coefficient de détermination R^2	0,916551341							
Erreur-type	9,203574871							
Observations	38							
ANALYSE DE VARIANCE								
	Degré de liberté	Somme des carrés	Moyenne des carrés	F	Valeur critique de F			
Régression	9	35185,6326	3909,514734	46,15404349	1,9169E-14			
Résidus	28	2371,76213	84,7057904					
Total	37	37557,3947						
	Coefficients	Erreur-type	Statistique t	Probabilité	Limite inférieure pour seuil de confiance = 95%	Limite supérieure pour seuil de confiance = 95%	Limite inférieure pour seuil de confiance = 95,0%	Limite supérieure pour seuil de confiance = 95,0%
Constante	259,4022266	16,6291499	15,59924763	2,44107E-15	225,338957	293,465496	225,338957	293,465496
Variable X 1	0,979244504	0,10973208	8,923958129	1,11609E-09	0,75446852	1,20402048	0,75446852	1,20402048
Variable X 2	0,567865346	2,96130595	0,191761796	0,849312696	-5,49809491	6,6338256	-5,49809491	6,6338256
Variable X 3	18,86759666	4,59650765	4,104767818	0,000317099	9,45207757	28,2831158	9,45207757	28,2831158
Variable X 4	32,46238128	64,7949137	0,501001999	0,62028567	-100,263983	165,188745	-100,263983	165,188745
Variable X 5	-45,69034937	62,2626231	-0,733832709	0,469151339	-173,229551	81,8488524	-173,229551	81,8488524
Variable X 6	35,75731751	80,1889442	0,445913309	0,659089024	-128,502289	200,016924	-128,502289	200,016924
Variable X 7	1,831365109	13,4073847	0,136593762	0,892329317	-25,6324174	29,2951477	-25,6324174	29,2951477
Variable X 8	0,00918234	0,00110059	8,343078622	4,46403E-09	0,00692788	0,0114368	0,00692788	0,0114368
Variable X 9	-5,580166069	4,29830844	-1,298223742	0,204802511	-14,3848518	3,22451964	-14,3848518	3,22451964

Ici, le lecteur va être surpris, parce que nous avons dit, et démontré, que le problème de l'ajustement était un problème d'analyse (projection sur un sous-espace vectoriel). Or toutes les annonces de Excel sont données en un vocabulaire probabiliste. Essayons de le comprendre et de voir s'il est pertinent.

Notons tout d'abord que l'utilisation d'Excel se fait sans aucune précaution : vous pouvez mettre absolument n'importe quoi dans chacune des colonnes, Excel renverra le tableau de résultats ci-dessus.

Statistiques de la régression :

Ligne 1 : Coefficient de détermination multiple, encore appelé coefficient de corrélation

C'est la racine carrée du coefficient de détermination R^2 donné ligne 2

(il est vraiment absurde que le même nom, à savoir "coefficient de détermination", soit utilisé à la fois pour une valeur x et son carré x^2 !)

Ligne 2 : Coefficient de détermination R^2

Il peut être défini comme $R^2 = \frac{v(P)}{v(Y)}$; il est toujours compris entre 0 et 1 (si les hypothèses de

traitement sont satisfaites) et vaut 1 si $v(R) = 0$.

Ligne 3 : Coefficient de détermination R^2

Le nom est le même, mais la valeur est différente. Il s'agit ici (Excel ne le dit pas) de ????

Ligne 4 : Erreur type

Pour Excel, c'est $\sqrt{\frac{1}{N-n} \sum_{i=1}^N r_i^2}$

On lit ensuite :

Somme des carrés : expression de la variabilité

SCE : somme des carrés expliqués. Exprime la variabilité expliquée, c'est-à-dire la variation que le paramètre explique. Il est défini par :

$$SCE = \|P - m_y\|_2^2 = \sum_{i=1}^N (p_i - m_y)^2$$

SCR : somme des carrés résiduels. Exprime la variabilité résiduelle, c'est-à-dire la variation que le paramètre n'arrive pas à expliquer. Il est défini par :

$$SCR = \|Y - P\|_2^2 = \sum_{i=1}^N (y_i - p_i)^2$$

SCT : somme des carrés totaux. Exprime la variabilité totale des observations. Il est défini par :

$$SCT = SCE + SCR = \|Y - m_y\|_2^2 = \sum_{i=1}^N (y_i - m_y)^2$$

Ici, Excel affirme que $SCT = SCE + SCR$; pour que ce soit vrai, il faut que $P - m_y$ et $Y - P$ sont orthogonaux : les normes au carré vont s'ajouter. Or :

$$\langle P - m_y, Y - P \rangle = \langle P - m_y, R \rangle = \langle P, R \rangle - \langle m_y, R \rangle$$

On a toujours, par construction, $\langle P, R \rangle = 0$; la formule d'additivité des variances ne sera donc exacte que si (et seulement si) $\langle m_y, R \rangle = 0$. Or $\langle m_y, R \rangle = \sum_{i=1}^N m_y r_i = N m_y m_R$; il faut donc que $m_y = 0$ ou $m_R = 0$; encore une fois, pour être valables, ces formules supposent que Y a préalablement été centrée, d'où nos recommandations initiales.

La colonne "coefficients" du second tableau donne les a_1, \dots, a_n ; il vaut mieux s'en tenir là et oublier tout le reste. On calculera directement les valeurs $r_i = y_i - (a_1 x_{i,1} + \dots + a_n x_{i,n})$, la somme des r_i^2 , etc.

Annexe

Notations probabilistes et notations dans l'espace euclidien

Une variable "aléatoire" Y , caractérisant N observations y_i , peut être vue d'un point de vue probabiliste :

$$\text{La moyenne est } EY = \frac{1}{N} \sum_{i=1}^N y_i$$

$$\text{De même, } EY^2 = \frac{1}{N} \sum_{i=1}^N y_i^2$$

$$\text{L'écart-type est } \sigma(Y) = \sqrt{E(Y - E(Y))^2} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - m_Y)^2}, \text{ avec } m_Y = \frac{1}{N} \sum_{i=1}^N y_i$$

Dans l'espace euclidien \mathbb{R}^N , Y peut être vu comme un vecteur ayant N composantes, les y_i , et la norme euclidienne est :

$$\|Y\|_2 = \sqrt{\sum_{i=1}^N y_i^2}$$

et par conséquent :

$$\|Y\|_2^2 = \sum_{i=1}^N y_i^2 = N \cdot EY^2$$

et de même :

$$\sigma_y = \frac{1}{\sqrt{N}} \|Y - m_Y\|_2$$

Considérons deux listes X, Y : $X = (x_i)_{i=1, \dots, N}$; $Y = (y_i)_{i=1, \dots, N}$ correspondant aux relevés de deux paramètres dans une même expérience. On définit le produit scalaire :

$$\langle X, Y \rangle = \sum_{i=1}^N x_i y_i$$

$$\text{et } \|X\|_2^2 = \langle X, X \rangle$$

$$\text{Par conséquent } E(XY) = \frac{1}{N} \langle X, Y \rangle$$

Notons 1 la variable dont les N occurrences sont toutes égales à 1 (variable constante, donc).
On a :

$$\langle X, 1 \rangle = \sum_{i=1}^N x_i = NE(X)$$

et plus généralement, pour toute constante λ :

$$\langle X, \lambda 1 \rangle = \lambda \sum_{i=1}^N x_i = N\lambda E(X)$$

Deux variables sont orthogonales dans l'espace euclidien si $\langle X, Y \rangle = 0$.

Au sens probabiliste, la covariance de deux variables est définie par :

$$\text{cov}(X, Y) = E((X - EX)(Y - EY)) = E(XY) - (EX)(EY)$$

$$\text{avec } E(XY) = \frac{1}{N} \sum_{i=1}^N x_i y_i.$$

On a donc :

$$\text{cov}(X, Y) = E((X - EX)(Y - EY)) = \frac{1}{N} \langle X - EX.1, Y - EY.1 \rangle$$

Le coefficient de corrélation s'écrit :

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{\frac{1}{N} \langle X - EX.1, Y - EY.1 \rangle}{\frac{1}{\sqrt{N}} \|X - m_X\|_2 \frac{1}{\sqrt{N}} \|Y - m_Y\|_2} = \frac{\langle X - EX.1, Y - EY.1 \rangle}{\|X - m_X\|_2 \|Y - m_Y\|_2}$$

Si les variables sont centrées réduites :

$$\text{corr}(X, Y) = \langle X, Y \rangle$$