

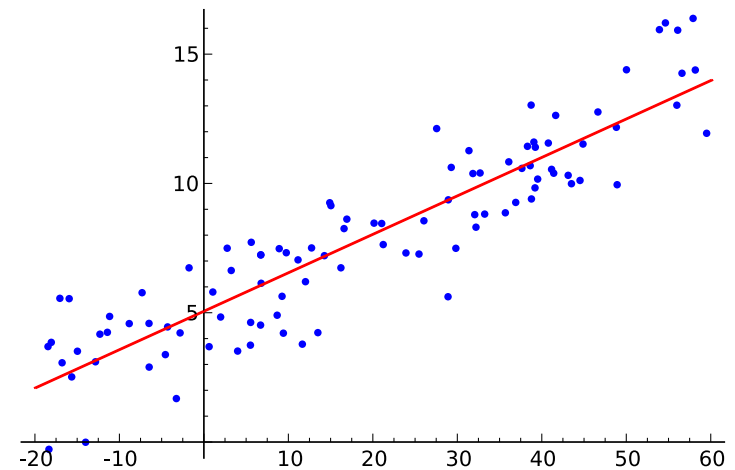
INFO601_CMI Algorithmique numérique

Introduction au domaine via la méthode des moindres carrés

[Jacques-Olivier Lachaud]

Modèle prédictif à partir de données

- données observées $[(x_1, y_1), \dots, (x_n, y_n)]$
- peuvent conduire à une *prédiction* ?
- *prédiction* : $f(x)$ avec $\forall i, f(x_i) \approx y_i$
- famille paramétrée de modèles f_α
- modèle prédictif le plus “*pertinent*” ?
⇒ Trouver α^* tel que f_{α^*} prédit au mieux les y à partir des x

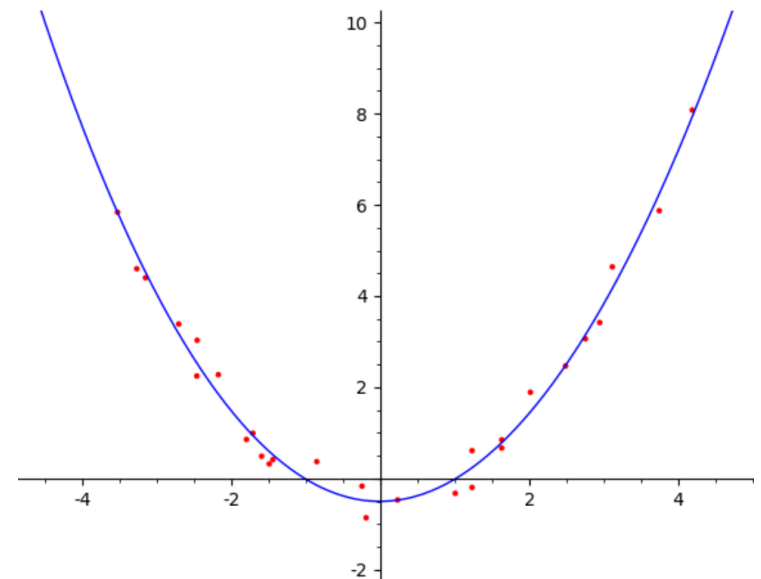


modélisation

- identifier une famille de fonction f_α
- caractériser les résultats pertinents

optimisation

- écrire le problème sans forme soluble
- chercher le meilleur paramètre α^*



Méthode des moindres carrés ordinaires

[Gauss, Legendre, circa 1800]

modélisation

- forme le vecteur des erreurs ou **résidu**

$$\mathbf{r} = \begin{pmatrix} y_1 - f_{\alpha}(x_1) \\ \vdots \\ y_n - f_{\alpha}(x_n) \end{pmatrix}$$

- erreur quadratique (**norme** euclidienne au carré)

$$E(\alpha) = \|\mathbf{r}\|^2 = \sum_{i=1}^n (y_i - f_{\alpha}(x_i))^2$$

- $E(\alpha)$ est une **fonction coût**
- modèle pertinent = petit coût

optimisation

- $\alpha^* := \arg \min_{\alpha \in D} E(\alpha)$
- problème **difficile** en général
- $\min E \Rightarrow$ dérivées nulles

\Rightarrow gradient $\nabla E = 0$

- trouver les zéros du **gradient**
- algorithmes de **descente en gradient**
- résolution de **systèmes linéaires**

Régression linéaire simple

- ajustement droite $f_{\alpha}(x) = \alpha_1 x + \alpha_2$
- erreur quadratique $E(\alpha_1, \alpha_2) = \sum_{i=1}^n r_i^2 = \sum_{i=1}^n (y_i - \alpha_1 x_i + \alpha_2)^2$

$$E(\alpha_1, \alpha_2) = (\alpha_1 \ \alpha_2) \underbrace{\begin{pmatrix} \sum x_i^2 & \sum x_i \\ \sum x_i & \sum 1 \end{pmatrix}}_M \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} - 2(\alpha_1 \ \alpha_2) \underbrace{\begin{pmatrix} \sum x_i y_i \\ \sum y_i \end{pmatrix}}_v + \underbrace{\sum y_i^2}_c$$
$$= \alpha^T M \alpha - 2\alpha^T v + c$$

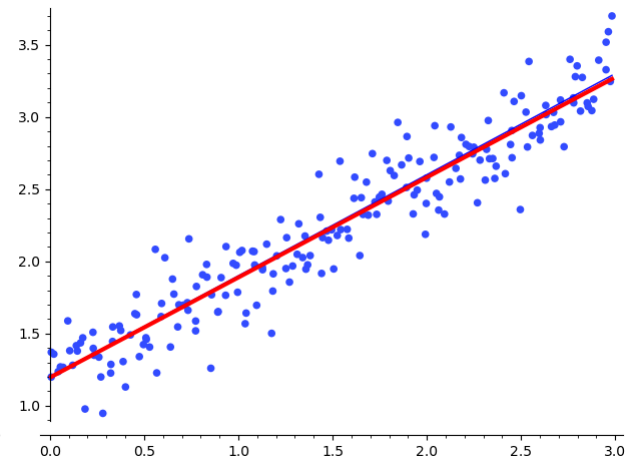
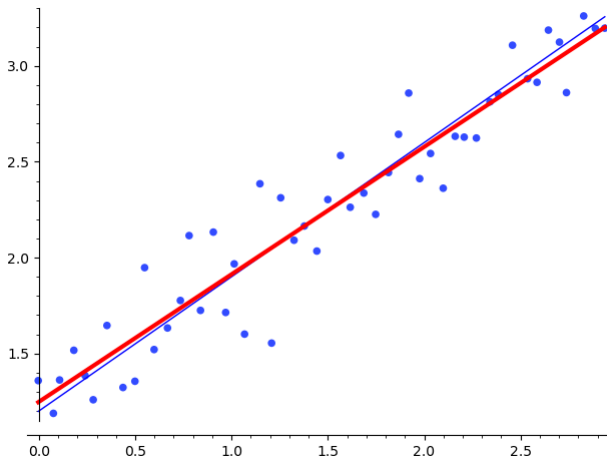
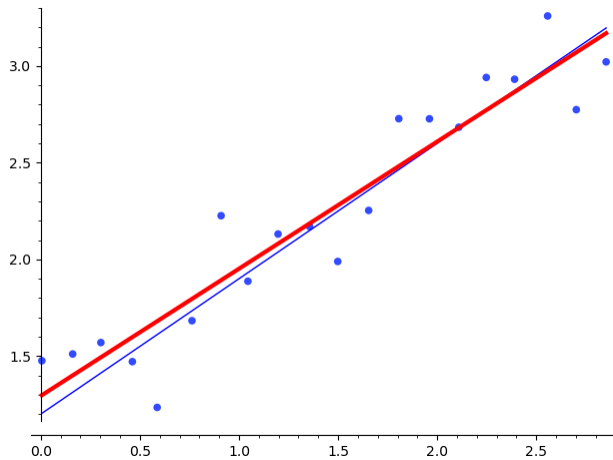
- *forme quadratique définie positive* $\Rightarrow E$ a un seul minimum
- **optimisation** mieux ici de chercher le zéro du gradient de E

$$\nabla E = \begin{pmatrix} \frac{\partial E}{\partial \alpha_1} \\ \frac{\partial E}{\partial \alpha_2} \end{pmatrix} = \begin{pmatrix} 2m_{11}\alpha_1 + (m_{12} + m_{21})\alpha_2 - 2v_1 \\ (m_{12} + m_{21})\alpha_1 + 2m_{22}\alpha_2 - 2v_2 \end{pmatrix} = \underbrace{(M + M^T)}_{2M} \alpha - 2v$$

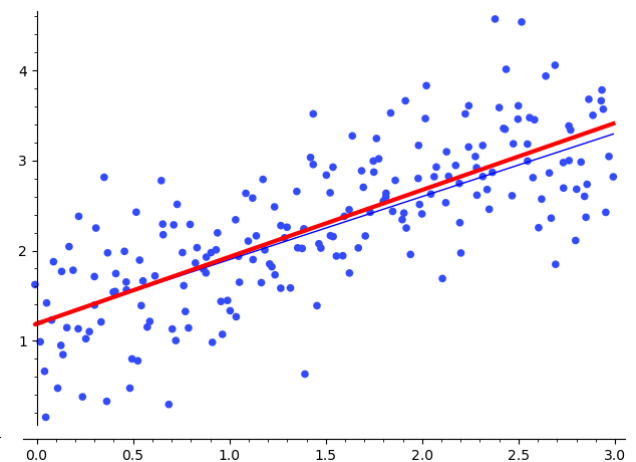
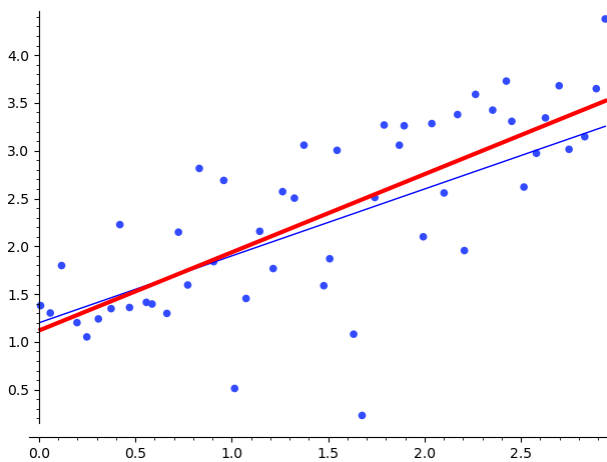
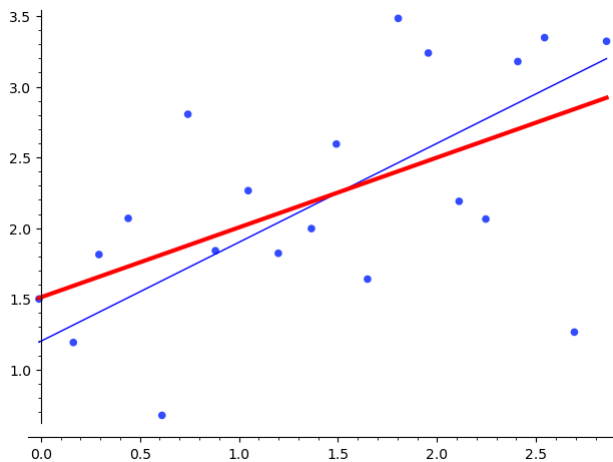
$$\alpha^* \text{ solution de } \nabla E = 0 \Leftrightarrow M\alpha = v, \text{ i.e. } \alpha^* = M^{-1}v$$

Régression linéaire simple

perturbation $\sigma = 0.2$



perturbation $\sigma = 0.6$



Des garanties théoriques

Si $y_i = \alpha_1 x_i + \alpha_2 + \varepsilon_i$, avec ε_i perturbation de la i -ème donnée

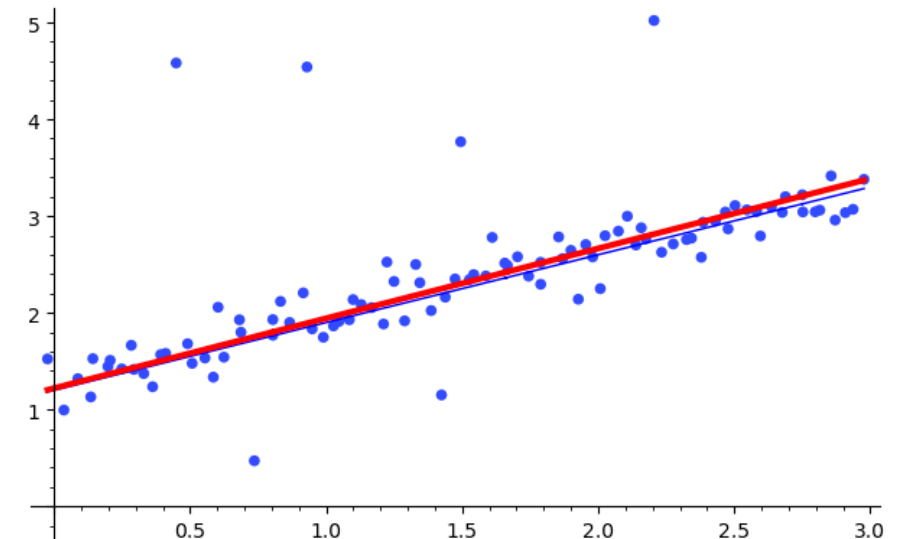
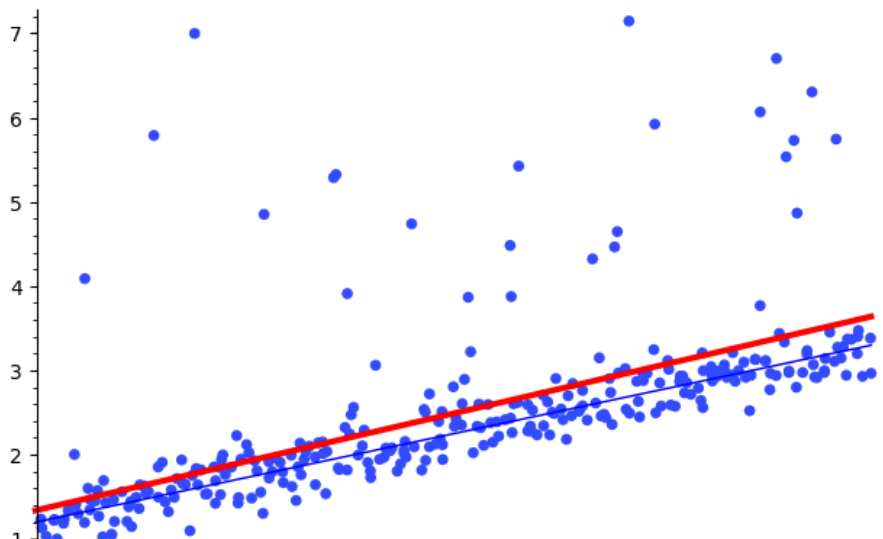
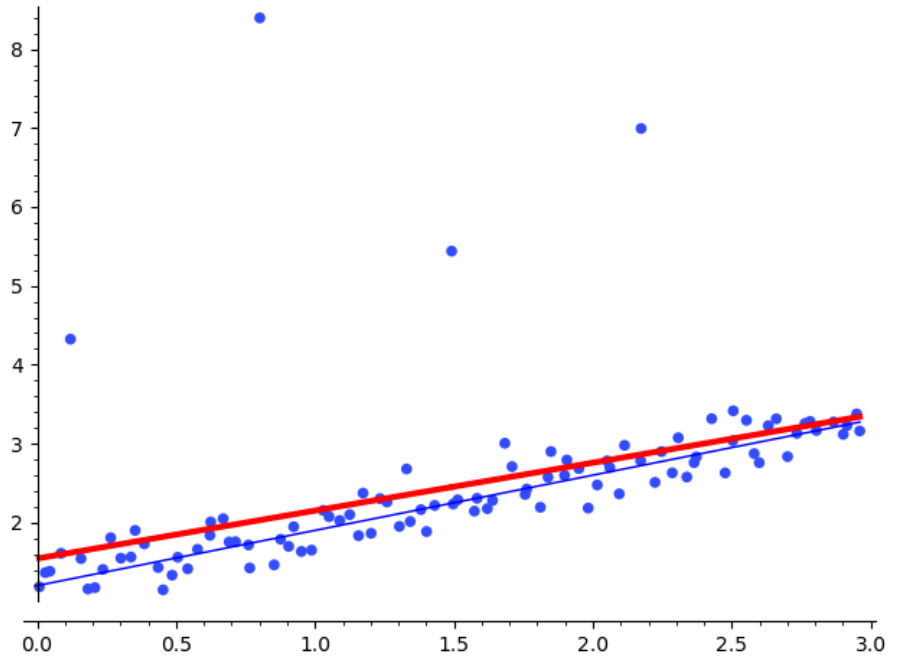
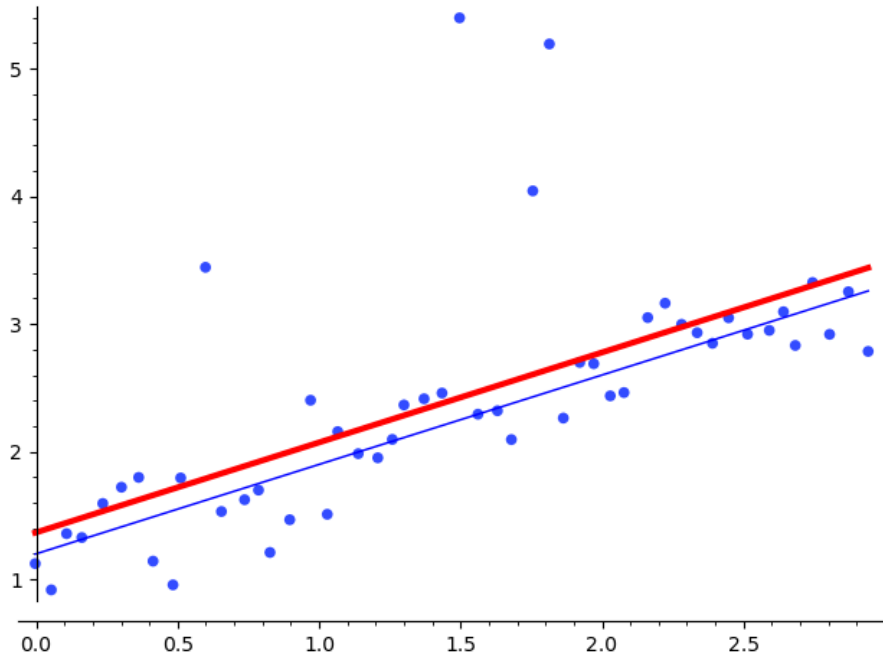
$$\underbrace{Y = X\alpha + \varepsilon}_{\text{écriture matricielle}} \Leftrightarrow \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_1 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

Thm [Gauss-Markov]: parmi tous les estimateurs linéaires non biaisés, l'estimateur par moindres carrés présente une variance minimale, i.e. $\|\alpha^* - \alpha\|^2$ est minimal. C'est le **BLUE** !

hypothèses

- $\mathbb{E}(\varepsilon_i) = 0$ "les erreurs sont sans biais"
- $\text{Var}(\varepsilon_i) = \sigma^2 < +\infty$ "les erreurs ont même variance, finie"
- $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0, i \neq j$ "les erreurs sont non corrélées"
- *estimateur linéaire* de α_j : $\hat{\alpha}_j = c_1 y_1 + \dots + c_n y_n$, avec c_i qui peuvent dépendre des x_i

Problème des données aberrantes (outliers)



Elimination des données aberrantes

utilisation des résidus

- enlever 10% des données ayant le plus gros résidu

heuristiques

- random forest, isolation forest, distance au barycentre des k-NN

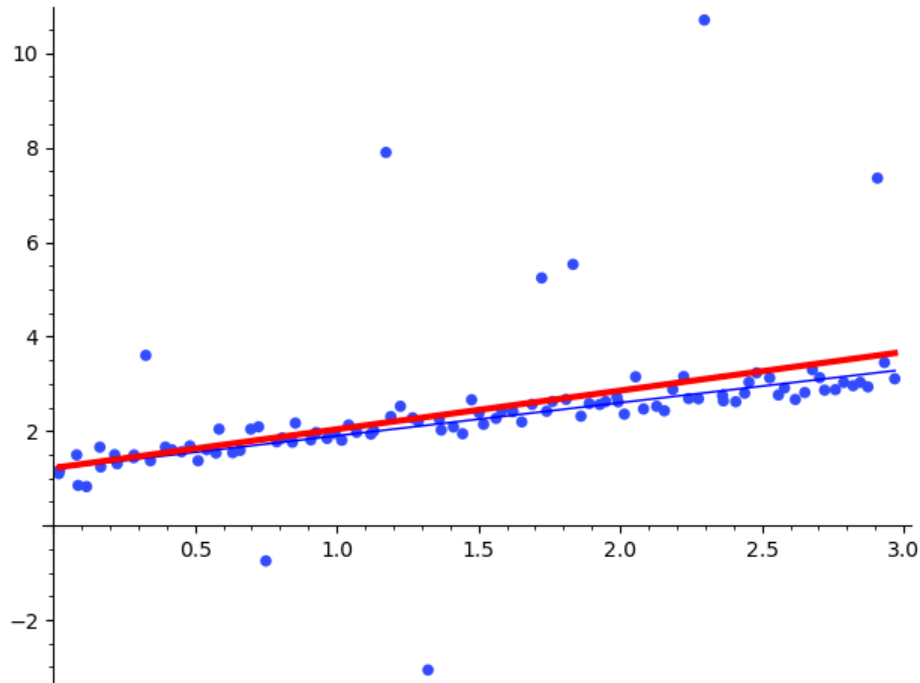
Autres méthodes statistiques

- Theil–Sen estimator: pente α_1 en prenant le médian des pentes de tous les couples de points.

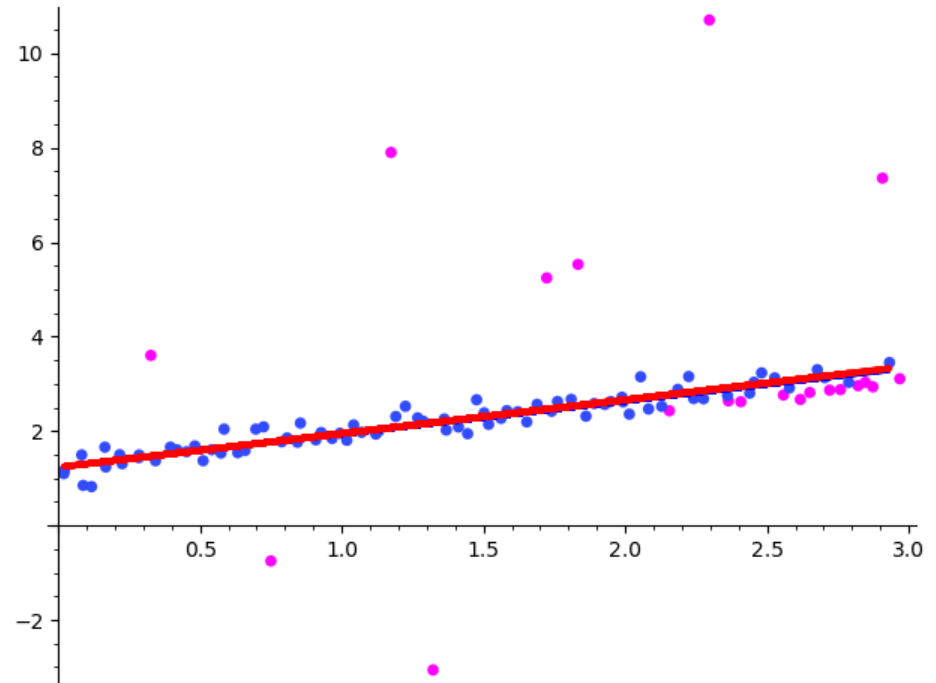
Problème de paramétrages, de choix de modèles, d'hypothèses sur les données, de temps de calcul, ...

Exercice : Complexité de la régression linéaire (n données) ? De l'estimateur Theil-Sen ?

Elimination des données aberrantes



sans filtrage

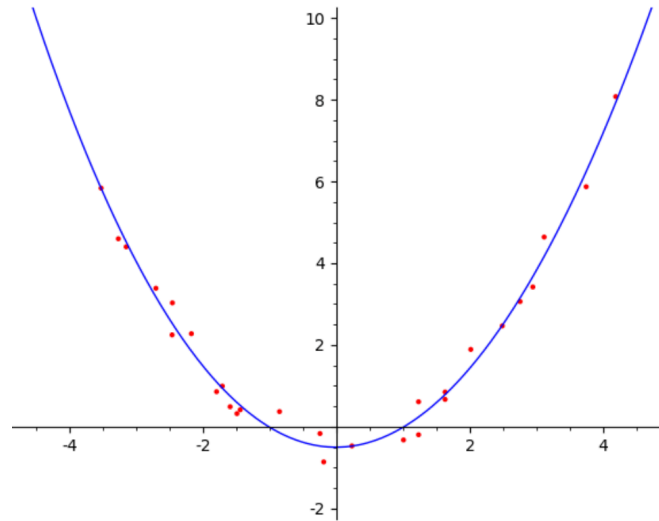


enlève 20% des données + recalcul

- on calcule α^* pour toutes les données
- on trie les données selon leur résidus r^2
- on garde les premiers 80% des données
- on recalcule α^* pour ces données

Et si on cherche des modèles plus complexes que
la droite ?!

Régression linéaire multiple



- on suppose que les données $(x_i, y_i)_{i=1, \dots, n}$ suivent une loi parabolique
- pour chaque donnée, $y_i = \alpha_1 x_i^2 + \alpha_2 x_i + \alpha_3 + \varepsilon_i$, avec ε_i une erreur

$$\underbrace{\mathbf{Y} = \mathbf{X}\boldsymbol{\alpha} + \boldsymbol{\varepsilon}}_{\text{écriture matricielle}} \Leftrightarrow \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_1^2 & x_1 & 1 \\ \vdots & \vdots & \\ x_n^2 & x_n & 1 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

La relation est encore linéaire !

Régression linéaire multiple

- vecteur résidu : $\mathbf{r} = \mathbf{Y} - \mathbf{X}\boldsymbol{\alpha}$, erreur quadratique (ou fonction coût):

$$\begin{aligned} E(\boldsymbol{\alpha}) &= \|\mathbf{r}\|^2 = \mathbf{r}^\top \mathbf{r} = (\mathbf{Y} - \mathbf{X}\boldsymbol{\alpha})^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\alpha}) \\ &= \boldsymbol{\alpha}^\top \underbrace{\mathbf{X}^\top \mathbf{X}}_M \boldsymbol{\alpha} - 2\boldsymbol{\alpha}^\top \underbrace{\mathbf{X}^\top \mathbf{Y}}_v + \underbrace{\mathbf{Y}^\top \mathbf{Y}}_c \end{aligned}$$

Rappel: régression linéaire simple

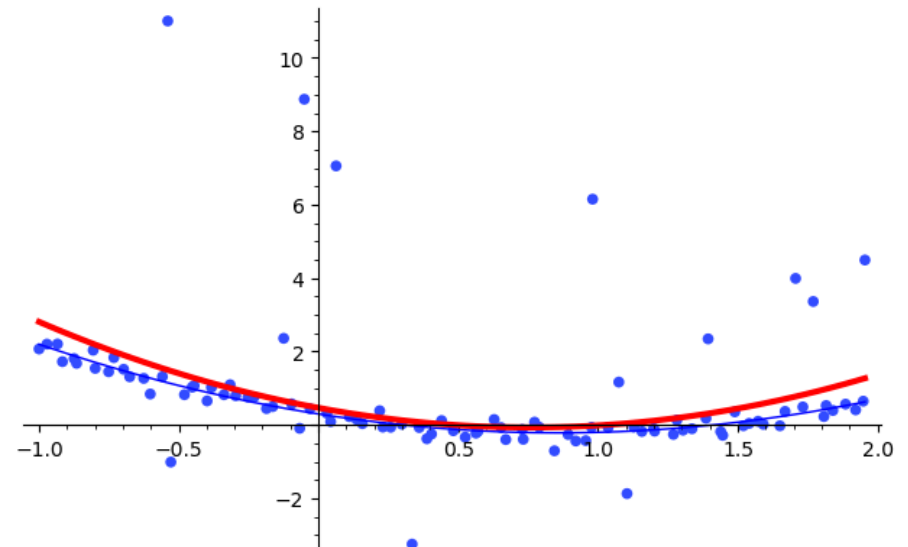
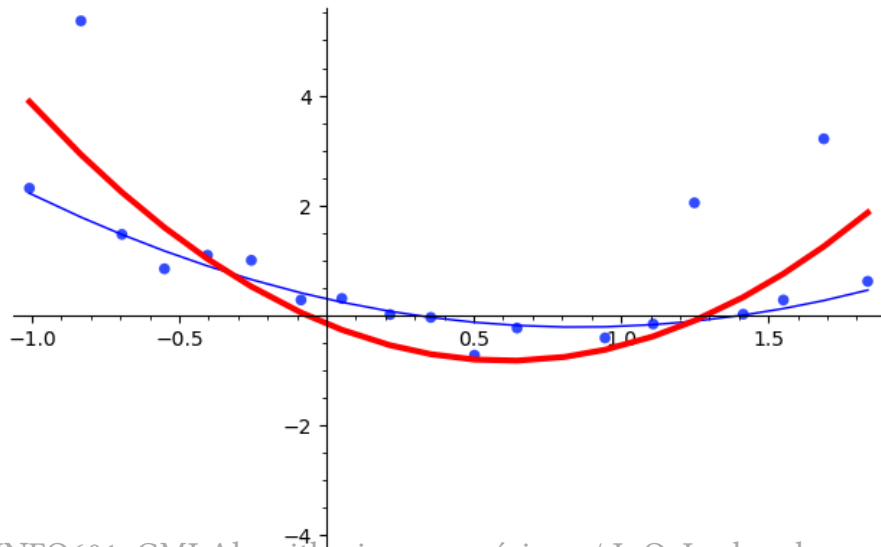
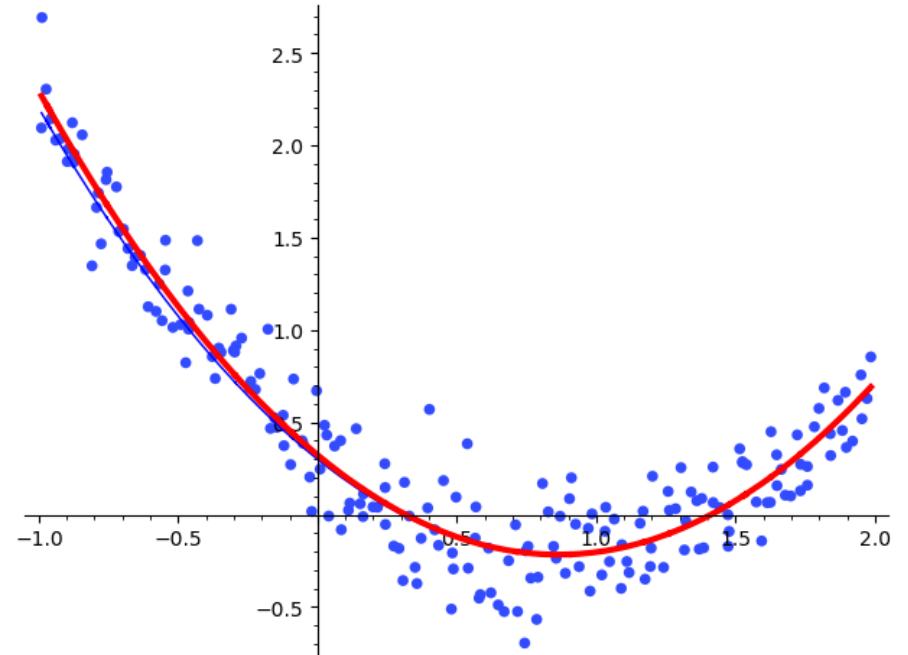
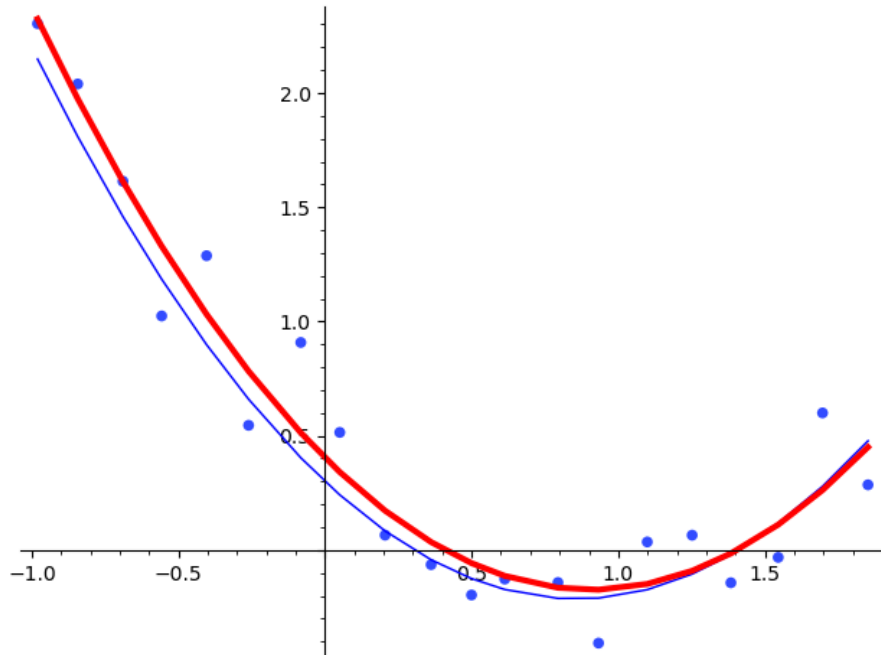
$$E(\alpha_1, \alpha_2) = (\alpha_1 \ \alpha_2) \underbrace{\begin{pmatrix} \sum x_i^2 & \sum x_i \\ \sum x_i & \sum 1 \end{pmatrix}}_M \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} - 2(\alpha_1 \ \alpha_2) \underbrace{\begin{pmatrix} \sum x_i y_i \\ \sum y_i \end{pmatrix}}_v + \underbrace{\sum y_i^2}_c$$

- **optimisation** en cherchant le zéro du gradient de E

$$\nabla E = \frac{\partial E}{\partial \boldsymbol{\alpha}} = 2\mathbf{X}^\top \mathbf{X}\boldsymbol{\alpha} - 2\mathbf{X}^\top \mathbf{Y} = 2M\boldsymbol{\alpha} - 2v$$

$$\boldsymbol{\alpha}^* \text{ solution de } \nabla E = 0 \Leftrightarrow \boldsymbol{\alpha}^* = M^{-1}v = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

Régression linéaire multiple



Méthode des moindres carrés

Régression linéaire multiple

- ajuster un polynôme de degré quelconque reste linéaire
- ajuster un plan à des données (x_i, y_i, z_i) reste linéaire
- ajuster une surface polynomiale à des données (x_i, y_i, z_i) reste linéaire

Moindres carrés généralisés

- donner un poids $w_i = \frac{1}{\sigma_i}$ différent à chaque donnée (typique d'une incertitude sur une mesure)
- donner une matrice de variance/covariance Σ entre les données

⇒ une approche similaire à la précédente fonctionne

Moindres carrés non linéaires

- ajuster une fonction $y = \cos(\alpha_1 x + \alpha_2)$ n'est pas linéaire

⇒ mais on peut ajuster $\arccos(y) = \alpha_1 x + \alpha_2$!

- ajuster une fonction $y = \exp(-\alpha_1(x - \alpha_2)^2)$ n'est pas linéaire

Quelques exercices

Exercice 1: On cherche à ajuster un polynôme de degré 3 à des données $(x_i, y_i)_{i=1, \dots, n}$. Explicitez la relation matricielle $r = Y - X\alpha$.

Exercice 2: On vous donne les données suivantes : $(-1, 2), (1, 3), (3, 5), (4, 6), (6, 8)$ Trouvez la droite de meilleur ajustement.

Exercice 3: On dispose d'un nuage de points $(x_i, y_i, z_i)_{i=1, \dots, n}$. Proposez une approche pour trouver un plan approchant ces points.

Exercice 4: Est-ce que l'on obtient la même droite d'ajustement si on ajuste plutôt les x par rapport aux y ? On suppose donc que les données suivent

$$x_i = \beta_0 y_i + \beta_1 + \varepsilon'_i$$