

## TP3 – Estimation robuste

### Influence des données aberrantes sur l'estimation aux moindres carrés

Il a été vu lors du TP2 comment estimer les paramètres de la droite de régression  $D_{YX}$  de  $Y$  en  $X$ , et ceux de la droite de régression  $D_{\perp}$  en distance orthogonale, à partir d'un nuage de points. Si parmi ces points, certains constituent des *données aberrantes*, cela risque fort d'altérer l'estimation. Lancez le script `donnees_aberrantes.m`, qui permet de visualiser l'influence de la présence de données aberrantes sur l'erreur angulaire moyenne commise sur la direction de chaque droite de régression. Bien sûr, ces erreurs augmentent avec la proportion de données aberrantes, mais vous pouvez aussi faire les observations suivantes :

- En l'absence de données aberrantes, la droite de régression  $D_{YX}$  est plus précise, en moyenne, que la droite de régression  $D_{\perp}$ . Au-delà de 5% de données aberrantes, la situation s'inverse.
- Si l'addition d'un bruit gaussien et le tirage de valeurs aberrantes (selon une loi uniforme) s'appliquent non seulement aux ordonnées, mais également aux abscisses des points du nuage, alors l'estimation de la droite  $D_{YX}$  se dégrade nettement. Au contraire, l'écart angulaire moyen de  $D_{\perp}$  diminue très légèrement.

Le but de ce TP est de vous montrer deux façons de limiter l'influence des données aberrantes sur l'estimation des droites de régression, en remettant en question l'estimation au sens des moindres carrés vue dans le TP2.

### Première parade aux données aberrantes : l'algorithme RANSAC

RANSAC (abréviation de *RANdom SAMple Consensus*) est un algorithme itératif d'*estimation robuste*, publié par Fischler et Bolles en 1981, qui consiste à effectuer une *partition des données* entre données aberrantes (*outliers*) et données conformes au modèle (*inliers*). Une caractéristique de cet algorithme est qu'il est non déterministe : le résultat n'est garanti qu'avec une certaine probabilité, qui croît avec le nombre d'itérations.

L'exemple qui illustre le mieux l'algorithme RANSAC est l'estimation robuste d'une droite de régression. Le principe de RANSAC consiste à tirer aléatoirement un sous-ensemble de données de cardinal égal au nombre minimal de données permettant d'estimer le modèle (nombre égal à 2 dans le cas d'une droite de régression). Ces données sont considérées comme des données conformes au modèle (cela reste à vérifier), puis la séquence suivante est répétée en boucle :

1. Les paramètres du modèle sont estimés à partir de ce sous-ensemble de données conformes.
2. Toutes les autres données sont testées relativement au modèle estimé, afin de détecter les données conformes (un point est jugé conforme si sa distance à la droite de régression est inférieure à un seuil).
3. Le modèle estimé est accepté si la proportion de données conformes est supérieur à un seuil.
4. Si le modèle est accepté, il est réestimé à partir de l'ensemble des données conformes.

Le modèle retenu est celui qui donne la plus faible moyenne des carrés des résidus des données conformes.

### Exercice 1 : estimation robuste par l'algorithme RANSAC

Complétez le script `exercice_1.m`, qui applique l'algorithme RANSAC à l'estimation de la droite de régression  $D_{\perp}$ , lorsque l'addition d'un bruit gaussien et le tirage de valeurs aberrantes s'appliquent à la fois aux abscisses et aux ordonnées des points du nuage. Si le bénéfice de cette nouvelle méthode d'estimation semble incontestable en termes de robustesse aux données aberrantes, elle pêche manifestement par la nécessité de régler correctement ses paramètres `k_max`, `seuil_distance` et `seuil_proportion`, et surtout par sa lenteur !

## Deuxième parade aux données aberrantes : autres critères à optimiser

Dans le TP2, l'estimation des paramètres  $(a, b)$  de la droite de régression  $D_{YX}$  et celle des paramètres  $(\theta, \rho)$  de la droite de régression  $D_{\perp}$  ont été effectuées en résolvant des problèmes d'optimisation du type suivant :

$$\min_{\theta=[\theta_1, \dots, \theta_p]} \left\{ \sum_{i=1}^n r_{\theta}(P_i)^2 \right\} \quad (1)$$

Or, les données aberrantes dégradent la précision de ces estimations à cause de résidus  $r_{\theta}(P_i)$  élevés, puisque ces points seront probablement éloignés de la droite de régression. Cela est encore amplifié par le fait que les résidus sont élevés au carré dans (1). Une deuxième façon de limiter l'influence des données aberrantes consiste donc à conserver toutes les données, mais à utiliser des *moindres carrés pondérés*, c'est-à-dire un nouveau critère  $\mathcal{W}(\theta) = \sum_{i=1}^n w_i r_{\theta}(P_i)^2$ , en faisant en sorte que le poids  $w_i \geq 0$  du point  $P_i$  soit d'autant plus faible que le résidu  $r_{\theta}(P_i)$  est plus élevé. Malheureusement, il semble que cette idée soit impossible à mettre en œuvre, car on ne sait pas quels points constituent des données aberrantes (le principe de cette approche est de ne pas effectuer de partition des données, contrairement à l'algorithme RANSAC). Néanmoins, l'optimalité de ce critère s'écrit :

$$\nabla \mathcal{W}(\theta) = 0 \iff \sum_{i=1}^n w_i r_{\theta}(P_i) \frac{\partial r_{\theta}}{\partial \theta_j}(P_i) = 0, \quad j \in [1, p] \quad (2)$$

D'autre part, toute fonction  $\phi$  d'une variable réelle à valeurs dans  $\mathbb{R}^+$ , dérivable, paire, croissante sur  $\mathbb{R}^+$ , permet de définir un critère  $\mathcal{H}(\theta) = \sum_{i=1}^n \phi(r_{\theta}(P_i))$  tout aussi valide que le critère utilisé dans (1), c'est-à-dire tout aussi valide que la somme des carrés des résidus. L'optimalité du critère  $\mathcal{H}(\theta)$  s'écrit :

$$\nabla \mathcal{H}(\theta) = 0 \iff \sum_{i=1}^n \phi'(r_{\theta}(P_i)) \frac{\partial r_{\theta}}{\partial \theta_j}(P_i) = 0, \quad j \in [1, p] \quad (3)$$

Par identification des équations (2) et (3), on trouve les poids suivants :

$$w_i = \frac{\phi'(r_{\theta}(P_i))}{r_{\theta}(P_i)}, \quad i \in [1, n] \quad (4)$$

Il existe de nombreuses fonctions  $\phi$  nulles en 0, telles que ces poids décroissent lorsque  $r_{\theta}(P_i)$  croît, comme par exemple  $\phi_1(x) = \sqrt{x^2 + \alpha^2} - \sqrt{\alpha^2}$  ou  $\phi_2(x) = \ln(x^2 + \beta^2) - \ln(\beta^2)$ , où le rôle des paramètres  $\alpha$  et  $\beta$  est de rendre ces fonctions dérivables en  $x = 0$ . Dans la pratique, il n'est pas question de résoudre les équations (3), qui n'admettent pas de solution analytique pour des fonctions telles que  $\phi_1$  ou  $\phi_2$ . En revanche, il est facile de minimiser le critère  $\mathcal{H}(\theta)$  par la méthode du maximum de vraisemblance. Comme la dérivabilité n'est alors plus requise, il suffit de résoudre un des problèmes suivants :

$$\begin{cases} \min_{\theta=[\theta_1, \dots, \theta_p]} \sum_{i=1}^n \left\{ \sqrt{r_{\theta}(P_i)^2 + \alpha^2} - |\alpha| \right\} \\ \min_{\theta=[\theta_1, \dots, \theta_p]} \sum_{i=1}^n \left\{ \ln(r_{\theta}(P_i)^2 + \beta^2) - 2 \ln |\beta| \right\} \end{cases} \quad (5)$$

Notez que pour  $\alpha = 0$ , le critère  $\sum_{i=1}^n \phi_1(r_{\theta}(P_i)) = \sum_{i=1}^n |r_{\theta}(P_i)|$  est très souvent rencontré : il s'agit de la « norme  $L_1$  du vecteur des résidus ».

## Exercice 2 : estimation robuste aux données aberrantes de la droite $D_{\perp}$

Faites une copie du script `donnees_aberrantes.m`, de nom `exercice_2.m`, que vous modifierez de manière à ne conserver que l'estimation de la droite de régression  $D_{\perp}$  lorsque l'addition d'un bruit gaussien et le tirage de valeurs aberrantes s'appliquent à la fois aux abscisses et aux ordonnées.

Comparez ensuite cette estimation au sens des moindres carrés avec la résolution des problèmes (5), en fixant  $\alpha = \beta = 1$ . Vous supposerez encore que toute droite de régression contient le centre de gravité, bien que cela ne soit plus le cas. Enfin, observez l'effet d'une modification des paramètres  $\alpha$  et  $\beta$ .