Data Science

M2 Actuariat

Session 2 : Modèles Linéaires Généralisés et Pénalisés

François HU

Responsable du pôle Intelligence Artificielle, Milliman R&D Enseignant ISFA 2025





Sommaire du cours Data Science

- 1. Apprentissage statistique et lien avec l'Actuariat
- 2. Modèles linéaires généralisés et pénalisés
- 3. Arbre de décision et méthodes ensemblistes
- 4. Interprétabilité des modèles d'apprentissage
- 5. IA de confiance et biais algorithmiques
- 6. Apprentissage non supervisé
- 7. Introduction aux données non structurées



Session modèles linéaires généralisés et pénalisés

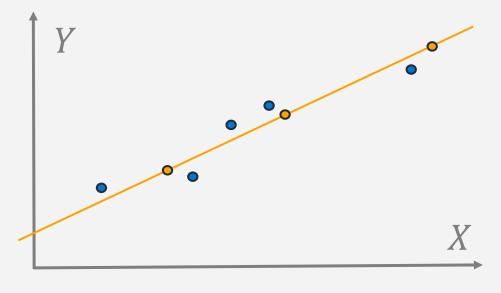
- 1. Modèles Linéaires : Rappel et Limites
- 2. Composants des modèles GLM : fonctions de lien et distributions de probabilité
- 3. Modèle Logistique et Poisson
- 4. Estimation des paramètres : La méthode de vraisemblance
- 5. Interprétation des coefficients
- Evaluation des modèles GLM
- 7. Sélection des variables explicatives
- 8. GLM pénalisé : Régression Lasso et Ridge
- 9. Extension des modèles GLM aux modèles GAM



1. Modèles Linéaires

Rappel

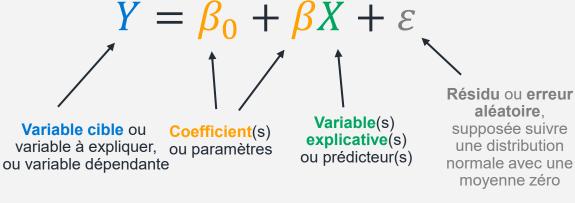
Limites





Rappel des modèles linéaires simples (1/5)

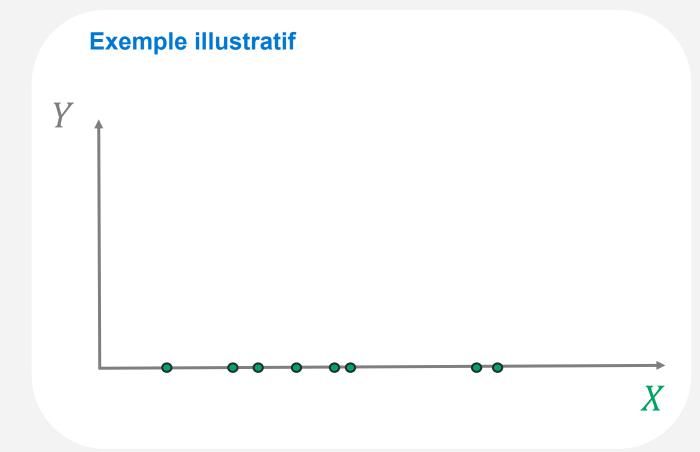
Définitions et notations



Les **modèles linéaires simples** permettent de décrire la relation entre deux variables quantitatives (**variable cible** *Y* et **variable explicative** *X*) par une équation linéaire.

Pour un **portefeuille automobile**, la variable cible peut être le coût du sinistre (en supposant que la personne est exposée) et la variable explicative l'âge du conducteur.

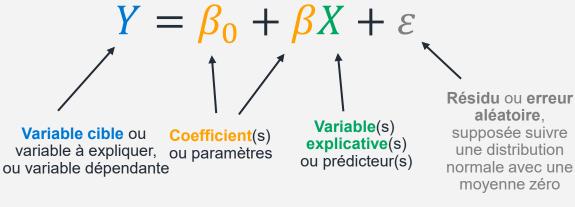
Le **principal objectif** du modèle linéaire est de **prédire** les valeurs Y à partir de X en estimant les **coefficients** β **et** β_0 . Les coefficients sont généralement estimés par la **méthode des moindres carrées**.





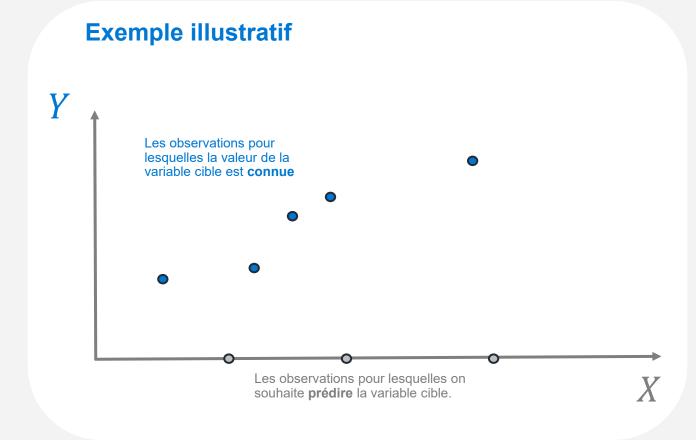
Rappel des modèles linéaires simples (2/5)

Evaluer la généralisation



Séparation train-test : La séparation entre les ensembles d'entraînement (train) et de test est importante pour évaluer la performance réelle d'un modèle linéaire

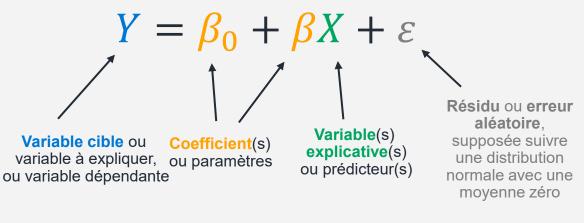
Évaluer la généralisation : Le modèle linéaire est ajusté sur l'ensemble d'entraînement (train), les données qu'il utilise pour apprendre. Cependant, cela ne garantit pas qu'il fonctionnera bien sur de nouvelles données, non vues auparavant. L'ensemble de test, composé de données que le modèle n'a jamais vues, permet de vérifier la capacité du modèle à généraliser à des données inconnues.





Rappel des modèles linéaires simples (3/5)

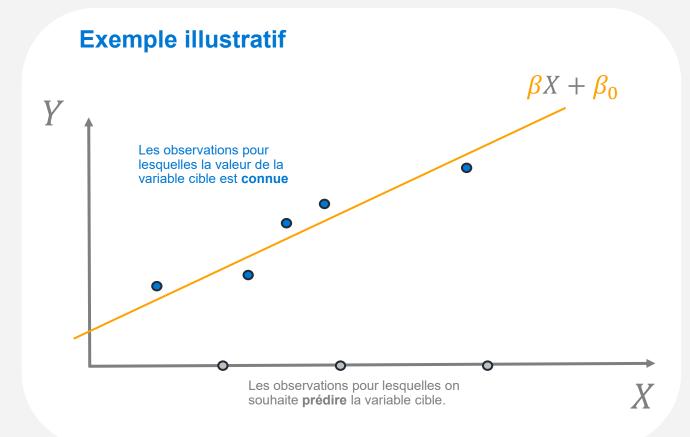
Différence entre entraînement et inférence



L'entraînement consiste à ajuster les **paramètres** d'un modèle à partir de **données connues** pour qu'il apprenne des patterns.

L'inférence, en revanche, est le processus où le modèle, une fois entraîné, fait des prédictions sur des nouvelles **données non vues**.

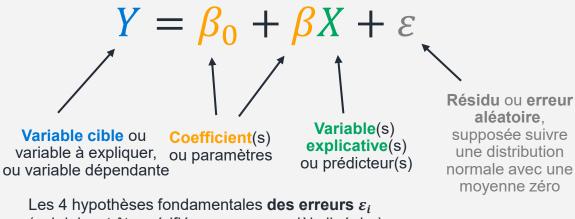
En résumé, l'entraînement est l'apprentissage, tandis que l'inférence est l'utilisation du modèle pour prédire.





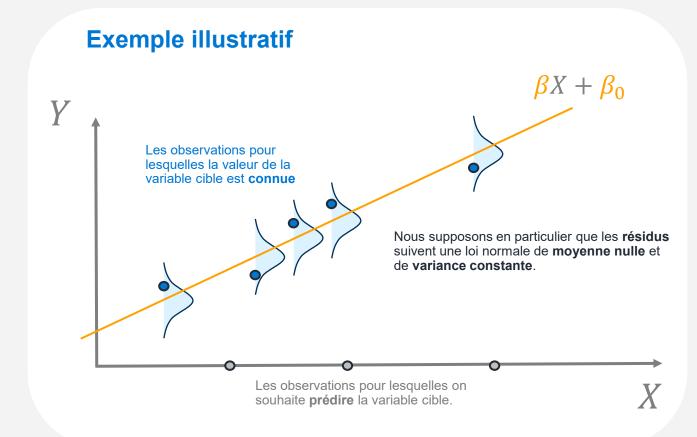
Rappel des modèles linéaires simples (4/5)

Hypothèses d'un modèle linéaire



Les 4 hypothèses fondamentales des erreurs ε_i (qui doivent être vérifiées pour un modèle linéaire) :

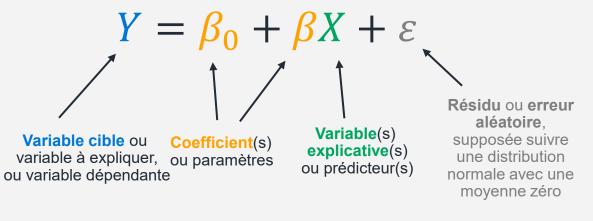
- Moyenne nulle (erreurs centrées) : le modèle ne doit pas présenter de biais systématique.
- Variance constante (homoscédasticité) : La dispersion des erreurs doit être constante, pour éviter que certaines parties du modèle soient plus imprécises que d'autres.
- Indépendance des erreurs : l'erreur commise pour une observation ne doit pas influencer les autres.
- Distribution normale (loi gaussienne): Les erreurs suivent une distribution normale, ce qui permet d'appliquer les tests statistiques et d'estimer la précision des prédictions.





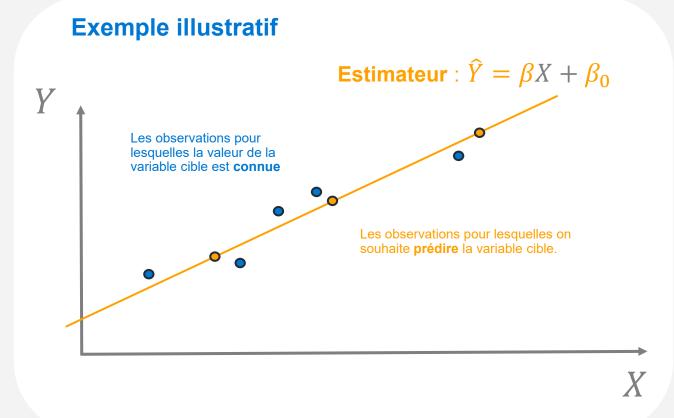
Rappel des modèles linéaires simples (5/5)

Performance du modèle



Plusieurs approches pour vérifier si un modèle linéaire est bien estimé :

- **Visuellement**: Utiliser des graphiques comme le QQ-plot pour vérifier si les erreurs suivent une distribution normale.
- **Empiriquement :** Utiliser le coefficient de détermination \mathbb{R}^2 pour mesurer la proportion de la variation expliquée par le modèle par rapport à la variation totale.
- **Statistiquement**: Effectuer une analyse de la variance (ANOVA) pour évaluer la significativité du modèle avec la statistique F et sa p-valeur.





Plusieurs variables explicatives : des modèles linéaires dites « multiples »

Cas avec plusieurs variables explicatives (souvent plus de 10 variables).

Coût sinistre	Catégorie	Ancienneté	Age
0	3	0	41
122.3	1	8	20
977.8	2	7	18
0	1	5	47
311.1	1	2	29



Modèles linéaires multiples (2/5)

Un peu de formalisme

Modèles linéaires multiples $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_d X_d + \varepsilon$ avec $\varepsilon_i iid \sim N(0, \sigma^2)$

Y	X_1	X_2	X_3	ε	
0	3	0	41	$arepsilon_{arepsilon_1}$	
122.3	1	8	20	ε_2	
977.8	2	7	18	\mathcal{E}_3	
0	1	5	47	\mathcal{E}_4	
311.1	1	2	29	\mathcal{E}_{5}	



Modèles linéaires multiples (3/5)

Ecriture matricielle

Modèles linéaires multiples
$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_d X_d + \varepsilon$$
 avec $\varepsilon_i iid \sim N(0, \sigma^2)$

Ecriture matricielle

$\begin{bmatrix} 0 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \end{bmatrix} \begin{bmatrix} 0 \\ 41 \end{bmatrix}$	β_0	$\begin{bmatrix} \epsilon_1 \end{bmatrix}$	
122.3 1 1 8 20 977.8 = 1 2 7 18 × 0 1 1 5 47 311.1 1 1 2 29	$egin{array}{c c} eta_1 \\ eta_2 \\ eta_3 \end{array} +$	$egin{array}{c} arepsilon_2 & & & & & & & & & & & & & & & & & & &$	



Modèles linéaires multiples (4/5)

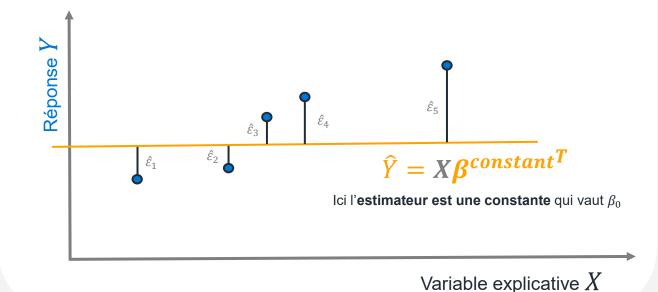
Estimateur constant

Modèle linéaire multiple $Y = X\beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2)$

Notre objectif est d'estimer les coefficients β . Plus précisément nous souhaitons trouver β qui minimise la méthode des moindres carrés, une approche statistique utilisée pour ajuster un modèle à des données en minimisant la somme des carrés des écarts entre les valeurs observées et les valeurs prédites par le modèle

$$R(\boldsymbol{\beta}) \coloneqq \sum_{i} (Y_i - X_i \beta_i^T)^2 = \sum_{i} \hat{\varepsilon}_i^2$$

Exemple d'un mauvais estimateur avec $R(\beta^{best})$ grand





Modèles linéaires multiples (5/5)

Estimateur optimal

Modèle linéaire multiple $Y = X\beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2)$

Notre objectif est d'estimer les coefficients β . Plus précisément nous souhaitons trouver β qui minimise la méthode des moindres carrés, une approche statistique utilisée pour ajuster un modèle à des données en minimisant la somme des carrés des écarts entre les valeurs observées et les valeurs prédites par le modèle

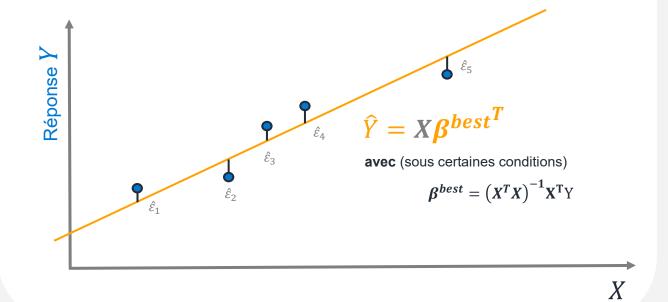
$$R(\boldsymbol{\beta}) \coloneqq \sum_{i} (Y_i - X_i \beta_i^T)^2 = \sum_{i} \hat{\varepsilon}_i^2$$

Expression explicite du paramètre optimal β^{best} : si toutes les variables explicatives dans X sont **linéairement** indépendantes entre elles alors

$$\boldsymbol{\beta}^{best} = \left(\boldsymbol{X}^T \boldsymbol{X} \right)^{-1} \boldsymbol{X}^T \boldsymbol{Y}$$

avec $\mathbb{E}[\boldsymbol{\beta}^{best}] = \boldsymbol{\beta}$ (on dit que l'estimateur est non biaisé) $\mathbb{V}[\boldsymbol{\beta}^{best}] = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$ la matrice variance-covariance

Exemple d'un bon estimateur avec $R(\beta^{best})$ petit

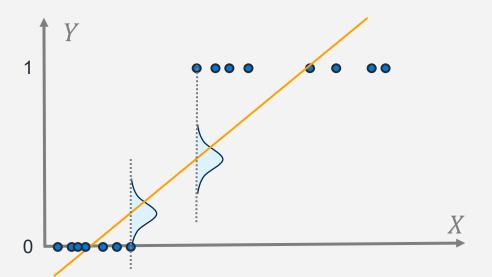




1. Modèles Linéaires

Rappel

Limites





Exemple illustratif

Limites des modèles linéaires classiques

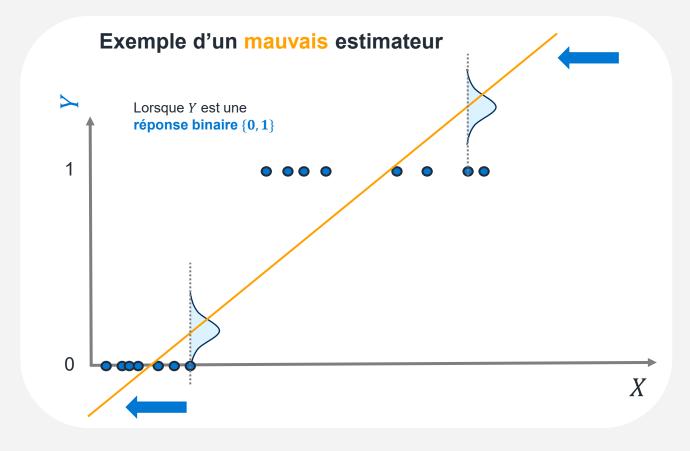
Cas variable cible binaire

▶ Modèle linéaire multiple $Y = X\beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2)$

Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
:	:	ŧ
1203.1	1	3

Lorsque *Y* n'est pas une variable cible continue alors le **modèle linéaire classique** n'a plus de sens :

- 1) L'estimateur à pour valeur dans un ensemble discret
- 2) Les **hypothèses sur les résidus** (notamment de variance constante) ne sont pas vérifiées





Cas variable cible binaire

Solution : Prédire la probabilité d'un événement binaire (0 ou 1) à partir de variables explicatives X = x.

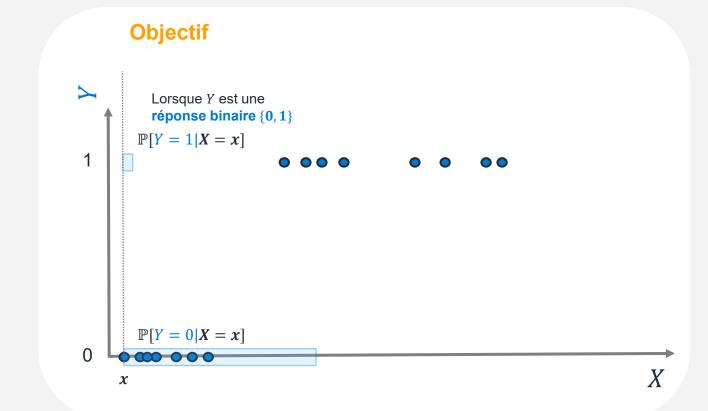
Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
÷	ŧ	÷
1203.1	1	3

Modèle linéaire : $\mathbb{E}[Y|X] = X\beta^T$

Réponse binaire : $\mathbb{E}[Y|X] = \mathbb{P}(Y=1 \mid X) = p$

Distribution « naturelle » : $Y \mid X \sim Ber(p)$

Lien ? $p \leftrightarrow \beta$



Cas variable cible binaire

Solution : Prédire la probabilité d'un événement binaire (0 ou 1) à partir de variables explicatives X = x.

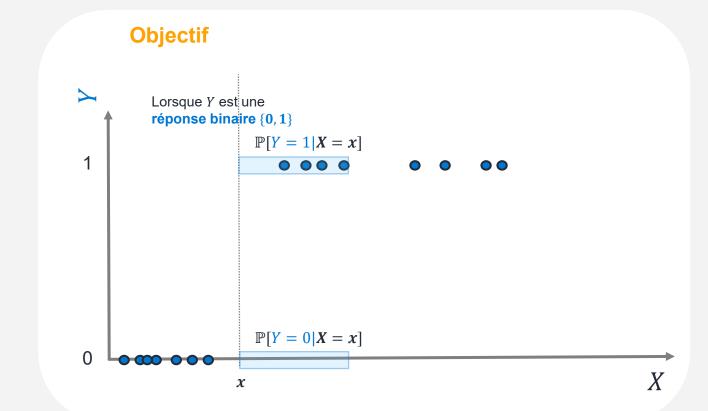
Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
ŧ	:	ŧ
1203.1	1	3

Modèle linéaire : $\mathbb{E}[Y|X] = X\beta^T$

Réponse binaire : $\mathbb{E}[Y|X] = \mathbb{P}(Y=1 \mid X) = p$

Distribution « naturelle » : $Y \mid X \sim Ber(p)$

Lien ? $p \leftrightarrow \beta$



Limites des modèles linéaires classiques

Exemple illustratif

Cas variable cible binaire

Solution : Prédire la probabilité d'un événement binaire (0 ou 1) à partir de variables explicatives X = x.

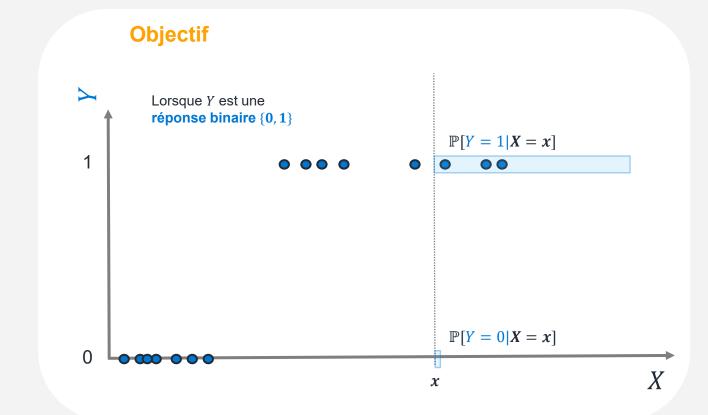
Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
÷	ŧ	÷
1203.1	1	3

Modèle linéaire : $\mathbb{E}[Y|X] = X\beta^T$

Réponse binaire : $\mathbb{E}[Y|X] = \mathbb{P}(Y=1 \mid X) = p$

Distribution « naturelle » : $Y \mid X \sim Ber(p)$

Lien ? $p \leftrightarrow \beta$



Exemple illustratif

Limites des modèles linéaires classiques

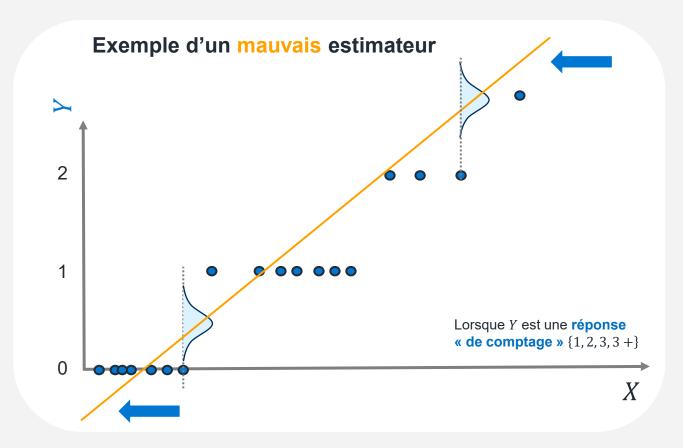
Cas variable cible discrète

Modèle linéaire multiple $Y = X \beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2)$

Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
i	:	ŧ
1203.1	1	3

Lorsque *Y* n'est pas une variable cible continue alors le **modèle linéaire classique** n'a plus de sens :

- 1) L'estimateur à pour valeur dans un ensemble discret
- Les hypothèses sur les résidus (notamment de variance constante) ne sont pas vérifiées





Limites des modèles linéaires classiques

Exemple illustratif

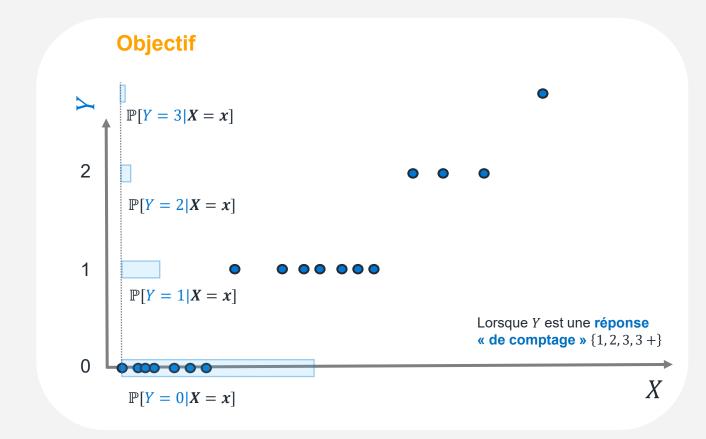
Cas variable cible discrète

Solution: Prédire la **probabilité d'un événement discret** à partir de variables explicatives X = x.

Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
:	:	ŧ
1203.1	1	3

Modèle linéaire : $\mathbb{E}[Y|X] = X\beta^T$

Réponse discrète positive : $\mathbb{P}(Y = k \mid X) = ?$





Cas variable cible discrète

Solution: Prédire la **probabilité d'un événement discret** à partir de variables explicatives X = x.

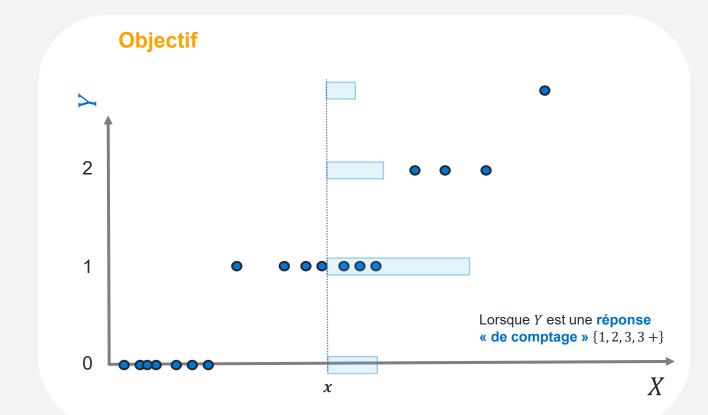
Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
:	:	:
1203.1	1	3

Modèle linéaire : $\mathbb{E}[Y|X] = X\beta^T$

Distribution « naturelle » : $Y \mid X \sim Poisson(\lambda)$

Réponse discrète positive : $\mathbb{P}(Y = k \mid X) = \frac{\lambda^k}{k!} e^{-\lambda}$

Lien ? $\lambda \leftrightarrow \beta$





Limites des modèles linéaires classiques

Exemple illustratif

Cas variable cible discrète

Solution: Prédire la **probabilité d'un événement discret** à partir de variables explicatives X = x.

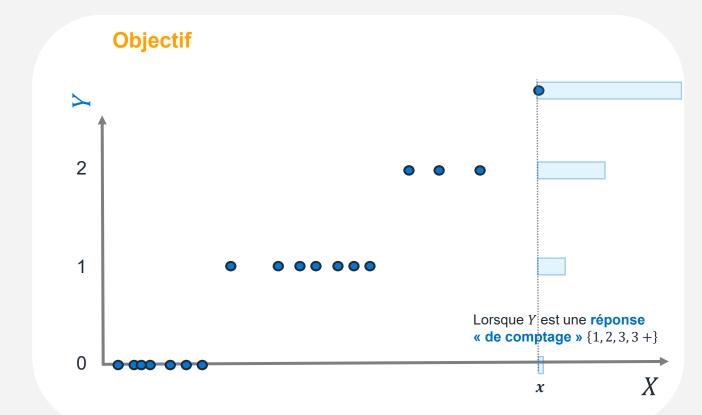
Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
:	:	ŧ
1203.1	1	3

Modèle linéaire : $\mathbb{E}[Y|X] = X\beta^T$

Distribution « naturelle » : $Y \mid X \sim Poisson(\lambda)$

Réponse discrète positive : $\mathbb{P}(Y = k \mid X) = \frac{\lambda^k}{k!} e^{-\lambda}$

Lien ? $\lambda \leftrightarrow \beta$





2. Distributions de la famille exponentielle



Introduction à la famille exponentielle (1/3)

Loi Normale : Modélisation des données continues avec variance constante

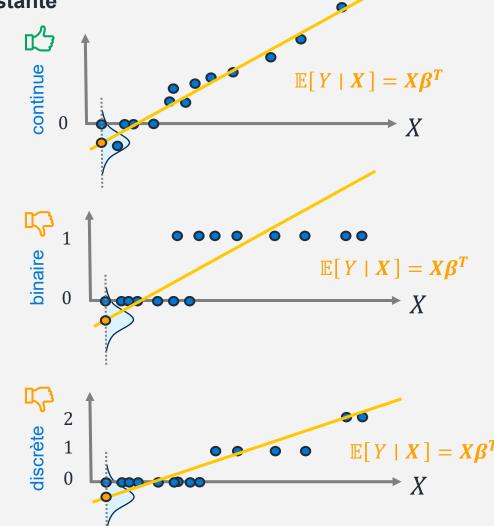
Modèle linéaire multiple

$$Y = X\beta^T + \varepsilon$$
 avec $\varepsilon_i iid \sim N(0, \sigma^2)$

donc: $Y \mid X \sim N(X\beta^T, \sigma^2 I)$

avec:

$$f(y) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{y - X\beta^T}{\sigma}\right)^2\right)$$





Introduction à la famille exponentielle (2/3)

Famille exponentielle : Des données discrètes ou continues avec variance constante ou non constante

Modèle linéaire multiple

$$Y = X\beta^T + \varepsilon$$
 avec $\varepsilon_i iid \sim N(0, \sigma^2)$

donc:
$$Y \mid X \sim N(X\beta^T, \sigma^2 I)$$

avec:

$$f(y) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{y - X\beta^T}{\sigma}\right)^2\right)$$

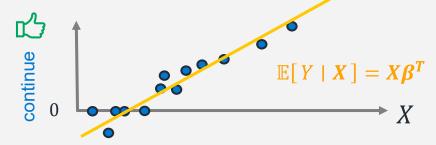
Extension : Les distributions de la famille exponentielle sont un ensemble de distributions paramétriques :

$$f(y) = c(y, \phi) \exp\left(\frac{y\beta - a(\beta)}{\phi}\right)$$

Incluent des distributions couramment utilisées pour modéliser divers types de données : **continues**, **discrètes**, et **binaires**.

$$\mathbb{E}[Y \mid X] = a'(\beta)$$

$$Var[Y|X] = \phi \times a''(\beta)$$









Introduction à la famille exponentielle (3/3)

Famille exponentielle : Des données discrètes ou continues avec variance constante ou non constante

Distribution de la famille exponentielle :

$$f(y) = c(y, \phi) \exp\left(\frac{y\beta - a(\beta)}{\phi}\right)$$

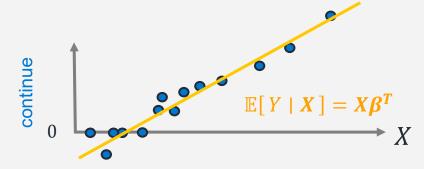
Pour ajuster les données en conséquence, nous devons les transformer à l'aide de fonctions appelées **fonctions de lien**, notées g:

$$g(\mathbb{E}[Y|X]) = X\boldsymbol{\beta}^T$$

Nous détaillerons dans **la section suivante** les graphiques et les résultats statistiques. Notamment :

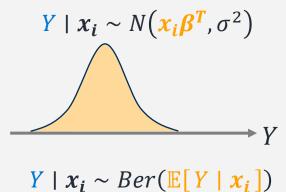
 $g_{pernoulli}$?

 $g_{poisson}$?

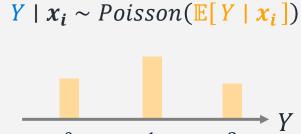














3. Modèles GLM

Modèle Logistique

Modèle de Poisson

Récapitulatif



Exemple illustratif

Variable cible binaire

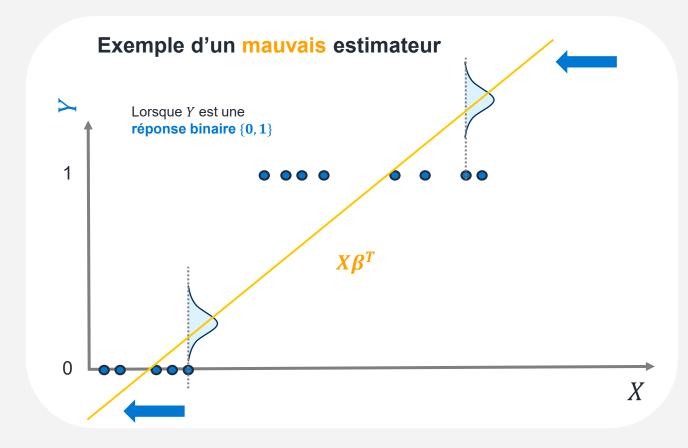
Limites des modèles linéaires classiques

▶ Modèle linéaire multiple $Y = X\beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2)$ $\mu \coloneqq \mathbb{E}[Y|X] = X\beta^T$ et $\sigma^2 \coloneqq Var[Y|X]$

Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
i	ŧ	÷
1203.1	1	3

Lorsque Y est une **réponse binaire** $\{0, 1\}$, nous avons clairement plusieurs problèmes:

- L'estimateur **a pour valeur dans** ℝ
- Les **hypothèses sur les résidus** (notamment de variance constante) ne sont pas vérifiées



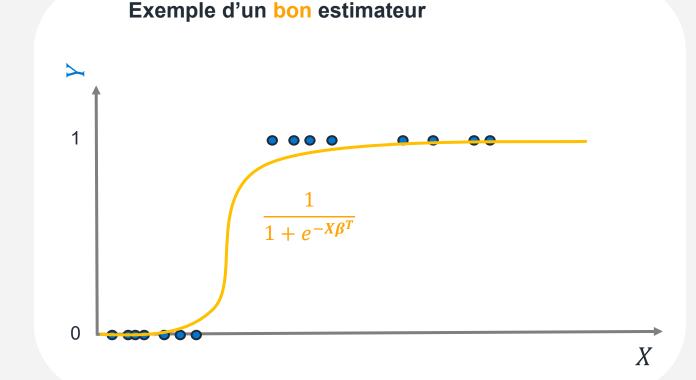


Modèle Logistique

Modèle linéaire multiple $Y = X\beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2) \ \mu \coloneqq \mathbb{E}[Y|X] = X\beta^T$ et $\sigma^2 \coloneqq Var[Y|X]$

Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
÷	ŧ	÷
1203.1	1	3

▶ Régression Logistique. Modèle statistique utilisée pour prédire la probabilité d'un événement binaire (0 ou 1) à partir de variables explicatives. Elle est souvent employée pour des tâches de classification.



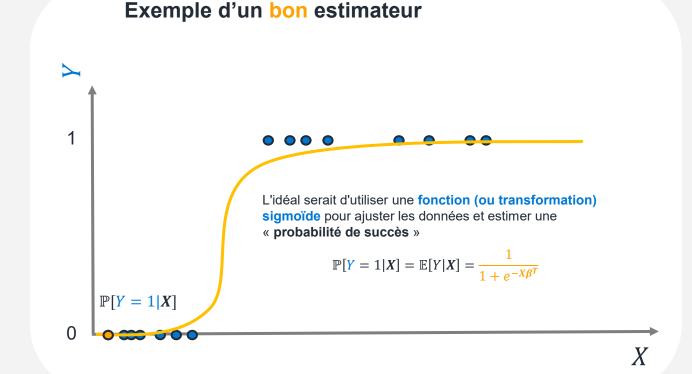


Modèle Logistique

Modèle linéaire multiple $Y = X\beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2) \ \mu \coloneqq \mathbb{E}[Y|X] = X\beta^T$ et $\sigma^2 \coloneqq Var[Y|X]$

Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
÷	ŧ	i
1203.1	1	3

Régression Logistique. Modèle statistique utilisée pour prédire la probabilité d'un événement binaire (0 ou 1) à partir de variables explicatives. Elle est souvent employée pour des tâches de classification.





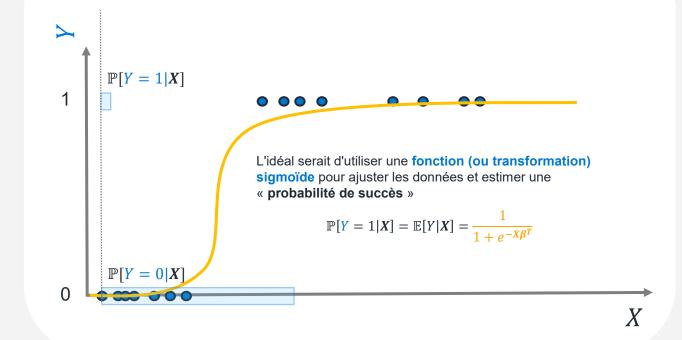
Modèle Logistique

Modèle linéaire multiple $Y = X\beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2) \ \mu \coloneqq \mathbb{E}[Y|X] = X\beta^T \ \text{et} \ \sigma^2 \coloneqq Var[Y|X]$

Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
ŧ	:	÷
1203.1	1	3

▶ Régression Logistique. Modèle statistique utilisée pour prédire la probabilité d'un événement binaire (0 ou 1) à partir de variables explicatives. Elle est souvent employée pour des tâches de classification.

Exemple d'un bon estimateur



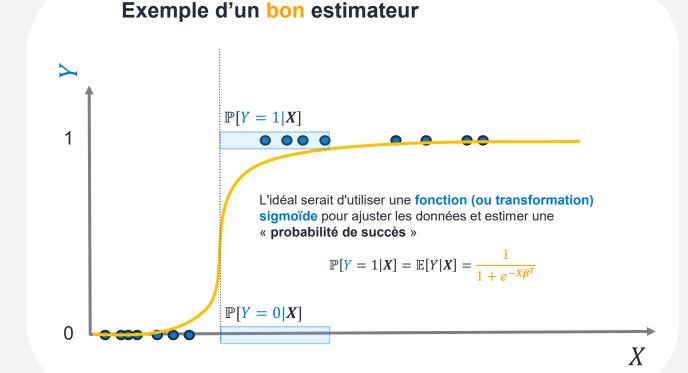


Modèle Logistique

Modèle linéaire multiple $Y = X\beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2) \ \mu \coloneqq \mathbb{E}[Y|X] = X\beta^T$ et $\sigma^2 \coloneqq Var[Y|X]$

Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
i	:	i
1203.1	1	3

Régression Logistique. Modèle statistique utilisée pour prédire la probabilité d'un événement binaire (0 ou 1) à partir de variables explicatives. Elle est souvent employée pour des tâches de classification.



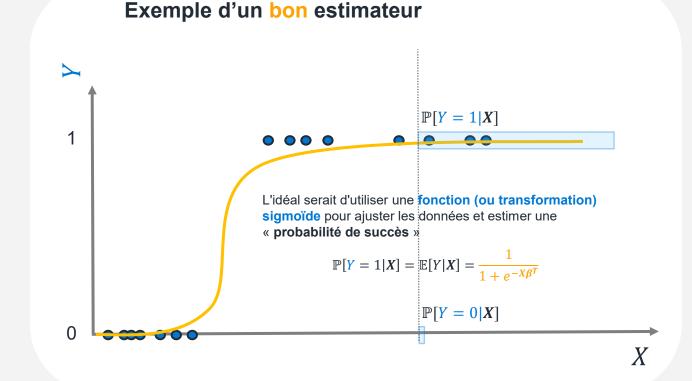


Modèle Logistique

Modèle linéaire multiple $Y = X\beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2) \ \mu \coloneqq \mathbb{E}[Y|X] = X\beta^T \ \text{et} \ \sigma^2 \coloneqq Var[Y|X]$

Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
ŧ	:	ŧ
1203.1	1	3

Régression Logistique. Modèle statistique utilisée pour prédire la probabilité d'un événement binaire (0 ou 1) à partir de variables explicatives. Elle est souvent employée pour des tâches de classification.





Modèle Logistique

Modèle linéaire multiple $Y = X\beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2) \ \mu \coloneqq \mathbb{E}[Y|X] = X\beta^T$ et $\sigma^2 \coloneqq Var[Y|X]$

Régression Logistique. Une technique statistique utilisée pour prédire la probabilité d'un événement binaire (0 ou 1) à partir de variables explicatives. Elle est souvent employée pour des tâches de classification.

Pour ajuster les données, nous devons les transformer à l'aide de fonctions appelées **fonctions de lien**, notées g:

$$g(\mathbb{E}[Y|X]) = X\beta^T$$

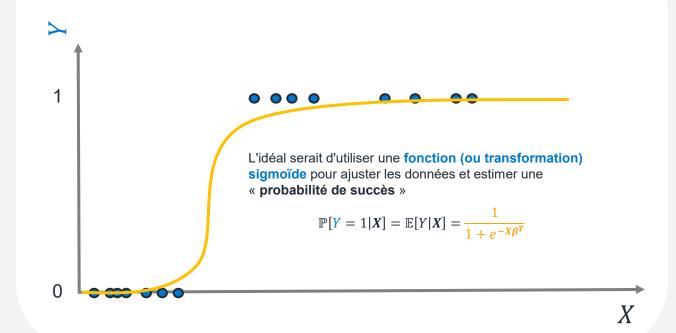
Fonction logistique comme fonction de lien :

$$g(\mathbb{E}[Y|X]) = \log\left(\frac{\mathbb{E}[Y|X]}{1 - \mathbb{E}[Y|X]}\right)$$

avec

$$\mathbb{E}[Y|X] = \frac{1}{1 + e^{-X\beta^T}}$$

Exemple d'un bon estimateur





3. Modèles GLM

Modèle Logistique

Modèle de Poisson

Récapitulatif



Exemple illustratif

variable cible de comptage

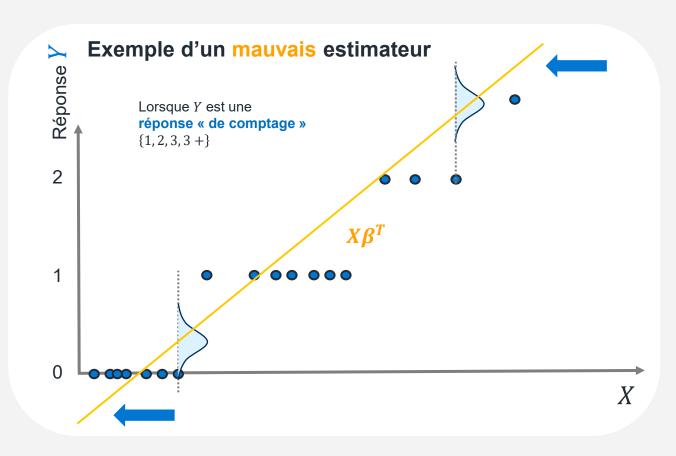
Cas variable cible discrète

▶ Modèle linéaire multiple $Y = X\beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2)$ $\mu \coloneqq \mathbb{E}[Y|X] = X\beta^T$ et $\sigma^2 \coloneqq \mathrm{Var}[Y|X]$

Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
:	:	÷
1203.1	1	3

Lorsque *Y* est une **réponse discrète (ici, de comptage)** nous avons clairement plusieurs problèmes:

- Les valeurs estimées devraient être des nombres entiers
- Les **hypothèses sur les résidus** (notamment de variance homogène, normalité etc. ...) ne sont pas vérifiées



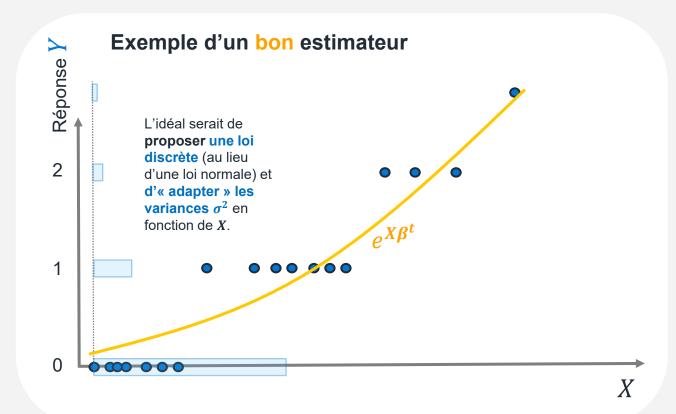


Cas variable cible discrète

Modèle linéaire multiple $Y = X\beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2)$ $\mu \coloneqq \mathbb{E}[Y|X] = X\beta^T$ et $\sigma^2 \coloneqq \mathrm{Var}[Y|X]$

Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
:	:	:
1203.1	1	3

Régression de Poisson. Pour modéliser le nombre de fois qu'un événement se produit dans un intervalle de temps ou d'espace, en fonction de variables explicatives. Elle est particulièrement adaptée aux données de comptage.



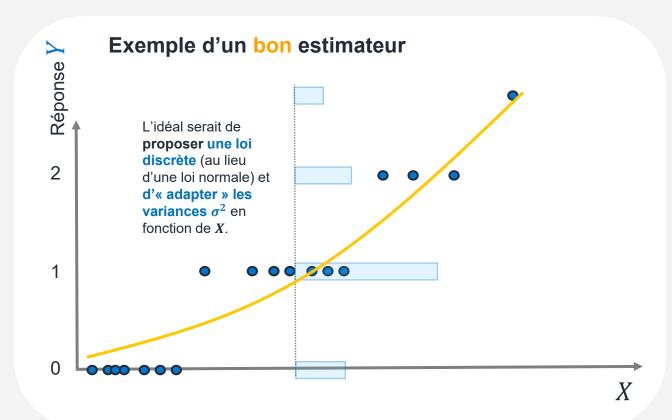


Cas variable cible discrète

Modèle linéaire multiple $Y = X\beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2)$ $\mu \coloneqq \mathbb{E}[Y|X] = X\beta^T$ et $\sigma^2 \coloneqq \mathrm{Var}[Y|X]$

Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
:	:	ŧ
1203.1	1	3

► Régression de Poisson. Pour modéliser le nombre de fois qu'un événement se produit dans un intervalle de temps ou d'espace, en fonction de variables explicatives. Elle est particulièrement adaptée aux données de comptage.



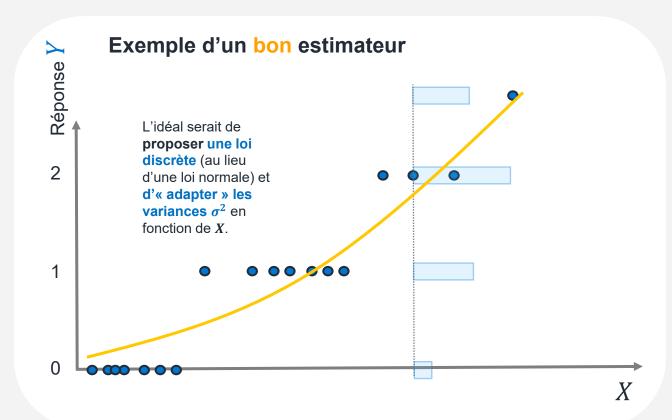


Cas variable cible discrète

Modèle linéaire multiple $Y = X\beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2)$ $\mu \coloneqq \mathbb{E}[Y|X] = X\beta^T$ et $\sigma^2 \coloneqq \mathrm{Var}[Y|X]$

Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
:	:	ŧ
1203.1	1	3

► Régression de Poisson. Pour modéliser le nombre de fois qu'un événement se produit dans un intervalle de temps ou d'espace, en fonction de variables explicatives. Elle est particulièrement adaptée aux données de comptage.



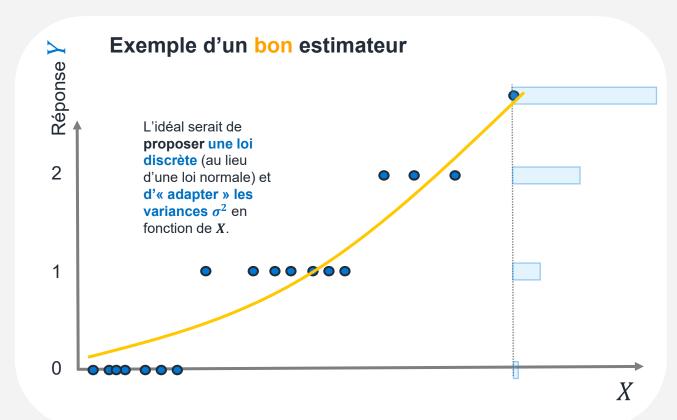


Cas variable cible discrète

Modèle linéaire multiple $Y = X\beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2)$ $\mu \coloneqq \mathbb{E}[Y|X] = X\beta^T$ et $\sigma^2 \coloneqq \mathrm{Var}[Y|X]$

Coût sinistre	Occurrence sinistre	Nombre sinistres
0	0	0
122.3	1	1
977.8	1	2
i	i	÷
1203.1	1	3

Régression de Poisson. Pour modéliser le nombre de fois qu'un événement se produit dans un intervalle de temps ou d'espace, en fonction de variables explicatives. Elle est particulièrement adaptée aux données de comptage.





Modèle de Poisson

Modèle linéaire multiple $Y = X\beta^T + \varepsilon$ avec $\varepsilon_i \ iid \sim N(0, \sigma^2)$ $\mu \coloneqq \mathbb{E}[Y|X] = X\beta^T$ et $\sigma^2 \coloneqq \mathrm{Var}[Y|X]$

Régression de Poisson. Une méthode statistique utilisée pour modéliser le nombre de fois qu'un événement se produit dans un intervalle de temps ou d'espace, en fonction de variables explicatives. Elle est particulièrement adaptée aux données de comptage:

$$Y \mid X \sim Pois(\lambda) \text{ avec } \lambda = \mathbb{E}[Y \mid X] = e^{X\beta^t} \text{ et}$$

$$\mathbb{P}(Y = k \mid X) = \frac{\lambda^k}{k!} e^{-k}$$

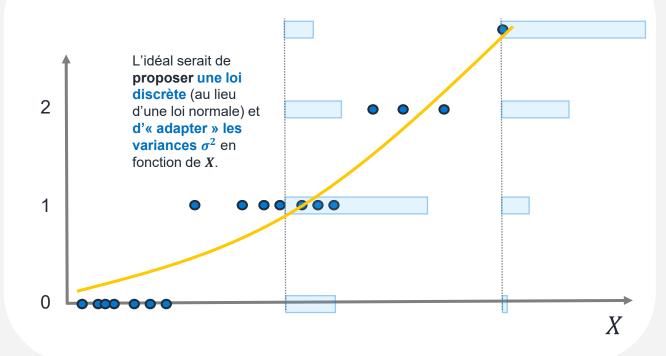
La fonction logarithmique comme fonction de lien :

$$g(\mathbb{E}[Y|X]) = \log(\mathbb{E}[Y|X]) = X\beta^{T}$$

avec

$$\mathbb{E}[Y|X] = e^{X\beta^t}$$

Exemple d'un bon estimateur





3. Modèles GLM

Modèle Logistique

Modèle de Poisson

Récapitulatif



Formulation générale

Modèles linéaires généralisés

Modèle linéaire généralisé ou GLM (pour Generalized Linear Model) étend le modèle linéaire « classique » en :

1) utilisant une fonction de lien (inversible) pour relier la variable de réponse au modèle

$$\mathbb{E}[Y|X] = g^{-1}(X\boldsymbol{\beta}^T)$$

2) et en permettant à la variance de dépendre de la valeur prédite selon la distribution choisie (V une fonction)

$$Var[Y|X] = \phi \times V(g^{-1}(X\beta^T))$$

avec V une fonction variance qui dépend $\mathbb{E}[Y \mid X]$ et ϕ un paramètre d'échelle.

Modèle classique

Distribution des erreurs : Normale

$$\mathbb{E}[Y|X] = X\boldsymbol{\beta}^T$$

$$Var[Y|X] = \sigma^2$$

Modèle logistique

Distribution : Bernoulli (ou Binomiale)

$$\mathbb{E}[Y|X] = \frac{1}{1 + e^{-X\beta^T}}$$

$$\mathsf{Var}[\underline{Y}|\underline{X}] = \mu(1-\mu)$$

Modèle de Poisson

Distribution des erreurs : Poisson

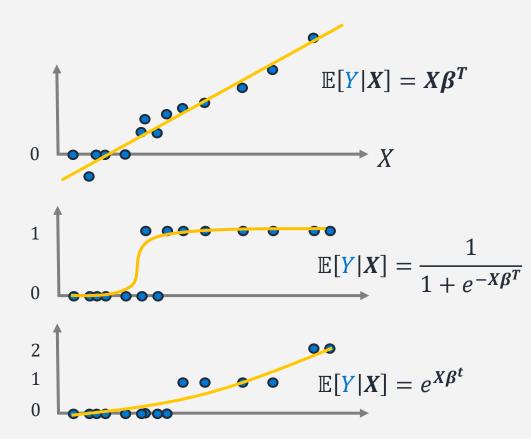
$$\mathbb{E}[Y|X] = e^{X\beta^t}$$

$$Var[Y|X] = \mathbb{E}[Y \mid X] = \mu$$



Modèles GLM et applications en assurance Récapitulatif

	Modèle Linéaire « classique »	Modèle Logistique	Modèle de Poisson
Variable continue (coût sinistre)	~		
Variable binaire (sinistre oui/non)		>	
Variable discrète (nombre ou fréquence)			~





4. Estimation des paramètres : Maximum de vraisemblance



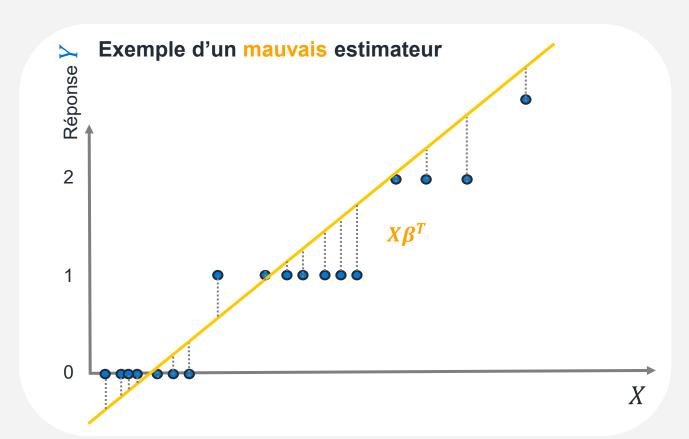
Fonction de vraisemblance (1/2)

Les limites de la méthode des moindres carrés

La méthode des moindres carrés

$$R(\boldsymbol{\beta}) \coloneqq \sum_{i} (Y_i - \widehat{Y}_i)^2 = \sum_{i} \hat{\varepsilon}_i^2$$

Adaptée aux relations linéaires entre les variables





Fonction de vraisemblance (1/2)

Les limites de la méthode des moindres carrés

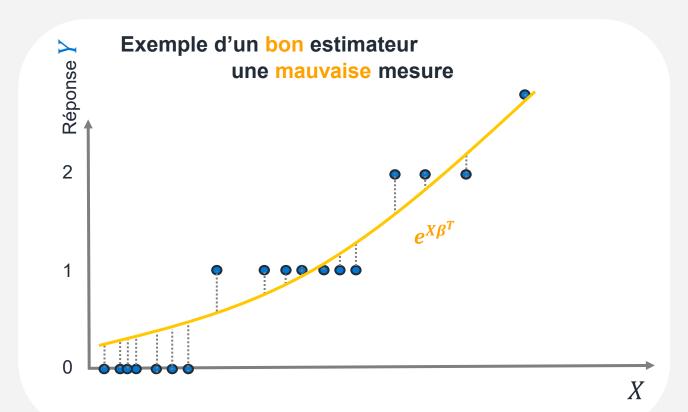
La méthode des moindres carrés

$$R(\boldsymbol{\beta}) \coloneqq \sum_{i} (Y_i - \widehat{Y}_i)^2 = \sum_{i} \hat{\varepsilon}_i^2$$

Adaptée aux relations linéaires entre les variables

以

En présence d'une **relation non linéaire**, son utilisation peut entraîner des erreurs significatives





Fonction de vraisemblance (2/2)

Les limites de la méthode des moindres carrés

La méthode des moindres carrés

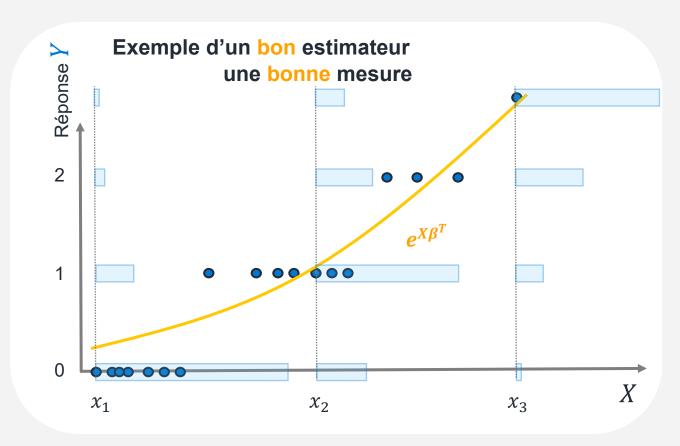
$$R(\boldsymbol{\beta}) \coloneqq \sum_{i} (Y_i - \widehat{Y}_i)^2 = \sum_{i} \hat{\varepsilon}_i^2$$

Adaptée aux relations linéaires entre les variables



Estimer les paramètres de la fonction de densité / probabilité

$$f_{\boldsymbol{\beta}}(y_i) \coloneqq \mathbb{P}_{\boldsymbol{\beta}}(Y_i = y_i \mid \boldsymbol{X})$$





Fonction de vraisemblance (2/2)

Les limites de la méthode des moindres carrés

La méthode des moindres carrés

$$R(\boldsymbol{\beta}) \coloneqq \sum_{i} (Y_i - \widehat{Y}_i)^2 = \sum_{i} \hat{\varepsilon}_i^2$$

Adaptée aux relations linéaires entre les variables



En présence d'une **relation non linéaire**, son utilisation peut entraîner des erreurs significatives

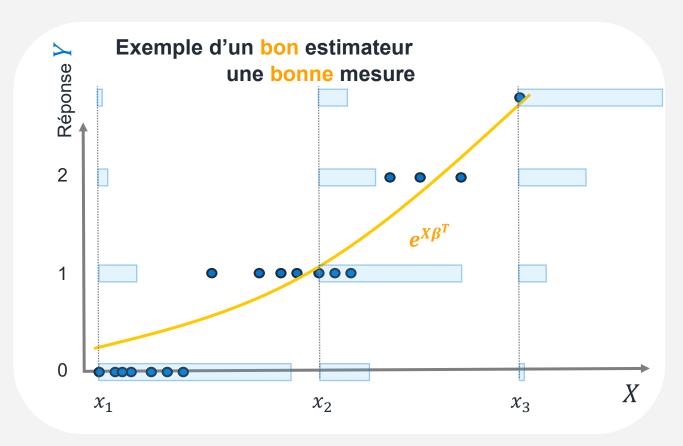
Estimer les paramètres de la fonction de densité / probabilité

$$f_{\boldsymbol{\beta}}(y_i) \coloneqq \mathbb{P}_{\boldsymbol{\beta}}(Y_i = y_i \mid \boldsymbol{X})$$

Fonction de vraisemblance (ou log de vraisemblance)

$$L(\mathbf{y}, \boldsymbol{\beta}) \coloneqq \prod_{i=1,\dots,n} f_{\boldsymbol{\beta}}(y_i)$$

$$l(\mathbf{y}, \boldsymbol{\beta}) \coloneqq \log L(\mathbf{y}, \boldsymbol{\beta}) = \sum_{i=1,\dots,n} \log f_{\boldsymbol{\beta}}(y_i)$$





Maximum de vraisemblance

Cas modèle de Poisson

Estimateur du maximum de vraisemblance (MLE en anglais) : estimateur statistique utilisé pour inférer les paramètres de la loi de probabilité d'un échantillon donné en recherchant les valeurs des paramètres maximisant la fonction de vraisemblance (ou log de vraisemblance)

$$\beta_{GLM} = \arg \max_{\beta} l(y, \beta) = \arg \max_{\beta} \sum_{i=1,\dots,n} \log f_{\beta}(y_i)$$

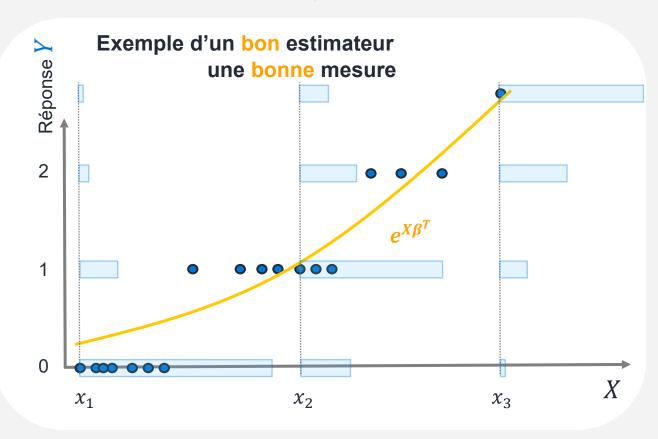
Modèle de Poisson

$$\sum_{i=1,\dots,n} \log f_{\boldsymbol{\beta}}(y_i) = \sum_{i=1,\dots,n} \log \frac{\left(\boldsymbol{x_i}\boldsymbol{\beta}^T\right)^{y_i}}{y_i!} e^{-y_i}$$
$$= \sum_{i=1,\dots,n} \left(y_i \boldsymbol{x_i}\boldsymbol{\beta}^T - e^{\boldsymbol{x_i}\boldsymbol{\beta}^T} - \log y_i! \right)$$

Comme nous cherchons β qui maximise $l(y, \beta)$, nous prenons la simplification suivante :

$$\boldsymbol{\beta}_{GLM} = \arg \max_{\boldsymbol{\beta}} \sum_{i=1,...,n} \left(y_i x_i \boldsymbol{\beta}^T - e^{x_i \boldsymbol{\beta}^T} \right)$$

Solution numérique : Descente de gradient





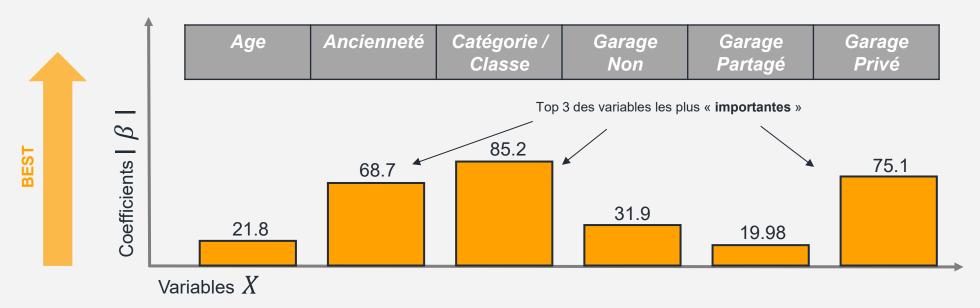
5. Interprétation des coefficients / paramètres



Importance des variables explicatives (1/2)

Des approches « simples » basées sur les coefficients

Exemple illustratif



Approche 1 : magnitude des coefficients

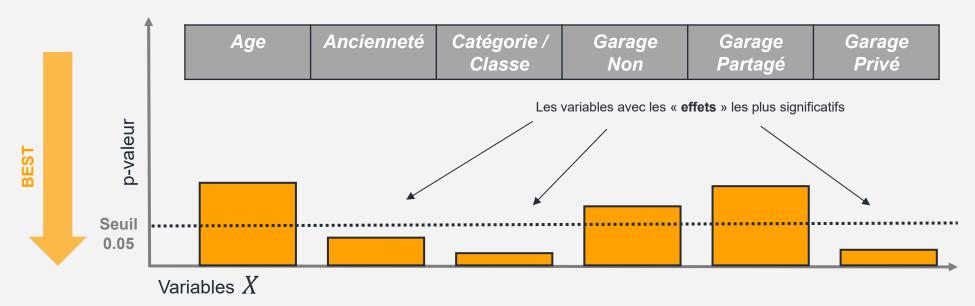
Les coefficients estimés mesurent l'effet d'une variation unitaire d'une variable explicative sur la variable cible, en maintenant constantes les autres variables. Dans certains cas, les coefficients peuvent être utilisés pour évaluer l'importance relative des variables. Plus un coefficient est grand (en valeur absolue), plus l'impact de la variable correspondante est important.



Importance des variables explicatives (2/2)

Des approches « simples » basées sur les coefficients

Exemple illustratif



Approche 1 : magnitude des coefficients

Les coefficients estimés mesurent l'effet d'une variation unitaire d'une variable explicative sur la variable cible, en maintenant constantes les autres variables. Dans certains cas, les coefficients peuvent être utilisés pour évaluer l'importance relative des variables. Plus un coefficient est grand (en valeur absolue), plus l'impact de la variable correspondante est important.

Approche 2 : signification statistique (p-valeur)

Une autre manière d'évaluer l'importance d'une variable est de tester si son coefficient est statistiquement différent de zéro. Si la **p-valeur** associée au coefficient est inférieure à un seuil conventionnel (**souvent 0.05**), la variable est considérée comme ayant un effet significatif sur la variable cible.

Nous détaillerons plus tard les tests statistiques



Zoom sur le modèle Logistique

Cas assurance auto

77.0 100	lustratif

Coût sinistre	Occurrence sinistre	Age	Ancienneté	Catégorie / Classe	Garage Non	Garage Partagé	Garage Privé
0	0	0.84	0	0.666	0	0	1
122.3	1	0	1	0.333	0	0	1
9/7.8	1	0.04	0.875	0.333	1	0	0
0	0	1	0.625	0.333	0	1	0
311.1	1	0.36	0.25	1	1	0	0



Tâche de classification avec une variable cible binaire

Rapport des cotes (Odds-Ratio ou OR) pour des **prédicteurs binaires** : C'est le rapport des cotes des probabilités que ce soit positif pour ceux qui ont la caractéristique X_i d'une part et de ceux qui ne l'ont pas d'autre part. Exemple avec

$$OR_{i} = \frac{\frac{\mathbb{P}(Y = 1 \mid X_{i} = 1)}{\mathbb{P}(Y = 0 \mid X_{i} = 1)}}{\frac{\mathbb{P}(Y = 1 \mid X_{i} = 0)}{\mathbb{P}(Y = 0 \mid X_{i} = 0)}} = \exp(\beta_{i})$$

Les cotes (ou *odds*) pour X_i :

$$e^{X\beta^T} = \frac{\mathbb{P}(Y=1 \mid X_i)}{1 - \mathbb{P}(Y=1 \mid X_i)} = \frac{\mathbb{P}(Y=1 \mid X_i)}{\mathbb{P}(Y=0 \mid X_i)}$$

- $OR_i = 1$, la variable cible est indépendantes de X_i
- $OR_i > 1$, la variable cible est plus fréquente pour les individus **qui ont** X_i
- $OR_i < 1$, la variable cible est plus fréquente pour les individus **qui n'ont pas** X_i

6. Evaluation des GLM



Rappel

Matrice de confusion et AUC

Occurrence sinistre	Score	Prédiction (seuil 0.5)	Prédiction (seuil 1)	Prédiction (seuil 0.65)
0	0.4	0	0	0
1	0.7	1	0	1
1	0.35	0	0	0
0	0.52	1	0	0
1	0.68	1	0	1
0	0.2	0	0	0
0	0.1	0	0	0
0	0.6	1	0	0
1	0.7	1	0	1
0	0.55	1	0	0

Accuracy: 0.6 Accuracy: 0.6 Accuracy: 0.9 (0, 0)

(TFP, TVP)

(0.25, 0.5)

(0, 0.75)

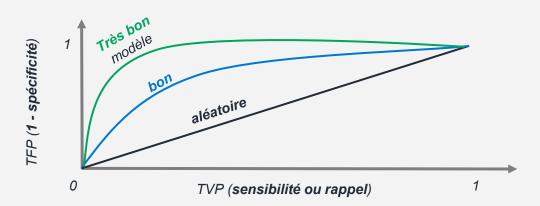
Exemple illustratif

Approche 1 : Matrice de confusion

Mesure la performance des modèles de classification à 2 classes ou plus. Dans le cas binaire, la matrice de confusion est un tableau à 4 valeurs représentant différentes combinaisons de valeurs réelles et valeurs prédites

Approche 2 : AUC-ROC

L'AUC-ROC (Area Under the Curve - Receiver Operating Characteristic) ou AUC évalue la capacité d'un modèle à distinguer entre les classes en traçant le TVP contre le TFP pour différents seuils de décision. L'AUC (aire sous la courbe) varie entre 0,5 (modèle aléatoire) et 1 (modèle parfait)



Indice de GINI = 2 x AUC - 1



Des mesures « agnostiques » aux modèles (s'appliquent à toutes techniques : GLM, arbres de décision, ...)

Evaluation empirique du modèle (1/3)

Coefficient de détermination linéaire de Pearson

Le coefficient de détermination \mathbb{R}^2 mesure la proportion de la variance de la variable cible expliquée par les variables indépendantes dans un modèle linéaire. Il est utilisé pour évaluer la qualité d'ajustement d'un modèle linéaire. Il varie entre 0 (mauvais ajustement) et 1 (ajustement parfait)

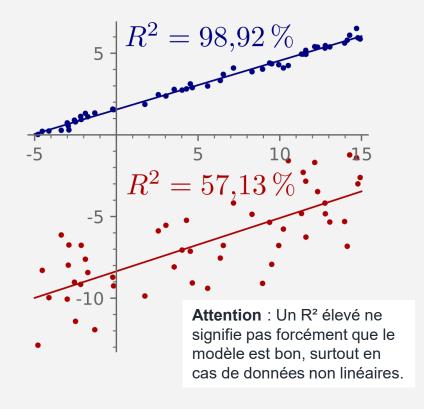
$$R^2 = 1 - \frac{\sum (Y_i - \widehat{Y}_i)^2}{\sum (Y_i - \overline{Y})^2}$$

Avec Y_i valeur observée, \widehat{Y}_i valeur estimée et \overline{Y} valeur moyenne

Interprétation:

- $R^2 = 1$: Le modèle explique **toute** la variance des données
- $R^2 = 0$: Le modèle n'explique **aucune** variance des données.
- $R^2 < 0$: Le modèle est **pire** qu'un modèle constant (moyenne).

 $R^2 = 0.85$ signifie que 85 % de la variance de la variable cible est expliquée par le modèle





Evaluation empirique du modèle (1/3)

Coefficient de détermination linéaire de Pearson

Le coefficient de détermination \mathbb{R}^2 mesure la proportion de la variance de la variable cible expliquée par les variables indépendantes dans un modèle linéaire. Il est utilisé pour évaluer la qualité d'ajustement d'un modèle linéaire. Il varie entre 0 (mauvais ajustement) et 1 (ajustement parfait)

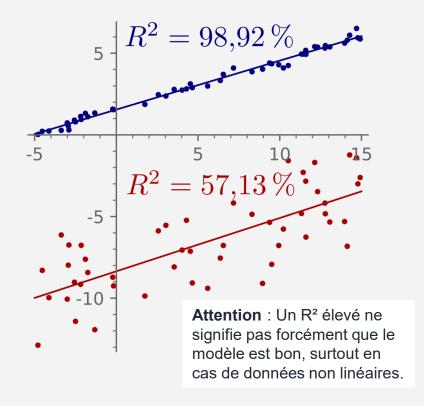
$$R^2 = 1 - \frac{\sum (Y_i - \widehat{Y}_i)^2}{\sum (Y_i - \overline{Y})^2}$$

Avec Y_i valeur observée, \widehat{Y}_i valeur estimée et \overline{Y} valeur moyenne

Interprétation:

- $R^2 = 1$: Le modèle explique **toute** la variance des données
- $R^2 = 0$: Le modèle n'explique **aucune** variance des données.
- $R^2 < 0$: Le modèle est **pire** qu'un modèle constant (moyenne).

 $R^2 = 0.85$ signifie que 85 % de la variance de la variable cible est expliquée par le modèle



Problème avec GLM: le R^2 des modèles linéaires est basé sur la somme des carrées des résidus et ne tient pas compte du **surajustement**. Il n'y a donc pas de R^2 dans les GLM qui sont basés sur le **Maximum de Vraisemblance**.

Les pseudos- R^2 sont des équivalents du R^2 dans les GLM

Pseudo-
$$R^2 = 1 - \frac{\text{déviance modèle}}{\text{déviance nulle}} = 1 - \frac{(-2 \log \text{vraisemblance modèle})}{(-2 \log \text{vraisemblance nulle})}$$



Evaluation empirique du modèle (2/3)

Les pseudos R^2 : quelques indicateurs

Comparer le modèle avec le modèle initial (trivial) constitué de la seule constante. Soit $L_{modèle}$ la vraisemblance maximale du modèle ajusté (avec les variables explicatives) et L_0 la vraisemblance maximale du modèle nul (sans variables explicatives, uniquement avec l'intercept).

R^2 de Cox et Snell

mesure pseudo-R² adaptée aux modèles de régression non linéaire, comme la régression logistique. Il est basé sur le rapport de vraisemblance et exprime la proportion de la variance expliquée par le modèle. Formulé comme :

$$R_{CS}^2 = 1 - \left(\frac{L_0}{L_{mod\`{e}le}}\right)^{\frac{2}{n}}$$

Ce R^2 est borné par 1, mais ne l'atteint jamais exactement, ce qui peut limiter son interprétation en comparaison avec d'autres mesures comme le R^2 de Nagelkerke.

$$R_N^2 = \frac{R_{CS}^2}{\max(R_{CS}^2)}$$

R² de McFadden

Evalue la qualité d'ajustement d'un modèle à partir de la logvraisemblance.

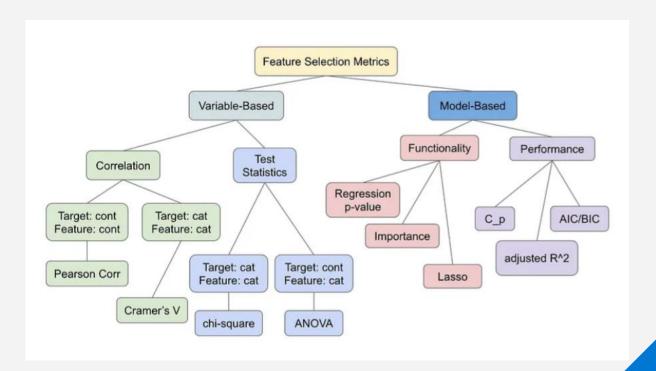
$$R_{MF}^2 = 1 - \frac{\log(L_{mod \ge le})}{\log(L_0)}$$

Par rapport au R^2 de Cox et Snell ou à celui de Nagelkerke, le R^2 de McFadden est **plus couramment utilisé pour les modèles logistiques** car il repose sur des principes de vraisemblance qui s'ajustent bien à ces modèles.

Souvent préféré pour la régression logistique en raison de son lien direct avec le test du rapport de vraisemblance.



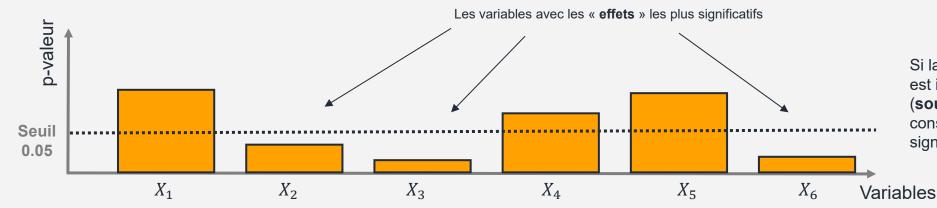
7. Sélection des variables





Test du rapport de vraisemblance

Exemple illustratif



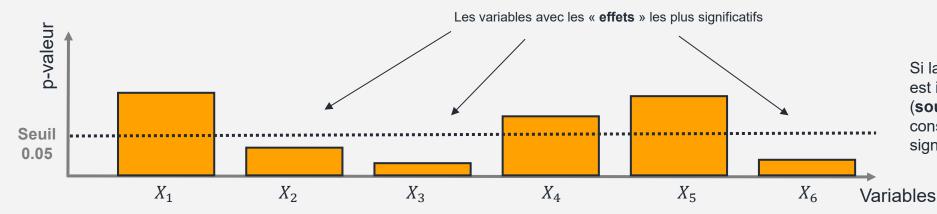
Si la **p-valeur** associée au coefficient est inférieure à un seuil conventionnel (**souvent 0.05**), la variable est considérée comme ayant un effet significatif sur la variable cible.

Test statistique : $m{\beta_j} = \mathbf{0}$ vs $m{\beta_j} \neq \mathbf{0}$ avec $j \in \{1, ..., d\}$ Le prédicteur n'est pas significatif dans le GLM



Test du rapport de vraisemblance

Exemple illustratif



Si la **p-valeur** associée au coefficient est inférieure à un seuil conventionnel (**souvent 0.05**), la variable est considérée comme ayant un effet significatif sur la variable cible.

Test statistique : $\beta_j = 0$ vs $\beta_j \neq 0$ avec $j \in \{1, ..., d\}$

► Rapport de vraisemblance

$$LR_i = D_{-i} - D_{modèle} \sim \chi^2(1)$$
 asymptotiquement

Avec la déviance du modèle ajusté (avec toutes les variables explicatives)

$$D_{Mod\`{e}le} = -2 \times log(L_{Mod\`{e}le})$$

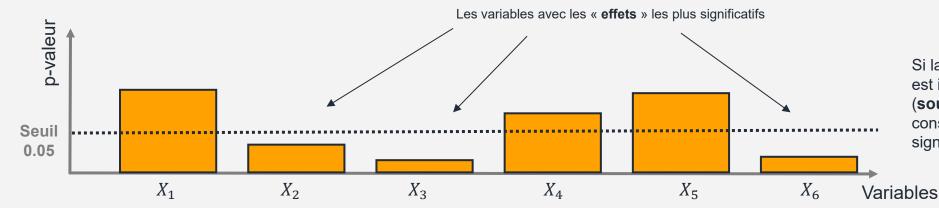
la **déviance du modèle réduit** (toutes les variables sauf X_j)

$$D_{-j} = -2 \times log(L_{-j})$$



Exemple illustratif

Test du rapport de vraisemblance



Si la **p-valeur** associée au coefficient est inférieure à un seuil conventionnel (**souvent 0.05**), la variable est considérée comme ayant un effet significatif sur la variable cible.

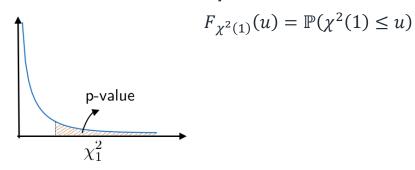
Test statistique : $\beta_j = 0$ vs $\beta_j \neq 0$ avec $j \in \{1, ..., d\}$

Rapport de vraisemblance

$$LR_j = D_{-j} - D_{mod\`{e}le} \sim \chi^2(1)$$
 asymptotiquement

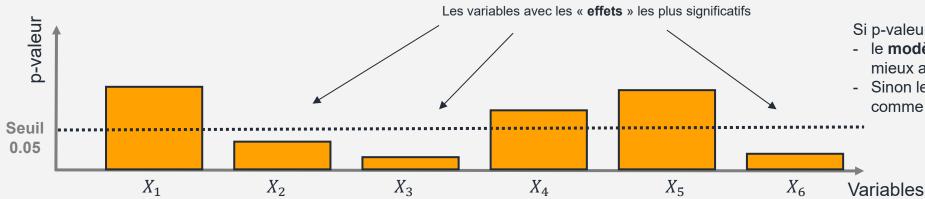
• p-valeur : $1 - F_{\chi^2(1)}(LR_j)$

Avec la fonction de répartition :





Test du rapport de vraisemblance, évaluer un groupe de prédicteurs



Exemple illustratif

Si p-valeur < 0.05 alors:

- le **modèle complet** s'ajuste significativement mieux aux données que le modèle restreint.
- Sinon le **modèle restreint** est considéré comme suffisant pour expliquer les données.

Test statistique : $\beta_j = 0$ vs $\beta_j \neq 0$ avec $j \in \{1, ..., d\}$

Rapport de vraisemblance

$$LR_j = D_{-j} - D_{mod\`{e}le} \sim \chi^2(1)$$
 asymptotiquement

p-valeur: $1 - F_{\chi^2(1)}(LR_j)$

Test statistique : $\beta_1, \beta_2, ..., \beta_r = 0$ vs $\beta_1, \beta_2, ..., \beta_r \neq 0$ avec $r \leq d$

Rapport de vraisemblance

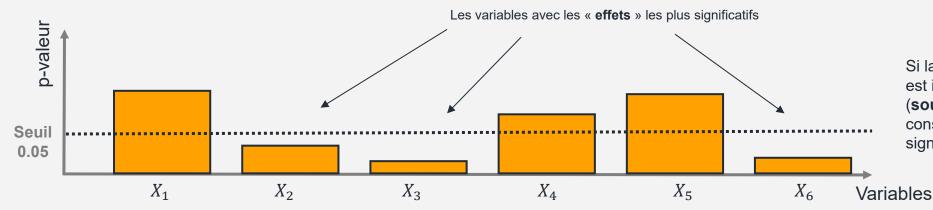
$$LR = D_{-(1,...,r)} - D_{Mod\`{e}le} \sim \chi^{2}(d-r)$$
 asymptotiquement

p-valeur : $1 - F_{\chi^2(d-r)}(LR)$



Exemple illustratif

Test de Wald : moins précis mais moins coûteux en ressources



Si la **p-valeur** associée au coefficient est inférieure à un seuil conventionnel (**souvent 0.05**), la variable est considérée comme ayant un effet significatif sur la variable cible.

Test statistique : $\beta_j = 0$ vs $\beta_j \neq 0$ avec $j \in \{1, ..., d\}$

Rapport de vraisemblance

$$LR_j = D_{-j} - D_{mod\`{e}le} \sim \chi^2(1)$$
 asymptotiquement

p-valeur: $1 - F_{\chi^2(1)}(LR_j)$

Approche moins précise mais en même temps moins coûteuse en ressources

Test de Wald, avec $\widehat{\sigma_i}$ l'erreur standard de $\widehat{\beta_i}$

$$W_j = \frac{\widehat{\beta_j}^2}{Var(\widehat{\beta_j})} \sim \chi^2(1)$$
 asymptotiquement

p-valeur : $1 - F_{\chi^{2}(1)}(W_{j})$

Critères AIC et BIC

Exemple illustratif

Des critères de sélection qui tiennent en compte le nombre de variables explicatives

Constat : L'ajout de variables dans un modèle réduit la déviance (ou augmente la vraisemblance),

$$D_{mod\`{e}le} = -2 \times log(L_{mod\`{e}le})$$

même si certaines variables sont non pertinentes, similaire à l'augmentation du R^2 en régression linéaire avec une diminution des degrés de liberté.

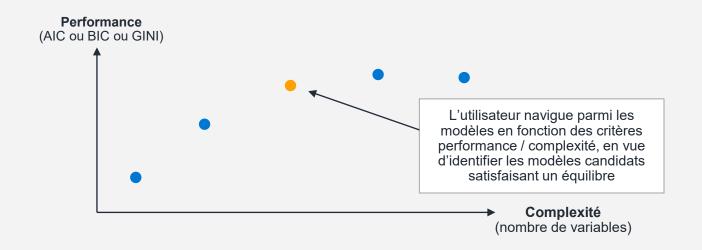
Solution : Contrebalancer la réduction de la déviance par une **pénalisation liée à la complexité du modèle**, comme avec l'AIC ou le BIC.

Critère AIKAKE:

$$AIC = -2\log(L_{mod\`{e}le}) + 2 \times (d+1)$$

Critère Schwartz:

$$BIC = -2 \log(L_{mod \ge le}) + \log(n) \times (d+1)$$



Nous allons évaluer des successions de modèles emboîtés :

- FORWARD : en les ajoutant au fur et à mesure

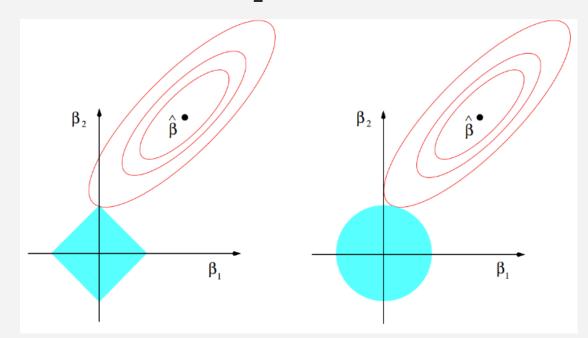
BACKWARD: en les retirant au fur et à mesure

STEPWISE: en alternant FORWARD / BACKWARD

Règle d'arrêt : l'ajout ou le retrait d'une variable n'améliore plus le critère



8. GLM pénalisé

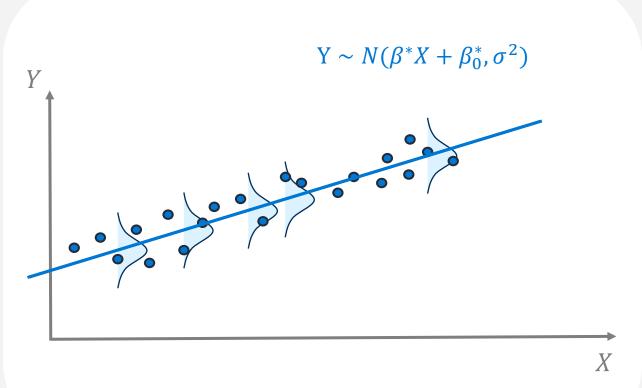




Cas des modèles linéaires classiques

Rappel : Estimateur des moindres carrés (*OLS*) : chercher les valeurs des paramètres minimisant

$$\boldsymbol{\beta}_{OLS} = \arg\min_{\boldsymbol{\beta}} R(\boldsymbol{\beta}) = \arg\min_{\boldsymbol{\beta}} \sum_{i=1,\dots,n} (Y_i - X_i \boldsymbol{\beta}_i^T)^2$$



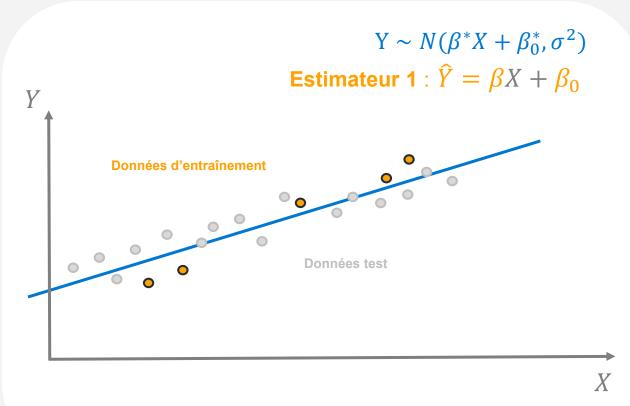


Cas des modèles linéaires classiques

Rappel : Estimateur des moindres carrés (*OLS*) : chercher les valeurs des paramètres minimisant

$$\boldsymbol{\beta}_{OLS} = \arg\min_{\boldsymbol{\beta}} R(\boldsymbol{\beta}) = \arg\min_{\boldsymbol{\beta}} \sum_{i=1,\dots,n} (Y_i - X_i \boldsymbol{\beta}_i^T)^2$$

Exemple illustratif



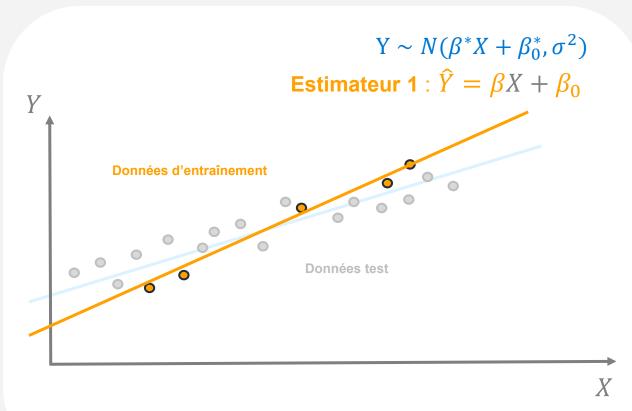


Cas des modèles linéaires classiques

Rappel : Estimateur des moindres carrés (*OLS*) : chercher les valeurs des paramètres minimisant

$$\boldsymbol{\beta}_{OLS} = \arg\min_{\boldsymbol{\beta}} R(\boldsymbol{\beta}) = \arg\min_{\boldsymbol{\beta}} \sum_{i=1,\dots,n} (Y_i - X_i \boldsymbol{\beta}_i^T)^2$$

Exemple illustratif



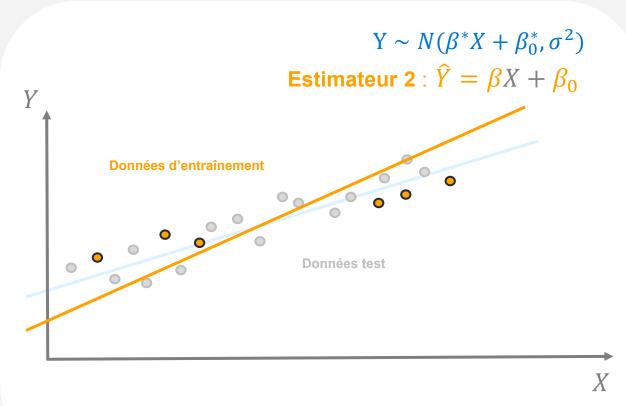


Cas des modèles linéaires classiques

Rappel : Estimateur des moindres carrés (*OLS*) : chercher les valeurs des paramètres minimisant

$$\boldsymbol{\beta}_{OLS} = \arg\min_{\boldsymbol{\beta}} R(\boldsymbol{\beta}) = \arg\min_{\boldsymbol{\beta}} \sum_{i=1,\dots,n} (Y_i - X_i \boldsymbol{\beta}_i^T)^2$$

Exemple illustratif



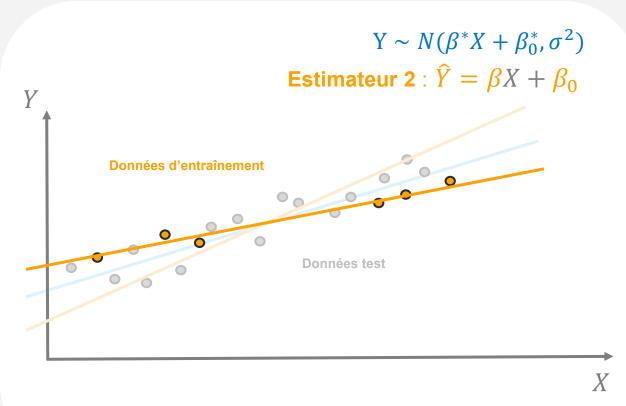


Cas des modèles linéaires classiques

Rappel : Estimateur des moindres carrés (*OLS*) : chercher les valeurs des paramètres minimisant

$$\boldsymbol{\beta}_{OLS} = \arg\min_{\boldsymbol{\beta}} R(\boldsymbol{\beta}) = \arg\min_{\boldsymbol{\beta}} \sum_{i=1,\dots,n} (Y_i - X_i \boldsymbol{\beta}_i^T)^2$$

Exemple illustratif



Problème : Incapacité à gérer le sur-apprentissage (ou surajustement)



Cas des modèles linéaires classiques

Rappel : Estimateur des moindres carrés (*OLS*) : chercher les valeurs des paramètres minimisant

$$\boldsymbol{\beta_{OLS}} = \arg\min_{\boldsymbol{\beta}} R(\boldsymbol{\beta}) = \arg\min_{\boldsymbol{\beta}} \sum_{i=1,\dots,n} (Y_i - X_i \boldsymbol{\beta_i^T})^2$$

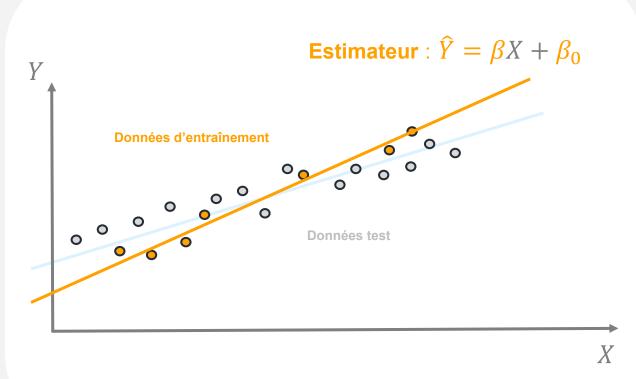
Modèle linéaire pénalisé (Ridge ou Lasso) : modèle linéaire qui ajoute une pénalité sur la magnitude des coefficients pour éviter le surajustement.

$$\beta_{ridge} = \arg\min_{\beta} R(\beta) + \lambda ||\beta||^{2}$$

$$= \arg\min_{\beta} R(\beta) + \lambda \sum_{i=1,\dots,n} \beta_{i}^{2}$$

$$\beta_{lasso} = \arg\min_{\beta} R(\beta) + \lambda ||\beta||^{1}$$

$$= \arg\min_{\beta} R(\beta) + \lambda \sum_{i=1,\dots,n} |\beta_{i}|$$



Problème : Incapacité à gérer le sur-apprentissage (ou surajustement)



Cas des modèles linéaires classiques

Rappel : Estimateur des moindres carrés (*OLS*) : chercher les valeurs des paramètres minimisant

$$\boldsymbol{\beta_{OLS}} = \arg\min_{\boldsymbol{\beta}} R(\boldsymbol{\beta}) = \arg\min_{\boldsymbol{\beta}} \sum_{i=1,\dots,n} (Y_i - X_i \boldsymbol{\beta_i^T})^2$$

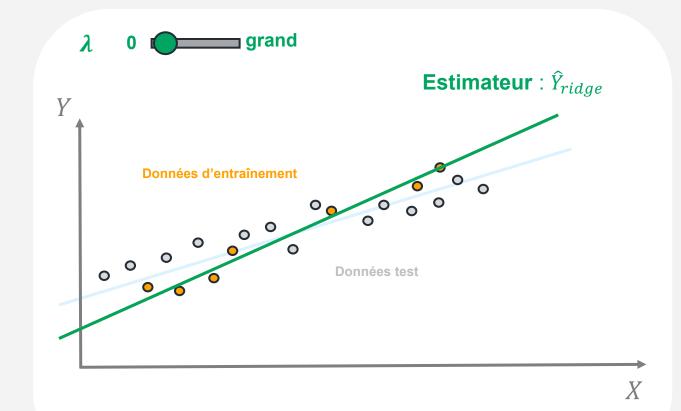
Modèle linéaire pénalisé (Ridge ou Lasso) : modèle linéaire qui ajoute une pénalité sur la magnitude des coefficients pour éviter le surajustement.

$$\beta_{ridge} = \arg\min_{\beta} R(\beta) + \lambda ||\beta||^{2}$$

$$= \arg\min_{\beta} R(\beta) + \lambda \sum_{i=1,\dots,n} \beta_{i}^{2}$$

$$\beta_{lasso} = \arg\min_{\beta} R(\beta) + \lambda ||\beta||^{1}$$

$$= \arg\min_{\beta} R(\beta) + \lambda \sum_{i=1,\dots,n} |\beta_{i}|$$





Cas des modèles linéaires classiques

Rappel : Estimateur des moindres carrés (*OLS*) : chercher les valeurs des paramètres minimisant

$$\beta_{OLS} = \arg\min_{\beta} R(\beta) = \arg\min_{\beta} \sum_{i=1,\dots,n} (Y_i - X_i \beta_i^T)^2$$

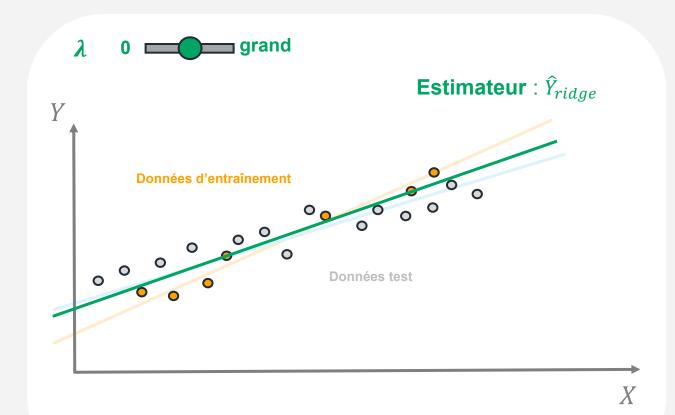
Modèle linéaire pénalisé (Ridge ou Lasso) : modèle linéaire qui ajoute une pénalité sur la magnitude des coefficients pour éviter le surajustement.

$$\beta_{ridge} = \arg\min_{\beta} R(\beta) + \lambda || \beta ||^{2}$$

$$= \arg\min_{\beta} R(\beta) + \lambda \sum_{i=1,\dots,n} \beta_{i}^{2}$$

$$\beta_{lasso} = \arg\min_{\beta} R(\beta) + \lambda || \beta ||^{1}$$

$$= \arg\min_{\beta} R(\beta) + \lambda \sum_{i=1,\dots,n} |\beta_{i}|$$





Cas des modèles linéaires classiques

Rappel : Estimateur des moindres carrés (*OLS*) : chercher les valeurs des paramètres minimisant

$$\boldsymbol{\beta_{OLS}} = \arg\min_{\boldsymbol{\beta}} R(\boldsymbol{\beta}) = \arg\min_{\boldsymbol{\beta}} \sum_{i=1,\dots,n} (Y_i - X_i \boldsymbol{\beta_i^T})^2$$

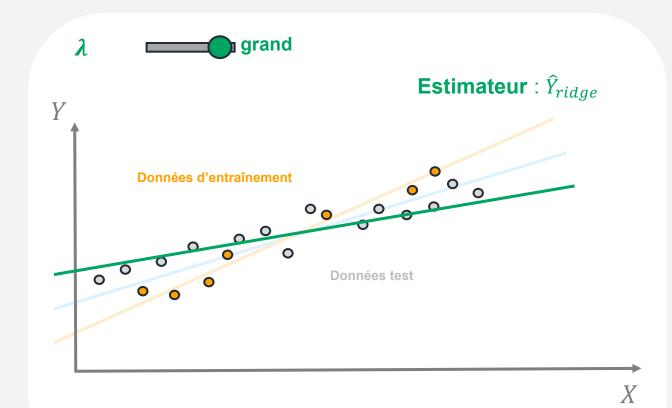
Modèle linéaire pénalisé (Ridge ou Lasso) : modèle linéaire qui ajoute une pénalité sur la magnitude des coefficients pour éviter le surajustement.

$$\beta_{ridge} = \arg\min_{\beta} R(\beta) + \lambda ||\beta||^{2}$$

$$= \arg\min_{\beta} R(\beta) + \lambda \sum_{i=1,\dots,n} \beta_{i}^{2}$$

$$\beta_{lasso} = \arg\min_{\beta} R(\beta) + \lambda ||\beta||^{1}$$

$$= \arg\min_{\beta} R(\beta) + \lambda \sum_{i=1,\dots,n} |\beta_{i}|$$



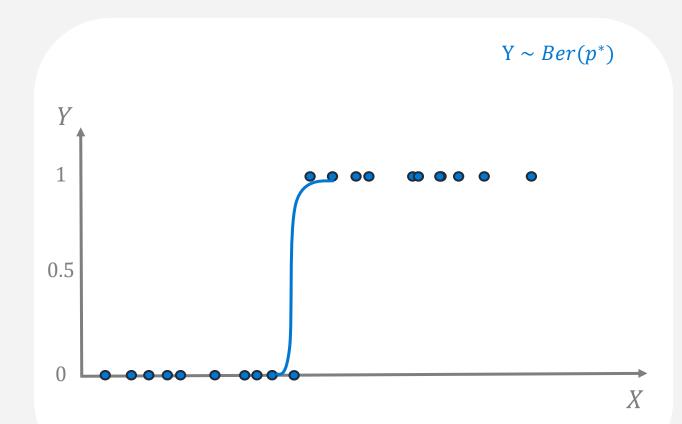
Nous sélectionnons la pénalisation optimale avec la validation croisée!



Cas des modèles linéaires généralisés

Rappel: Estimateur du maximum de vraisemblance (*MLE*): chercher les valeurs des paramètres maximisant la fonction de vraisemblance (ou log de vraisemblance)

$$\beta_{MLE} = \arg \max_{\beta} l(y, \beta) = \arg \max_{\beta} \sum_{i=1,...,n} \log f_{\beta}(y_i)$$

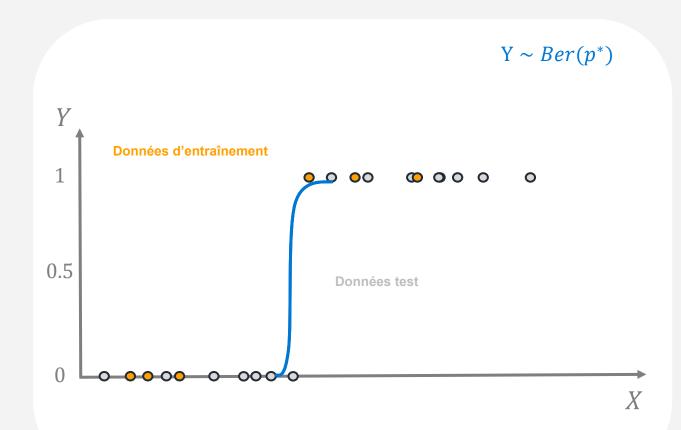




Cas des modèles linéaires généralisés

Rappel: Estimateur du maximum de vraisemblance (*MLE*): chercher les valeurs des paramètres maximisant la fonction de vraisemblance (ou log de vraisemblance)

$$\beta_{MLE} = \arg \max_{\beta} l(y, \beta) = \arg \max_{\beta} \sum_{i=1,...,n} \log f_{\beta}(y_i)$$

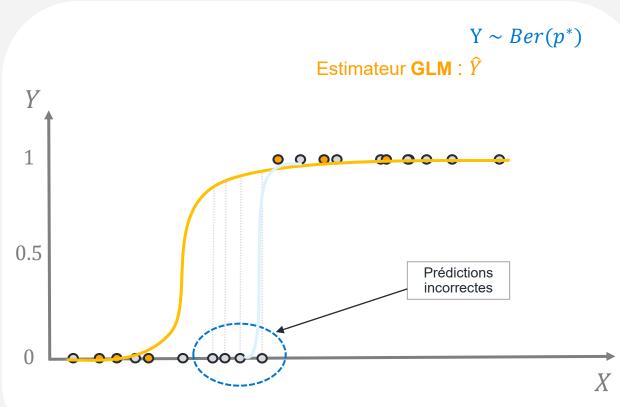




Rappel: Estimateur du maximum de vraisemblance (*MLE*): chercher les valeurs des paramètres maximisant la fonction de vraisemblance (ou log de vraisemblance)

$$\beta_{MLE} = \arg \max_{\beta} l(y, \beta) = \arg \max_{\beta} \sum_{i=1,...,n} \log f_{\beta}(y_i)$$

Exemple illustratif



Problème : Gérer le sur-apprentissage



Rappel: Estimateur du maximum de vraisemblance (*MLE*): chercher les valeurs des paramètres maximisant la fonction de vraisemblance (ou log de vraisemblance)

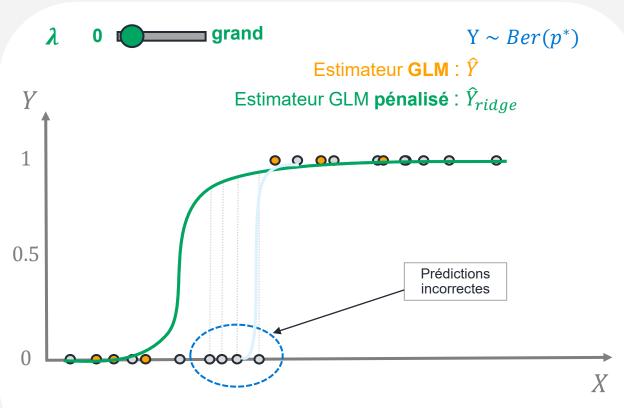
$$\beta_{MLE} = \arg \max_{\beta} l(y, \beta) = \arg \max_{\beta} \sum_{i=1,...,n} \log f_{\beta}(y_i)$$

GLM avec une pénalité (Ridge ou Lasso) : GLM qui ajoute une pénalité sur la magnitude des coefficients pour éviter le surajustement.

$$\beta_{MLE-ridge} = \arg \max_{\beta} l(y, \beta) - \lambda || \beta ||^2$$

$$\beta_{MLE-lasso} = \arg \max_{\beta} l(y, \beta) - \lambda || \beta ||^{1}$$

Exemple illustratif



Problème : Gérer le sur-apprentissage



Rappel: Estimateur du maximum de vraisemblance (*MLE*): chercher les valeurs des paramètres maximisant la fonction de vraisemblance (ou log de vraisemblance)

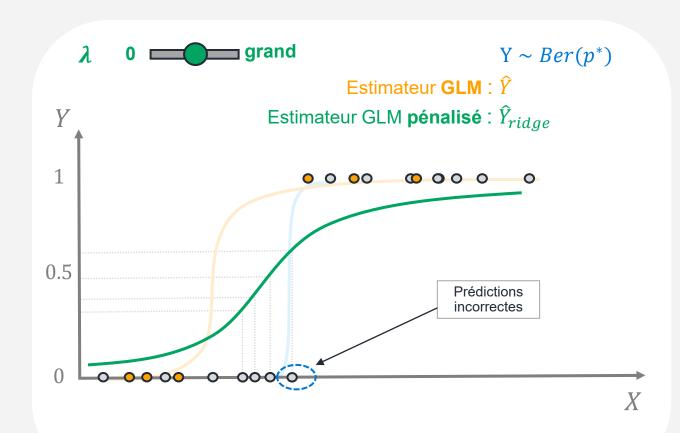
$$\beta_{MLE} = \arg \max_{\beta} l(y, \beta) = \arg \max_{\beta} \sum_{i=1,...,n} \log f_{\beta}(y_i)$$

► GLM avec une pénalité (Ridge ou Lasso) : GLM qui ajoute une pénalité sur la magnitude des coefficients pour éviter le surajustement.

$$\beta_{MLE-ridge} = \arg \max_{\beta} l(y, \beta) - \lambda || \beta ||^2$$

$$\beta_{MLE-lasso} = \arg \max_{\beta} l(y, \beta) - \lambda || \beta ||^{1}$$

Exemple illustratif





$$\beta_{MLE} = \arg \max_{\beta} l(y, \beta) = \arg \max_{\beta} \sum_{i=1,...,n} \log f_{\beta}(y_i)$$

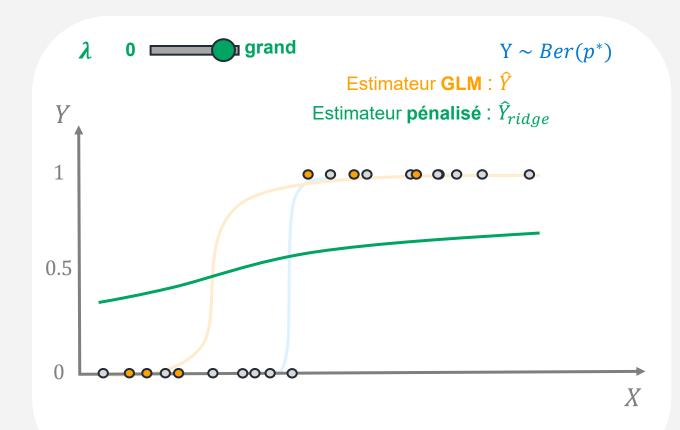
GLM avec une pénalité (Ridge ou Lasso) : GLM qui ajoute une pénalité sur la magnitude des coefficients pour éviter le surajustement.

$$\beta_{MLE-ridge} = \arg \max_{\beta} l(y, \beta) - \lambda ||\beta||^2$$

$$\beta_{MLE-lasso} = \arg \max_{\beta} l(y, \beta) - \lambda || \beta ||^{1}$$

Nous sélectionnons la pénalisation optimale avec la validation croisée!

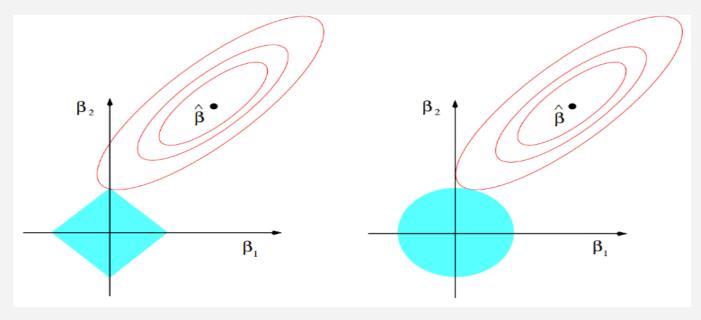
Exemple illustratif





Différences entre les modèles Ridge et Lasso (1/3)

Eviter le sur-apprentissage durant l'apprentissage



Lasso: La pénalité Lasso utilise la norme L1, poussant certains coefficients à zéro, ce qui permet une sélection automatique des variables. Il est utile pour simplifier les modèles en réduisant la dimensionnalité et en améliorant la généralisation. VS

Ridge: La pénalité Ridge utilise la norme L2 pour réduire la taille des coefficients, aidant à améliorer la généralisation, surtout en présence de variables corrélées.



Différences entre les modèles Ridge et Lasso (2/3)

Tableau récapitulatif

	Ridge	Lasso
Type pénalité	Norme L2 (somme des carrés des coefficients)	Norme L1 (somme des valeurs absolues des coefficients)
Effet sur les coefficients	Réduit les coefficients mais ne les annule pas	Peut amener certains coefficients à zéro
Sélection de variables		✓
Traiter la multi-colinéarité	✓	Moins efficace que Ridge
Stabilité des coefficients	✓	Moins stable que Ridge, (des coefficients sont nuls)
Utilisation optimale	Idéal lorsque les variables sont corrélées et toutes utiles	Idéal lorsqu'il y a beaucoup de variables , dont certaines non pertinentes
Quand utiliser ?	Réduire la variance sans éliminer de variables	Simplifier le modèle et à réduire la dimensionnalité



Différences entre les modèles Ridge et Lasso (3/3)

Extension avec la méthode Elastic-Net

Elastic Net combine les pénalités Ridge (L2) et Lasso (L1) :

$$\arg\max_{\pmb{\beta}} l(y,\pmb{\beta}) - \pmb{\lambda_2} \mid\mid \pmb{\beta} \mid\mid^2 - \pmb{\lambda_1} \mid\mid \pmb{\beta} \mid\mid^1$$

- Régularise les coefficients tout en effectuant une sélection de variables
- Particulièrement utile lorsque les variables explicatives sont corrélées et qu'on souhaite à la fois réduire la dimensionnalité et éviter la multi-colinéarité.

	Ridge	Lasso
Type pénalité	Norme L2 (somme des carrés des coefficients)	Norme L1 (somme des valeurs absolues des coefficients)
Effet sur les coefficients	Réduit les coefficients mais ne les annule pas	Peut amener certains coefficients à zéro
Sélection de variables		✓
Traiter la multi-colinéarité	✓	Moins efficace que Ridge
Stabilité des coefficients	✓	Moins stable que Ridge, (des coefficients sont nuls)
Utilisation optimale	Idéal lorsque les variables sont corrélées et toutes utiles	Idéal lorsqu'il y a beaucoup de variables , dont certaines non pertinentes
Quand utiliser ?	Réduire la variance sans éliminer de variables	Simplifier le modèle et à réduire la dimensionnalité



9. Extension aux GAM

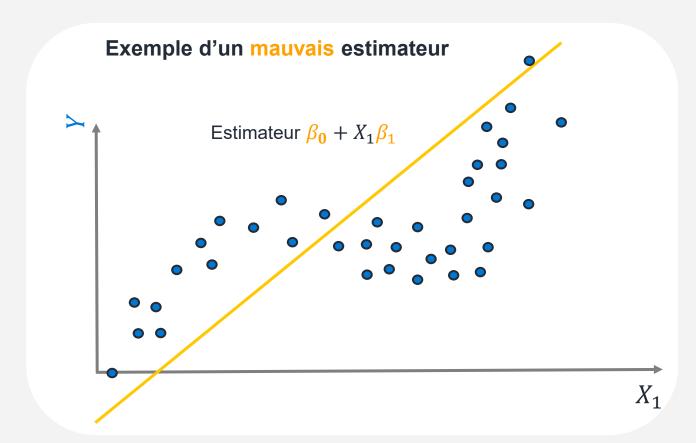


Limites aux modèles linéaires

Cas relation non-linéaire entre variables explicatives et variable cible

Modèle Linéaire Généralisé (GLM)

$$g(\mathbb{E}[Y \mid X]) = \frac{\beta_0}{\beta_0} + \frac{\beta_1}{\beta_1} X_1 + \frac{\beta_2}{\beta_2} X_2 + \dots + \frac{\beta_d}{\beta_d} X_d$$



Limites aux modèles linéaires

Cas relation non-linéaire entre variables explicatives et variable cible

Modèle Linéaire Généralisé (GLM)

$$g(\mathbb{E}[Y \mid X]) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_d X_d$$

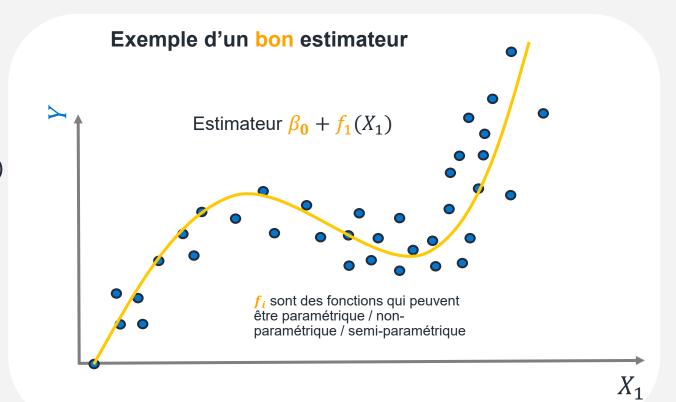
Modèle Additif Généralisé (GAM)

$$g(\mathbb{E}[Y \mid X]) = \beta_0 + f_1(X_1) + f_2(X_2) + \dots + f_d(X_d)$$

Les **GAM** (Modèles Additifs Généralisés) sont **une extension des GLM** qui permettent une plus grande flexibilité dans la modélisation des relations entre X et Y. Contrairement aux modèles GLM qui sont linéaires entre $g(\mathbb{E}[Y \mid X])$ et X, les GAM permettent à chaque variable explicative d'avoir une **relation non linéaire avec la réponse**.

Exemple de fonction :

$$f(x) = \beta_1 x + \beta_2 x^2 + \dots + \beta_r x^r$$









Merci

François HU

Francois.hu@milliman.com