

# À la découverte des prédictions conformes

(et des statistiques)

# Contexte ? La consommation électrique de l'IF



Taux d'occupation



Température extérieure



Heure de la journée

...



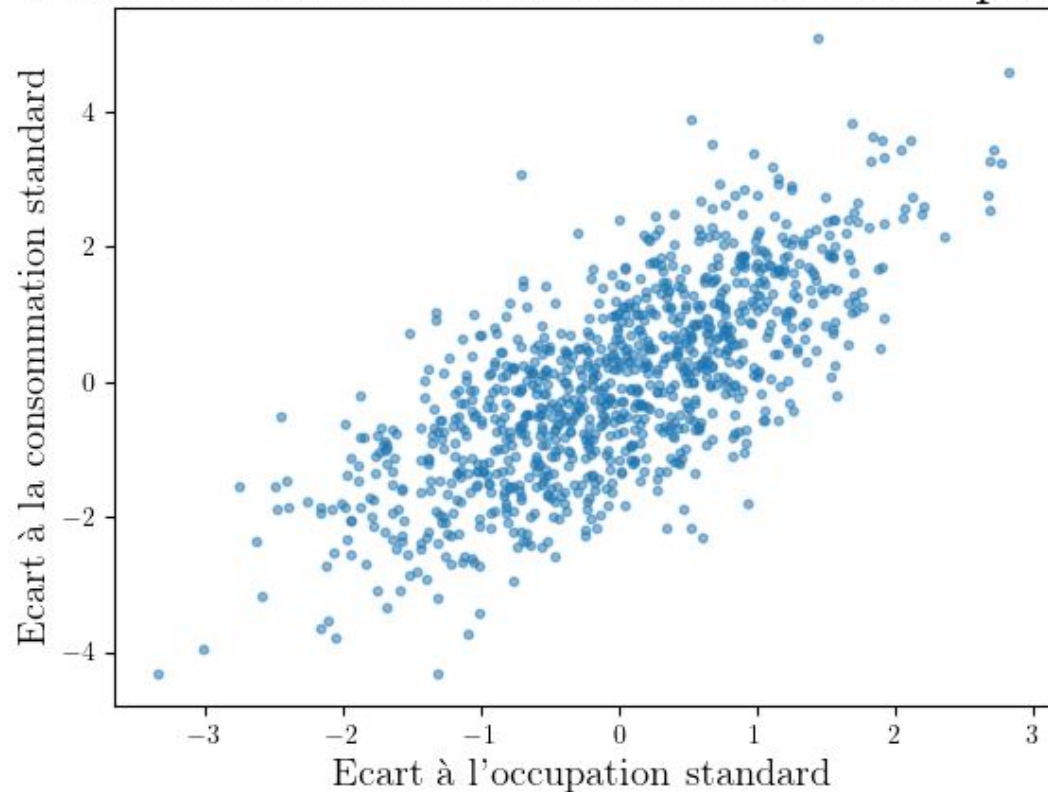
Consommation électrique

**Variables explicatives** 👁️



**Variable expliquée** 🎯

## Consommation en fonction de l'occupation

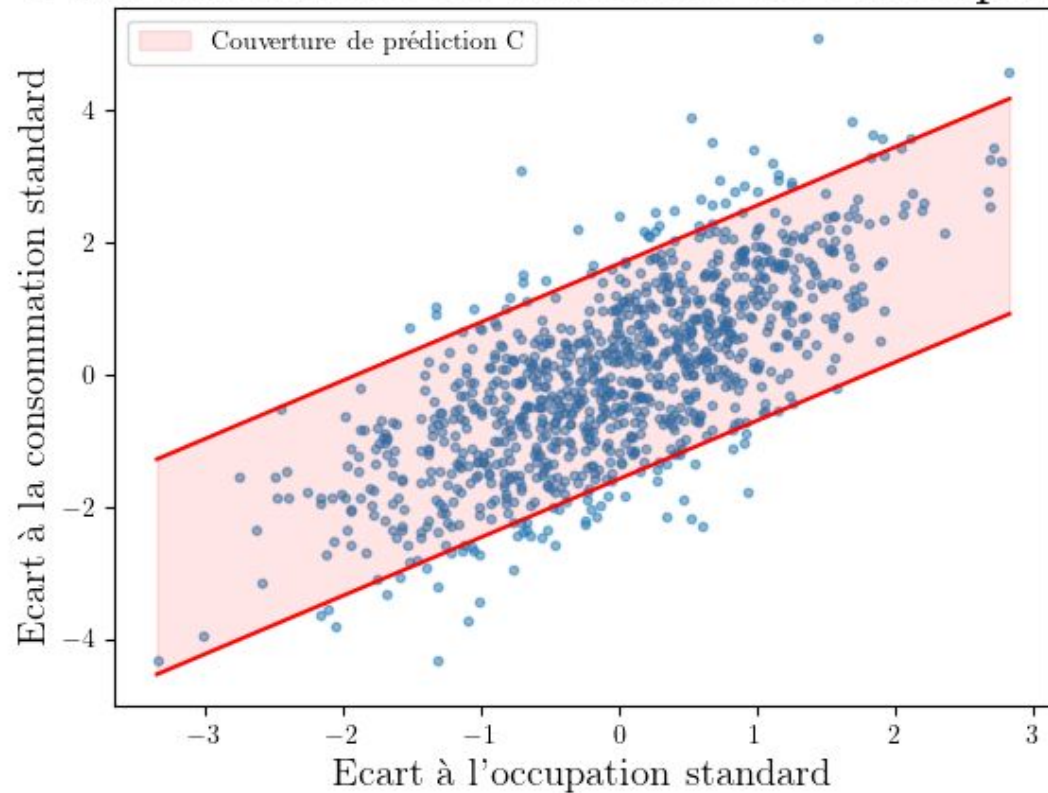


$$X_i \sim \mathcal{N}(0, 1)$$

$$\epsilon_i \sim \mathcal{N}(0, 1)$$

$$Y_i = X_i + \epsilon_i$$

# Consommation en fonction de l'occupation



$$\alpha \in ]0, 1[$$

$$\mathbb{P}[Y \in C(x) | X = x] \geq 1 - \alpha$$

Nouvelle observation

**Définition** (Statistique d'ordre). Soient  $X_1, \dots, X_n$  des variables aléatoires réelles. On appelle statistiques d'ordre les variables aléatoires  $X_{(1)}, \dots, X_{(n)}$  obtenues en réarrangeant les  $X_i$  par ordre croissant :

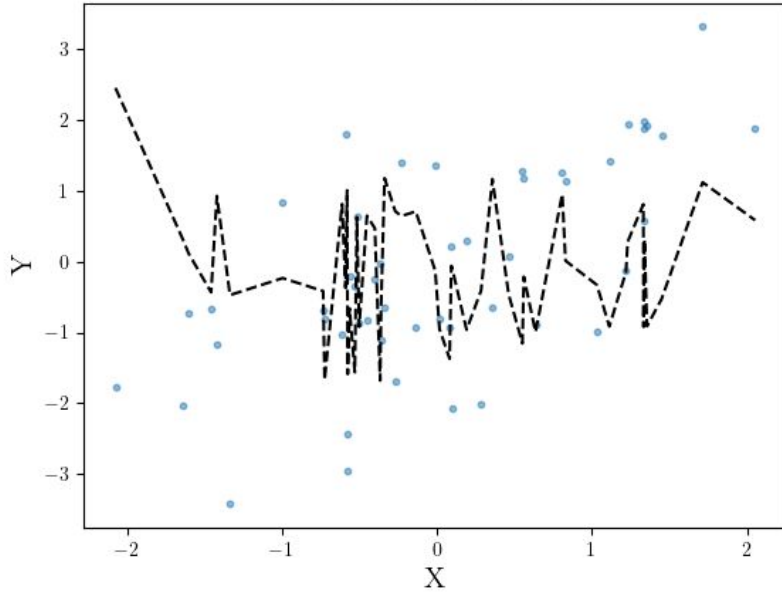
$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$$

Ainsi, pour tout  $k \in \llbracket 1; n \rrbracket$ ,  $X_{(k)}$  représente la  $k$ -ième plus petite valeur parmi les  $X_i$ .

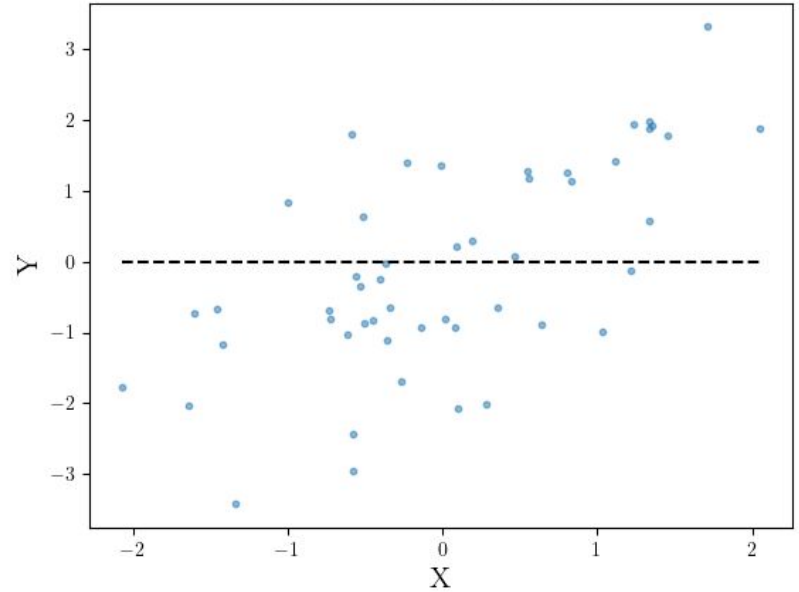
**Proposition 1.** Soient  $U_1, \dots, U_{n+1}$  des variables aléatoires réelles échangeables. On note  $F$  leur fonction de répartition commune qu'on suppose continue (pas de cas d'égalité) et  $f$  leur densité. On pose  $U_{(1)}, \dots, U_{(n+1)}$  leur statistique d'ordre. Le rang de  $U_{n+1}$  parmi  $U_1, \dots, U_{n+1}$  suit la loi uniforme sur  $\llbracket 1; n+1 \rrbracket$

# La problématique

Sur-entraînement



Sous-entraînement



$\hat{\mu}$  la fonction de régression empirique

## Données i.i.d.

Les données observées  $Z_1, \dots, Z_d$  sont i.i.d à valeur dans  $\mathbb{R}^d \times \mathbb{R}$  issues d'une distribution  $P$  avec une fonction moyenne  $\mu : x \mapsto \mathbb{E}[Y|X = x]$ .

## Bruit symétrique et indépendant)

Le bruit est symétrique et indépendant. C'est-à-dire que pour,  $(X, Y) \sim P$  on a  $\epsilon = Y - \mu(X)$  indépendant de  $X$ . De plus, la fonction densité de  $\epsilon$  est symétrique et décroissante sur  $[0, \infty[$ .

## Stabilité de l'échantillonnage

Soit  $\hat{\mu}_n$  la fonction de régression approchée entraînée sur  $n$  observations. Il existe  $(\eta_n)$  et  $(\rho_n)$  deux suites réelles convergeant vers 0 et une fonction  $\tilde{\mu} : \mathbb{R}^d \rightarrow \mathbb{R}$  telles qu'à partir d'un certain rang  $n$  :

$$\mathbb{P}[\|\hat{\mu}_n - \tilde{\mu}\|_\infty \geq \eta_n] \leq \rho_n$$

## Sensibilité à la perturbation

Pour  $\mathcal{Y}$  un segment, en notant pour  $y \in \mathcal{Y}$ ,  $\hat{\mu}_{n,(X,y)}$  la fonction de régression entraînée sur  $n$  observation agrémentée du point  $(X, y)$ . Il existe  $(\eta_n)$  et  $(\rho_n)$  deux suites réelles convergeant vers 0 telles que

$$\mathbb{P}[\sup_{y \in \mathcal{Y}} \|\hat{\mu}_n - \hat{\mu}_{n,(X,y)}\|_\infty \geq \eta_n] \leq \rho_n$$

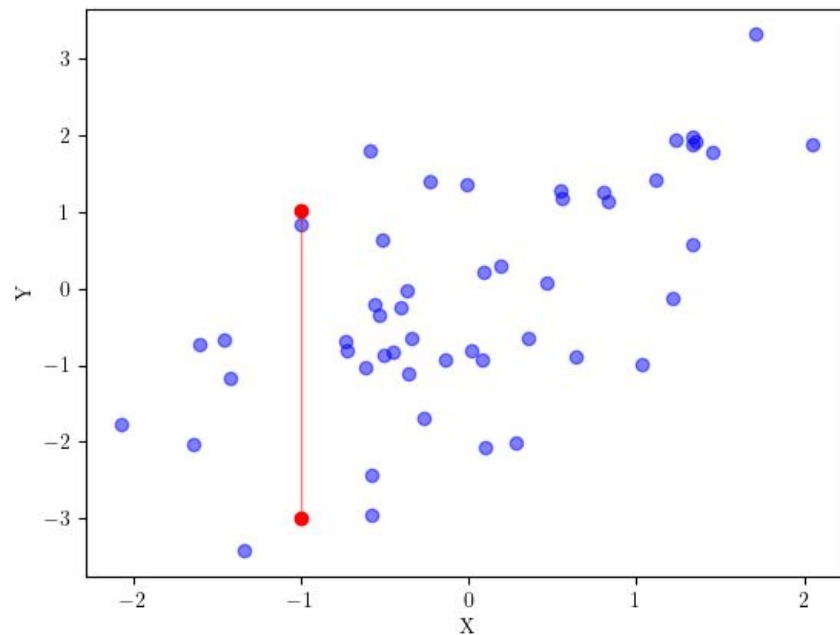
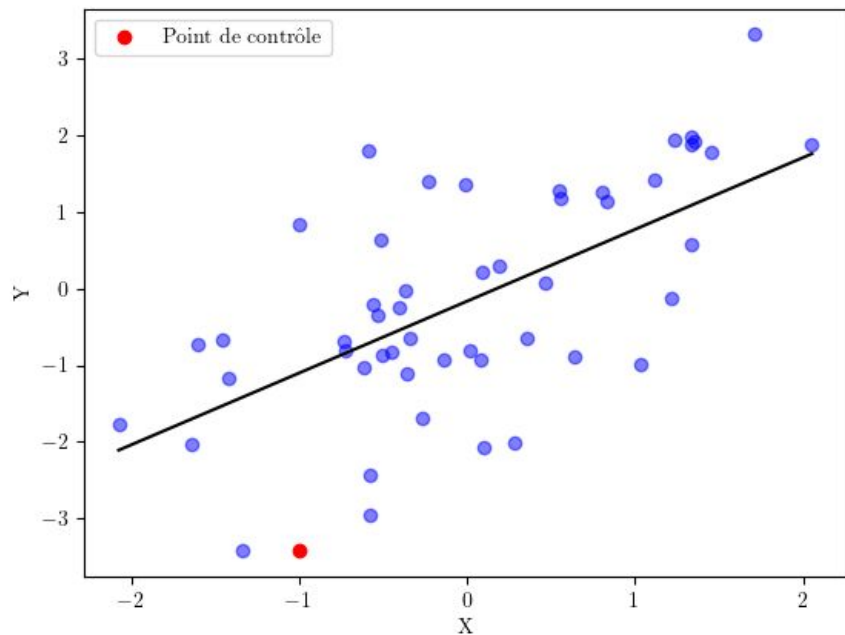
## Cohérence de l'estimateur

Il existe  $(\eta_n)$  et  $(\rho_n)$  deux suites réelles convergentes vers 0 telles qu'à partir d'un certain rang

$$\mathbb{P}[\mathbb{E}[(\hat{\mu}_n(X) - \mu(X))^2 | \hat{\mu}_n] \geq \eta_n] \leq \rho_n$$

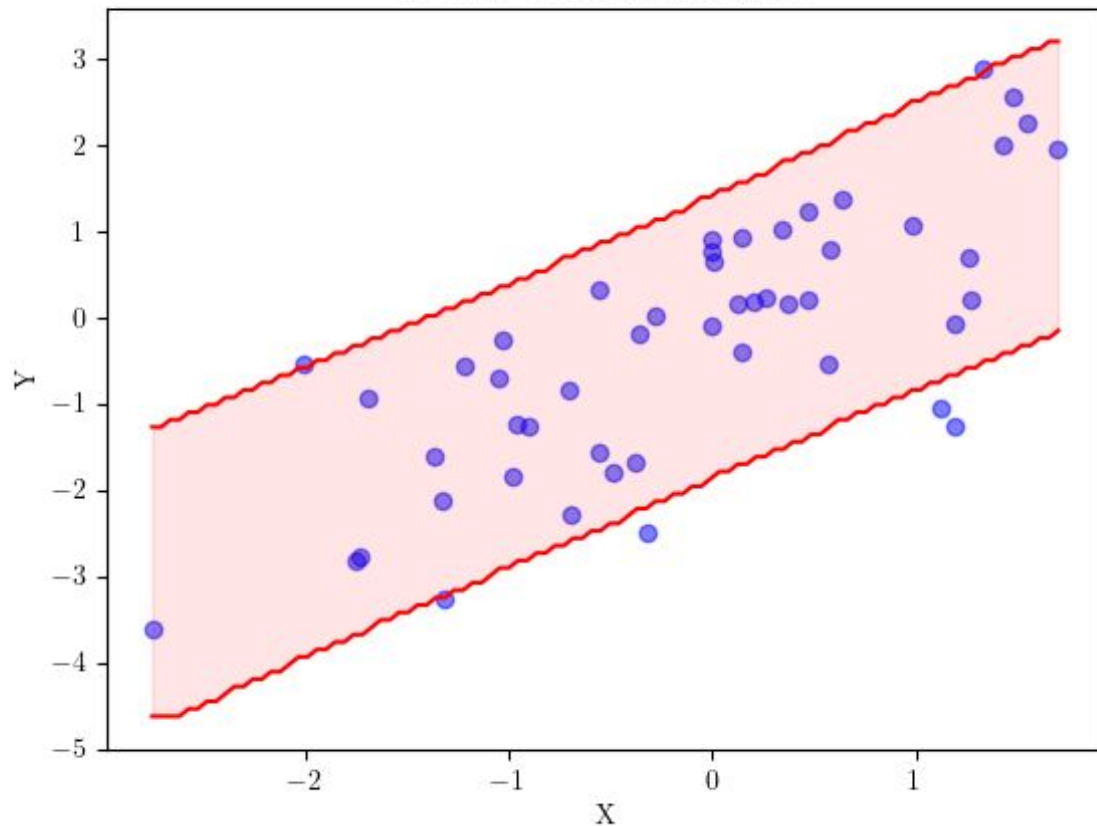


# Full conformal



$R_{i,y} := |Y_i - \hat{\mu}_y(X_i)|$ , le résidu absolu ajusté

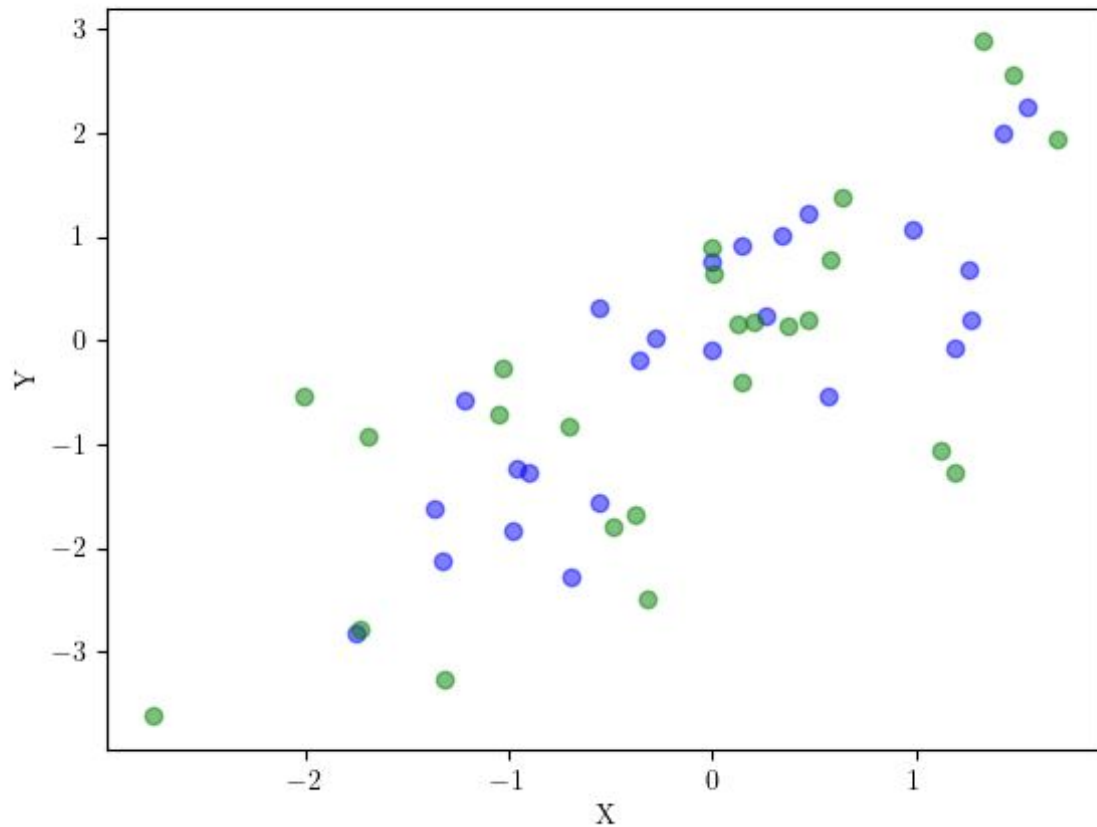
## Full conformal



$$\#Y_{trial} = 100$$

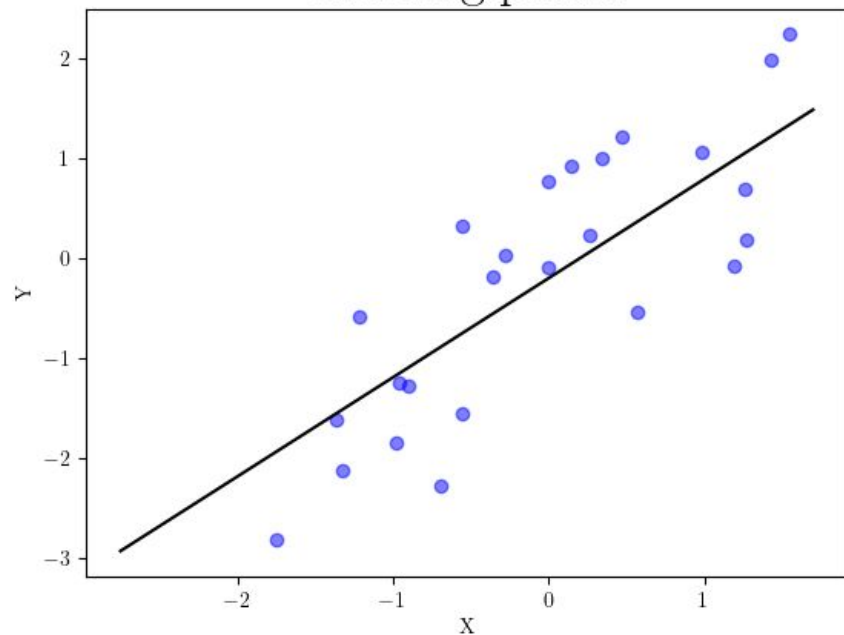
$$\#X_{conf} = 100$$

# Split Conformal

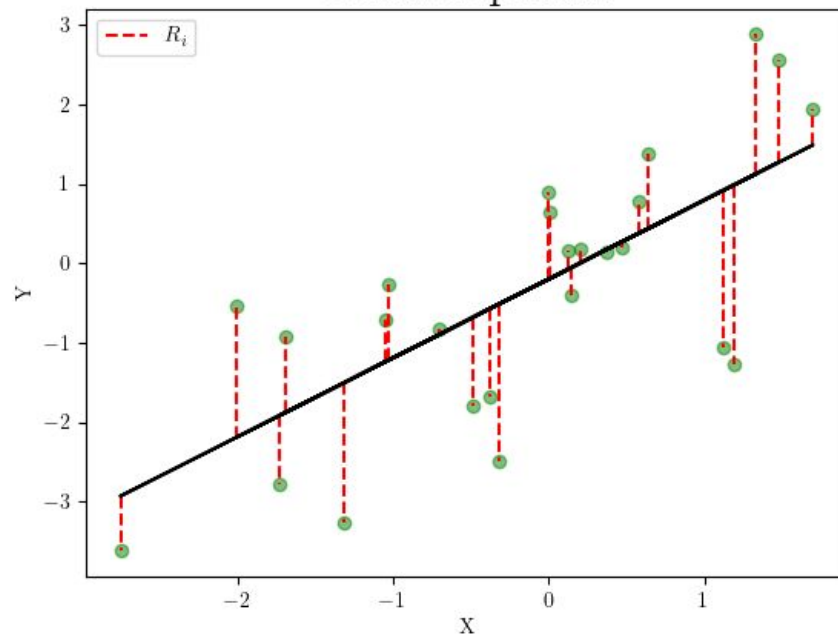


# Split Conformal

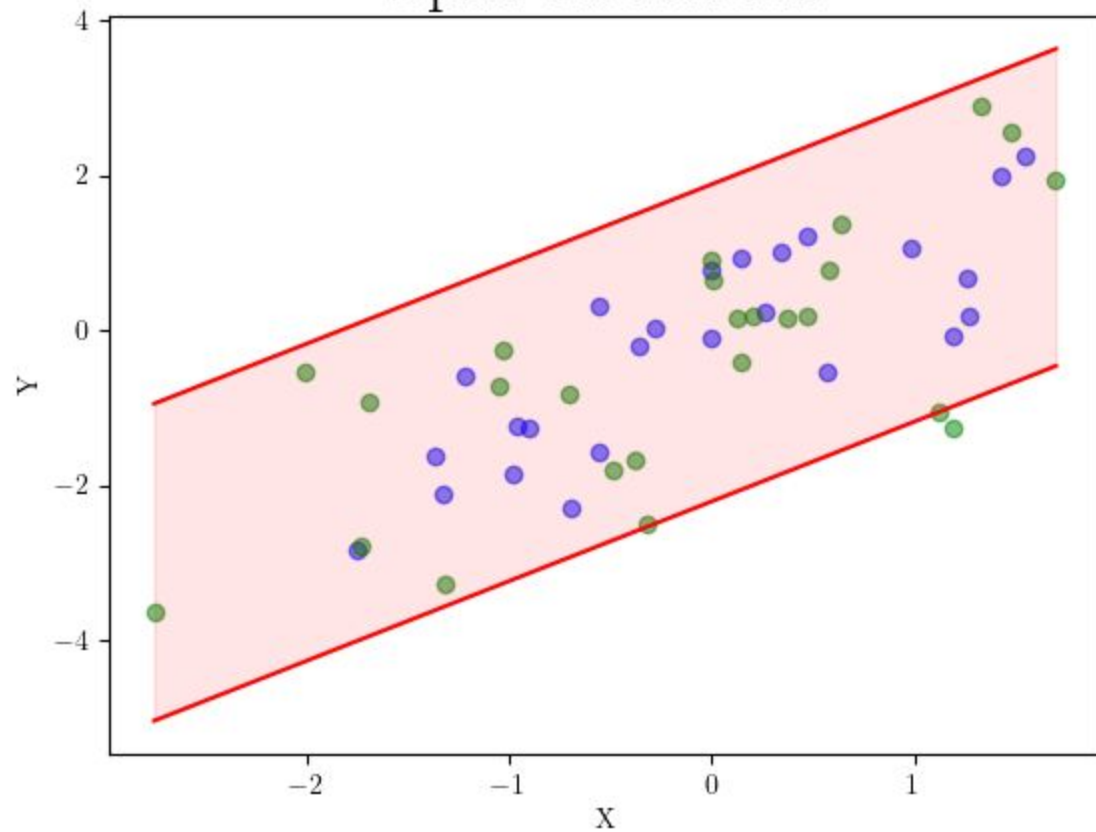
Training points

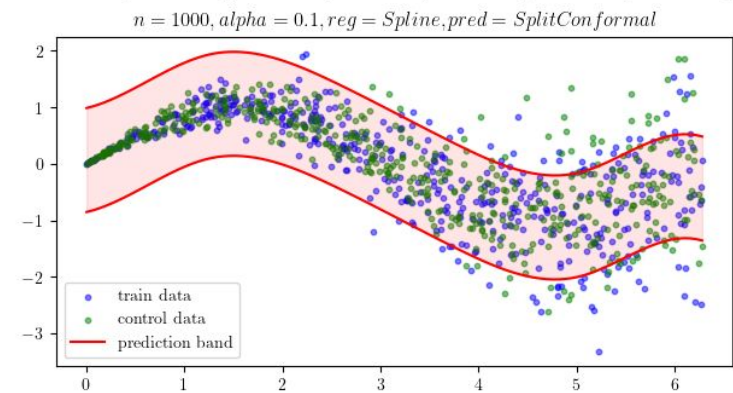
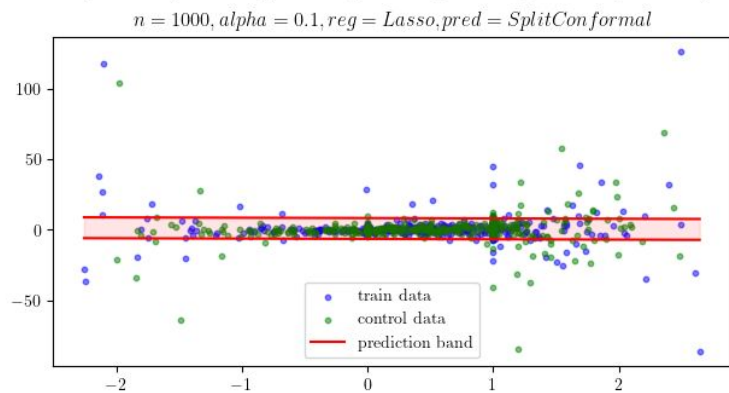
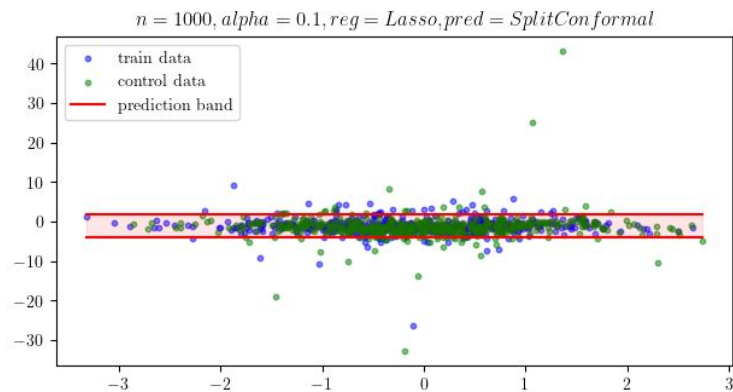
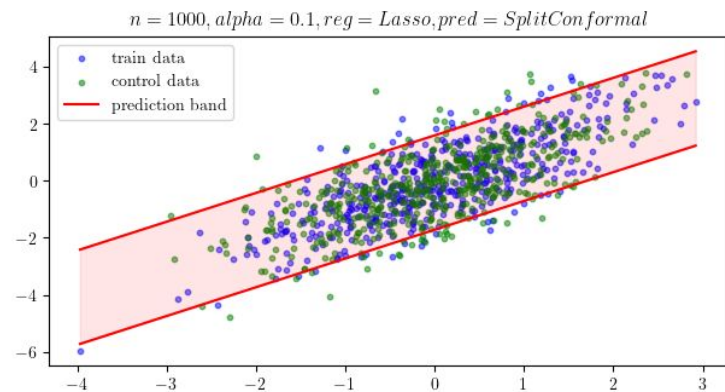


Control points



## Split conformal





	Setting A	Setting B	Setting C	Setting Partic 5
Coverage	89.89	89.90	90.02	89.93

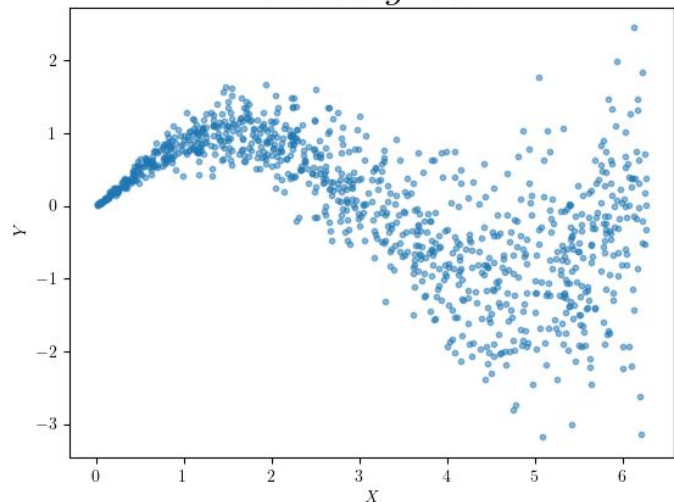
**Définition** (Hétéroscédasticité). Soit un modèle de régression

$$Y_i = \mu(X_i) + \varepsilon_i, \quad i = 1, \dots, n,$$

où les erreurs  $\varepsilon_i$  vérifient  $\mathbb{E}[\varepsilon_i] = 0$ . On dit que le modèle présente de l'hétéroscédasticité s'il existe une fonction non constante  $f$  telle que

$$\text{Var}(\varepsilon_i) = f(X_i) \quad \text{pour tout } i.$$

Setting P5

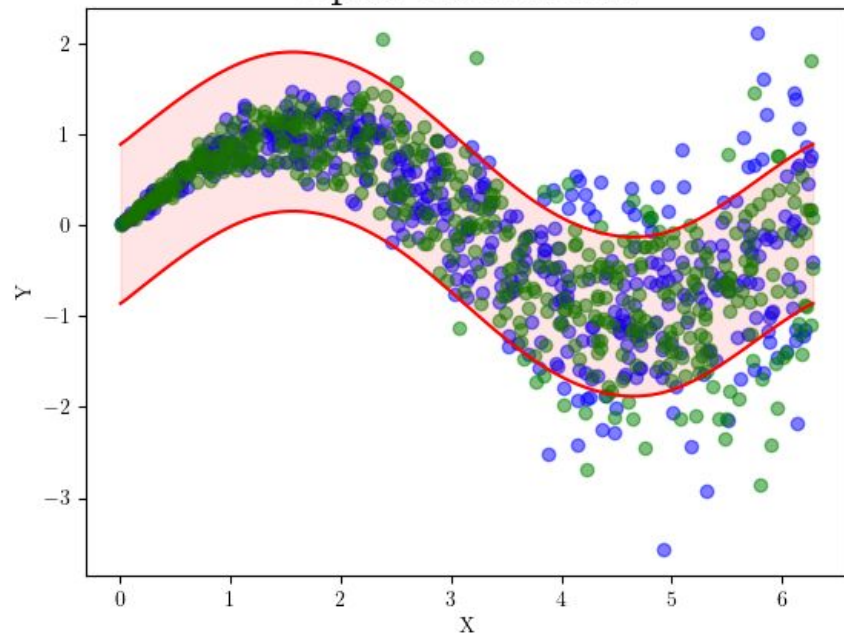


Exemple :

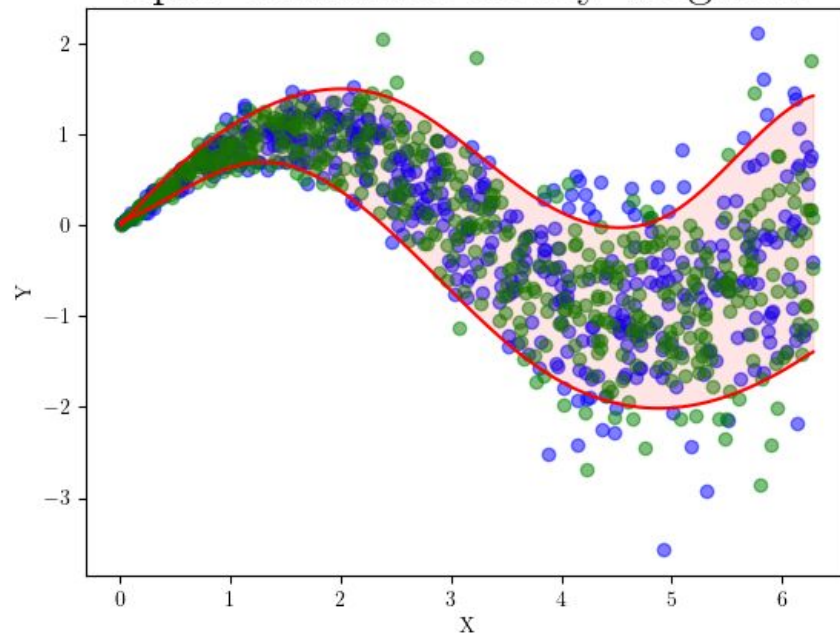
$$X_i \sim \mathcal{U}(0, 2\pi) \quad \varepsilon_i \sim \mathcal{N}(0, 1)$$

$$Y_i = \sin(X_i) + \frac{\pi|X_i|}{20}\varepsilon_i$$

Split conformal



Split conformal locally-weighted



Hétéroscédasticité



