

A U-Turn on "A U-Turn on Double Descent"

①

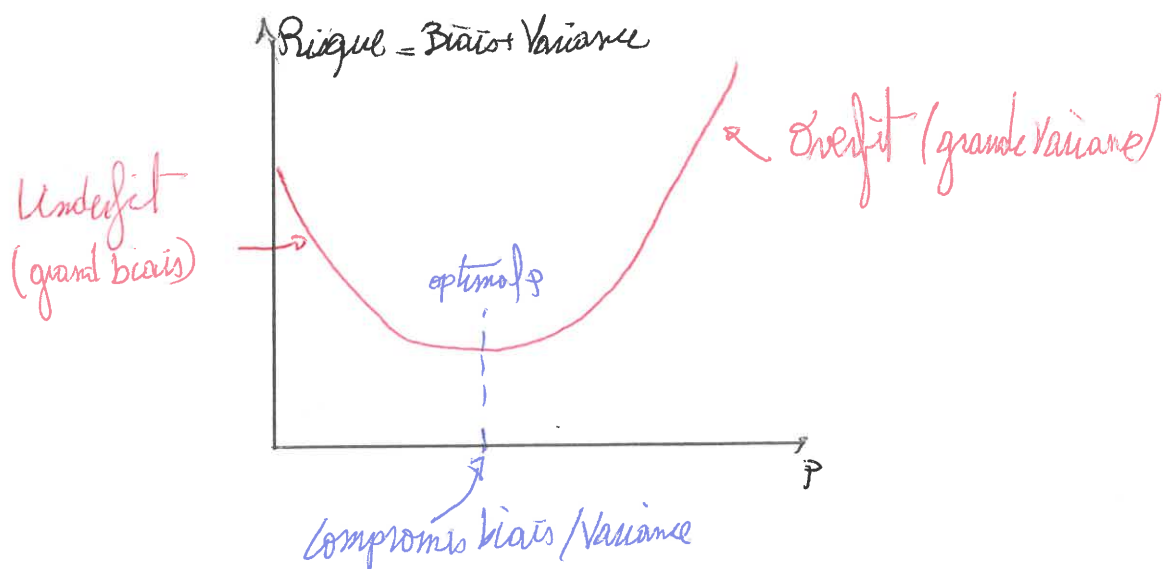
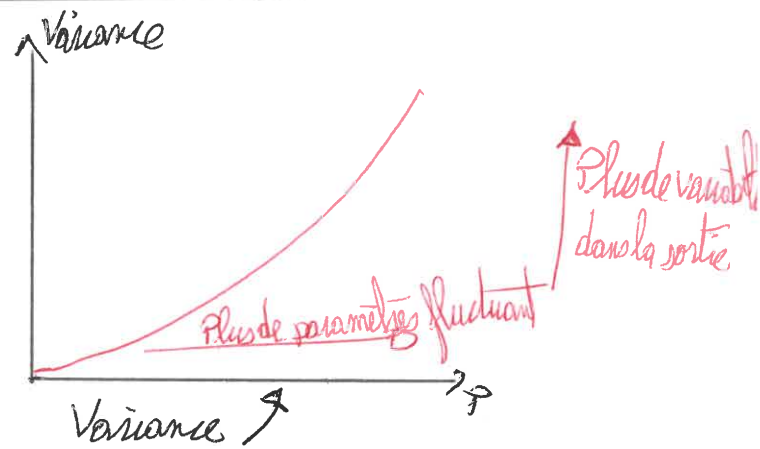
"A U-Turn on Double Descent: Rethinking Parameter Counting in Statistical Learning"
[Curth, Jeffares, v.d. Schaar - 2023]

I Biais, Variance, Double descent

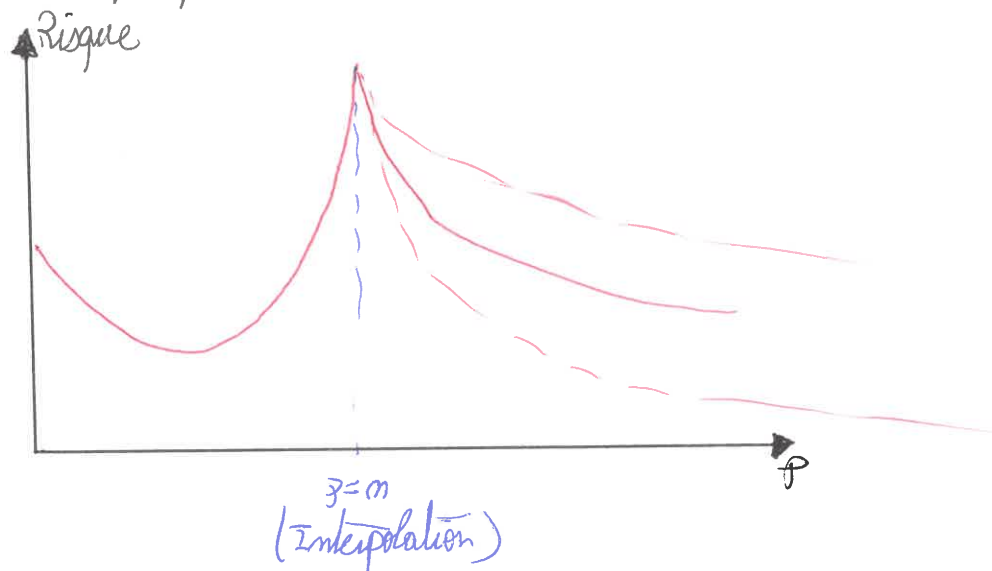
1) Phénomène de double descente

En théorie classique d'estimation, $\text{Risque} = \text{Biais} + \text{Variance}$

On note $p := \{\text{\# paramètres dans le modèle}\}$



Pourtant, des observations empiriques (notamment sur des réseaux de neurones) donnent (2)



→ Phénomène de "Double Descente" (Belkin et al 2019)
(A Brief History of Double Descent)

Programme :

- Exposer un résultat d'obligation au compromis biais / variance
- Donner des pistes sur une notion de dimension effective p_{eff} en régression linéaire featurisée (Regression RKHS)
- Tentative de faire l'exégèse du papier.

2) Le compromis biais / Variance est incontournable

"On lower bounds for the bias-variance trade-off"

[Derumigny, Schmitt-Hieber, 2023]

Dans des régimes surparamétrés, biais et variance ne s'expriment pas de façon simple en fonction de p (\Rightarrow pas la courbe en forme de U classique), mais ça ne veut pas dire que ce compromis est caduc.

Rappel: Si $\left\{ \begin{array}{l} \theta \in \mathbb{H} \subset \mathbb{R} \\ \hat{\theta} \text{ est un estimateur} \\ \text{modèle } p_\theta(x) dx = P(dx) \text{ régulières} \end{array} \right.$ on note $\left\{ \begin{array}{l} B(\theta) = \mathbb{E}_\theta \hat{\theta} - \theta \\ V(\theta) = \text{Var}_\theta(\hat{\theta}) \end{array} \right.$, alors

Cramer-Rao
(\approx Cauchy-Schwarz)

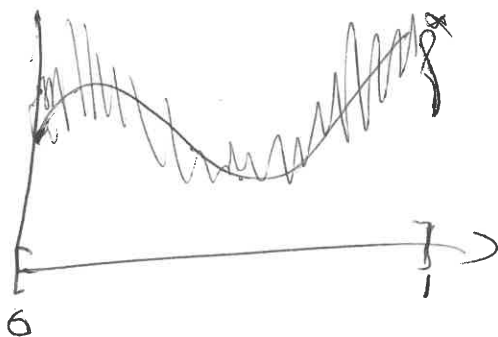
$$V(\theta) \geq \frac{(1 + B'(\theta))^2}{F(\theta)}$$

\nwarrow Information de Fisher $\mathbb{E}_\theta [\partial^2 \log p_\theta(X)]$

\swarrow Variance minimisée si $B(\theta) = 0$
 \searrow si $B'(\theta) \leq \frac{1}{2}$

Pour des problèmes non-paramétriques (régression), l'existence d'estimateurs non-biaisés n'est pas toujours garantie.

On considère le modèle $\boxed{dY_x = f^*(x) dx + \frac{1}{\sqrt{n}} dW_x}$, $x \in [0, 1]$



Pour $\{x_0 \in [0,1] \text{ tels que } 0 \leq x_0 - h \leq x_0 + h \leq 1, \text{ on considère } h > 0\}$

(4)

$$\hat{f}_h(x_0) := \frac{1}{2h} \int_{x_0-h}^{x_0+h} dY_t$$

$$\text{Bias}_f(\hat{f}_h) = \frac{1}{2h} \int_{x_0-h}^{x_0+h} (f(u) - f(x_0)) du$$

→ Dépend de la régularité de f .

$$\text{Var}_f(\hat{f}_h) = \left(\frac{1}{2h}\right)^2 \text{Var}\left(\int_{x_0-h}^{x_0+h} \frac{dW}{\sqrt{m}}\right) = \frac{1}{2mh}$$

→ Ne dépend pas de f

$$\frac{1}{\sqrt{m}}(B_{x_0+h} - B_{x_0-h}) \sim N(0, \frac{2h}{m})$$

Si $f \in \mathcal{C}^\beta$, $\text{Bias}(\hat{f}_h)^2 \lesssim h^{2\beta}$, de sorte que

$$\left(\sup_{f \in \mathcal{C}^\beta} |\text{Bias}_f(\hat{f}_h)|^{1/\beta}\right) \left(\sup_{f \in \mathcal{C}^\beta} \text{Var}_f(\hat{f}_h)\right) \gtrsim \frac{1}{m}$$

→ Derumigny & Schmidt-Hieber démontrent que c'est vrai pour tout estimateur \hat{f} de biais borné (Sinon $\hat{f} = 0$ a variance nulle $\forall f$)

→ Obligation d'avoir les ordres de grandeur classiques pour être minimax

→ Fonction en régression ponctuelle et $L^2([0,1])$.

→ la démonstration utilise;

- Des inégalités d'information

- Une réduction au modèle de la séquence gaussienne

- aux estimateurs invariants par rotation (\Rightarrow Pao de Stern)
 (pas bon en $\sigma = 0$)

- un argument de symétrie type Borsuk (épaisseur de Bernstein?)

II Double descent en régression featurisée (kernel regression)

1) Features et moindres carrés

On observe $\{Y_i = f^*(X_i) + \varepsilon_i\}$ avec $\begin{cases} X_i \in \mathcal{X} = \mathbb{R}^d \text{ fixes} \\ \varepsilon_i \text{ iid } \mathbb{R} \text{ } \mathbb{E}_\varepsilon[\varepsilon] = 0 \end{cases}$ Échantillon d'entraînement (train)
 $i \leq m$

Ici $f^*: \mathcal{X} \rightarrow \mathbb{R}$ est la fonction de régression qu'on cherche à estimer



Comme on va se placer dans des régimes sur-paramétrés où $f^*(X_i) = Y_i$,
 on va évaluer la qualité des estimateurs $\hat{f}: \mathcal{X} \rightarrow \mathbb{R}$ sur

$Y_j = f^*(x_j) + \varepsilon_j$ avec $\begin{cases} x_j \in \mathcal{X} = \mathbb{R}^d \text{ fixes} \\ \varepsilon_j \text{ iid } \mathbb{R} \text{ } \mathbb{E}_\varepsilon[\varepsilon] = 0 \end{cases}$ Échantillon de test (test)
 $j \leq m$

Risque de généralisation: $\forall g: \mathcal{X} \rightarrow \mathbb{R}$,

$$\begin{aligned} R_{\text{test}}(g) &:= \sum_{j=1}^m \mathbb{E}_{\varepsilon_j} (Y_j - g(x_j))^2 \\ &= \mathbb{E}_{\varepsilon} \|Y_{\text{test}} - g(X_{\text{test}})\|_{\ell^2(\mathbb{R}^m)}^2 \end{aligned}$$

$\varepsilon \in \mathbb{R}^m$ $\varepsilon \in (\mathbb{R}^m)^m$

Pour construire un estimateur, on se base sur l'échantillon empirique d'entraînement ⑥

$$\hat{L}_{\text{train}}(g) := \left\| \underset{\in \mathbb{R}^n}{Y_{\text{train}}} - g \left(\underset{\in (\mathbb{R}^d)^n}{X_{\text{train}}} \right) \right\|_{\ell^2(\mathbb{R}^{n_{\text{train}}})}^2$$

Pour tout vectoriser et prendre en compte les non-linéarités, on se donne une

feature map $\boxed{\phi: \mathcal{X} = \mathbb{R}^d \longrightarrow \mathbb{R}^p}$

! $\phi = \phi_p$ avec p qui variera

Ex: $\mathcal{X} = \mathbb{R}$, $\phi(x) = (1, x, x^2, \dots, x^{p-1})$

monomes

$\mathcal{X} = [-1, 1]$, $\phi(x) = (T_0(x), \dots, T_{p-1}(x))$

Chebyshev

$\mathcal{X} = \mathbb{R}^d$, $\phi(x) = (\cos(\langle v_1, x \rangle), \dots, \cos(\langle v_p, x \rangle))$

Random Fourier Features

$\mathcal{X} = \mathbb{R}^d$, $\phi(x) = (\langle v_1, x - a_1 \rangle_+, \dots, \langle v_p, x - a_p \rangle_+)$

$v_i \stackrel{\text{i.i.d.}}{\sim} N(0, 1)$

Réseaux de neurones ReLU à 1 couche (Perceptron)

Pour tout $\beta \in \mathbb{R}^p$,
 $x \in \mathcal{X}$

$$\boxed{g_\beta(x) := \langle \phi(x), \beta \rangle_{\ell^2(\mathbb{R}^p)} = \phi(x)^T \beta}$$

est le régresseur associé.

En notant $\Phi_{\text{train}} := \begin{pmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_m)^T \end{pmatrix} \in \mathbb{R}^{m \times p}$, on obtient l'expression matricielle

$$\hat{L}_{\text{train}}(g_\beta) = \left\| Y_{\text{train}} - \Phi_{\text{train}} \beta \right\|_{\ell^2(\mathbb{R}^m)}^2$$

2) Moindres carrés explicites et biais implicites

$\beta \mapsto \hat{L}_{\text{train}}(g_\beta) = \|Y_{\text{train}} - \Phi_{\text{train}} \beta\|_{\ell^2(\mathbb{R}^n)}^2$ est une forme quadratique minorée par 0, elle admet donc un minimum global (point critique). De plus,

$$\begin{aligned} \hat{L} &= \langle \Phi_{\text{train}} \beta - Y_{\text{train}}, \Phi_{\text{train}} \beta - Y_{\text{train}} \rangle \\ &= \langle \beta, \Phi_{\text{train}}^T \Phi_{\text{train}} \beta \rangle - \langle 2 \Phi_{\text{train}}^T Y_{\text{train}}, \beta \rangle + \langle Y_{\text{train}}, Y_{\text{train}} \rangle \end{aligned}$$

$$\text{D'où } \nabla_\beta \hat{L} = 2 \Phi_{\text{train}}^T (\Phi_{\text{train}} \beta - Y_{\text{train}}) = 0$$

$$\Leftrightarrow \Phi_{\text{train}}^T \Phi_{\text{train}} \beta = \Phi_{\text{train}}^T Y_{\text{train}}$$

Aparté algèbre linéaire

Si $A \in \mathbb{R}^{u \times v}$, $A^+ \in \mathbb{R}^{v \times u}$ est appelée *inverse généralisée* de A lorsque

de Moore-Penrose

$$(i) \quad AA^+A = A$$

$$(ii) \quad A^+AA^+ = A^+$$

) Existe toujours

$$(iii) \quad AA^+ \text{ est symétrique } (= \pi_{\text{Im}(A)})$$

$$(iv) \quad A^+A \text{ est symétrique } (= \pi_{\text{Im}(A^T)} = \pi_{\text{Ker}(A)^\perp})$$

) Rend A^+ unique

Si $Ax = b$ admet au moins une solution, alors $x = A^+b$ est la solution de norme minimale, et toute solution s'écrit

$$x = A^+b + \underbrace{(I_{v \times v} - A^+A)}_{\pi_{\text{Ker}(A)}} w, \quad w \in \mathbb{R}^v$$

Construction avec SVD: Si $A = U \Sigma V^T$ avec $\left| \begin{array}{l} U \in \mathcal{O}_{u \times u} \\ V \in \mathcal{O}_{v \times v} \\ \Sigma \text{ diagonale } \mathbb{R}^{u \times v} (\sigma_p)_p \end{array} \right.$

alors $A^+ = V \Sigma^+ U$

où $\Sigma^+ = (\sigma_p^+)_p, \sigma_p^+ = \begin{cases} \sigma_p^{-1} & \text{si } \sigma_p \neq 0 \\ 0 & \text{sinon} \end{cases}$

Construction analytique $A^+ = \lim_{\lambda \searrow 0} (A^T A + \lambda I_{v \times v})^{-1} A^T$

$= \lim_{\lambda \searrow 0} A^T (A A^T + \lambda I_{u \times u})^{-1}$

Cette dernière expression permet de voir que l'inverse généralisé intervient naturellement lorsqu'on obtient $\hat{\beta}$ par descente de gradient:

Appliquons une descente de gradient à pas fixe $\eta > 0$ partant de $\hat{\beta}^0 = 0$ pour $\beta \mapsto \|Y_{\text{train}} - \Phi_{\text{train}} \beta\|^2$:

$\beta^{t+1} \stackrel{\text{GD}}{=} \beta^t - \eta \times 2 \Phi_{\text{train}}^T (\Phi_{\text{train}} \hat{\beta}^t - Y_{\text{train}})$

$\stackrel{\substack{\lambda \searrow 0 \\ \hat{\beta}^0 = 0}}{=} \sum_{k=0}^t (\mathbf{I}_{p \times p} - 2\eta \Phi_{\text{train}}^T \Phi_{\text{train}})^k \cdot 2\eta \Phi_{\text{train}}^T Y_{\text{train}}$

Série de Von Neuman

$2\eta < \|\Phi_{\text{train}}\|_{\text{op}}^2 \xrightarrow{t \rightarrow \infty} (\mathbf{I} - (\mathbf{I} - 2\eta \Phi_{\text{train}}^T \Phi_{\text{train}}))^+ \cdot 2\eta \Phi_{\text{train}}^T Y_{\text{train}}$

$= (\Phi_{\text{train}}^T \Phi_{\text{train}})^+ \Phi_{\text{train}}^T Y_{\text{train}} = \Phi_{\text{train}}^+ \underbrace{(\Phi_{\text{train}}^T)^+ \Phi_{\text{train}}^T}_{\pi_{\text{Ker}(\Phi_{\text{train}}^T)^\perp} = \text{Ker}(\Phi_{\text{train}}^+)^+} Y_{\text{train}}$

$= \Phi_{\text{train}}^+ Y_{\text{train}}$

Fin de l'aparté.

3) Étude de l'erreur de généralisation

• Le régresseur moindres carrés est $g_{\hat{\beta}_{LS}}(x) = \langle \hat{\beta}_{LS}, \phi(x) \rangle$
 $= \langle (\Phi_{\text{train}}^+)^T \phi(x), Y_{\text{train}} \rangle_{\ell^2(\mathbb{R}^p)}$
! indépendant de Y_{train} !

• Ce régresseur ne peut pas faire mieux que le meilleur g_{β} pour l'erreur de test

$$\beta^* \in \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \mathbb{E}_{\varepsilon_{\text{test}}} \|Y_{\text{test}} - g_{\beta}(X_{\text{test}})\|_{\ell^2(\mathbb{R}^m)}^2$$

$$= \mathbb{E}_{\varepsilon_{\text{test}}} [\|X_{\text{test}} - g_{\beta}(X_{\text{test}})\|_{\ell^2(\mathbb{R}^m)}^2] + \mathbb{E}_{\varepsilon_{\text{test}}} [\|\varepsilon_{\text{test}}\|_{\ell^2(\mathbb{R}^m)}^2]$$

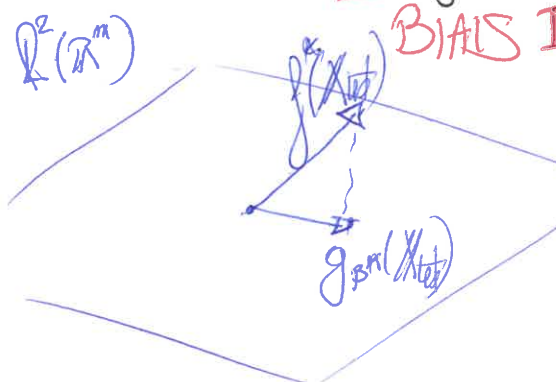
On étudie donc l'excès de risque, qui s'écrit par polarisation

$$R(g_{\beta}) - R(g_{\beta^*}) = \mathbb{E}_{\varepsilon} \left[\langle g_{\beta^*}(X_{\text{test}}) - g_{\beta}(X_{\text{test}}), \underbrace{2f^*(X_{\text{test}})}_{\text{centre}} + 2\varepsilon_{\text{test}} - g_{\beta^*}(X_{\text{test}}) - g_{\beta}(X_{\text{test}}) \rangle_{\ell^2(\mathbb{R}^m)} \right]$$

BIAS I (modèle) (nul si $p \gg m$)

$$= \mathbb{E} \left[\langle g_{\beta^*}(X_{\text{test}}) - g_{\beta}(X_{\text{test}}), 2(f^* - g_{\beta^*}) + (g_{\beta^*} - g_{\beta}) \rangle \right]$$

$$= \|g_{\beta^*}(X_{\text{test}}) - g_{\beta}(X_{\text{test}})\|_{\ell^2(\mathbb{R}^m)}^2$$



D'autre part, on a $\boxed{\beta^* = \Phi_{\text{test}}^+ f^*(X_{\text{test}})}$ si l'on considère l'élément de norme minimale

Comme $g_{\beta}(X_{\text{test}}) = \Phi_{\text{test}} \beta$ par définition, on obtient

$$R(g_{\hat{\beta}_{LS}}) - R(g_{\beta^*}) = \left\| \Phi_{\text{test}} (\Phi_{\text{test}}^+)^* (X_{\text{test}}) - \Phi_{\text{train}}^+ Y_{\text{train}} \right\|_{\ell^2(\mathbb{R}^m)}^2$$

On peut enfin recoller au papier "U-Turn on Double Descent" en effectuant une décomposition biais variance en l'alea ϵ_{train} (i.e Y_{train})

$$\mathbb{E}_{\epsilon_{\text{train}}} (R(g_{\hat{\beta}_{LS}}) - R(g_{\beta^*})) = \underbrace{\left\| \Phi_{\text{test}} (\Phi_{\text{test}}^+)^* (X_{\text{test}}) - \Phi_{\text{train}}^+ (X_{\text{test}}) \right\|_{\ell^2(\mathbb{R}^m)}^2}_{\text{Biais II (Transfer)}} + \underbrace{\mathbb{E}_{\epsilon_{\text{train}}} \left\| \Phi_{\text{test}} \Phi_{\text{train}}^+ \epsilon_{\text{train}} \right\|_{\ell^2(\mathbb{R}^m)}^2}_{\text{Variance}}$$

Hypothèse: $\text{Cov}(\epsilon) = \sigma^2 I_{m \times m}$.

Alors la variance s'écrit

$$\mathbb{E}_{\epsilon_{\text{train}}} \left\| \Phi_{\text{test}} \Phi_{\text{train}}^+ \epsilon_{\text{train}} \right\|^2 = \sigma^2 \left\| \Phi_{\text{test}} \Phi_{\text{train}}^+ \right\|_{\text{Frobenius}}^2$$

$p_{\text{effective}}$: dimension effective du modèle.

$$\begin{aligned} \mathbb{E} \|A\epsilon\|^2 &= \mathbb{E} \text{Tr}(A\epsilon\epsilon^T A^T) \\ &= \text{Tr}(A \mathbb{E}(\epsilon\epsilon^T) A^T) \\ &= \sigma^2 \text{Tr}(A A^T) \\ &= \sigma^2 \|A\|_F^2 \end{aligned}$$

→ Somme sur tous les $(x_j)_j$ (du test) de

$$\begin{aligned} \text{Var}_{\epsilon} (g_{\hat{\beta}_{LS}}(x_j)) &= \text{Var}_{\epsilon} (\langle \Phi_{\text{train}}^+ Y_{\text{train}}, \phi(x_j) \rangle) \\ &= \text{Var}_{\epsilon} (\langle \epsilon_{\text{train}}, (\Phi_{\text{train}}^+)^T \phi(x_j) \rangle) \\ &= \sigma^2 \|(\Phi_{\text{train}}^+)^T \phi(x_j)\|^2 \end{aligned}$$

4) Interpretation

2. Géométrie de la variance

Pour le terme de variance uniquement, on a $\forall x_0 \in \mathcal{X}$ fixe

$$\text{Var}_{\mathcal{E}}(g_{\mathcal{P}_{\mathcal{S}}}^{\perp}(x_0)) = 0 \iff \phi(x_0) \in \text{Ker}((\Phi_{\text{train}}^+)^T)$$

$$\iff \phi(x_0) \in \text{Im}(\Phi_{\text{train}}^+)^{\perp}$$

$$\iff \phi(x_0) \in \text{Im}(\Phi_{\text{train}}^T)^{\perp}$$

$$\iff \phi(x_0) \in \text{Ker}(\Phi_{\text{train}})$$

$$\implies g_{\mathcal{P}_{\mathcal{S}}}^{\perp}(x_0) = 0$$

↳ Pour une donnée dans une région de l'espace jamais vue, on répond 0 et donc
 } Variance nulle
 } possible grand biais

↳ De même, $\phi(x_{\text{test}})$ quasi orthogonal à $\text{Vect}(\phi(x_1), \dots, \phi(x_n))$ renvoie 0

En sommant sur $\mathcal{X}_{\text{test}}$, on mélange alors des régions de l'espace avec
 moyennise

possiblement / petit biais / grande variance
 grand biais / petite variance

B Grande dimension et petite variance

(Si $A: E \rightarrow F$, $\text{rang } A + \dim \text{Ker } A = \dim E$) (12)

Par théorème du rang,

$$\dim \text{Ker}(\Phi_{\text{train}}) + \underbrace{\text{rang } \Phi_{\text{train}}}_{\leq n} = p$$

$$\Rightarrow \boxed{\dim \text{Ker}(\Phi_{\text{train}}) \geq p - n}$$

\Rightarrow Quand $p \nearrow$, la taille de l'espace où la variance est petite \nearrow

Cela explique pourquoi la variance décroît après $p=n$, mais ça néglige complètement le biais (☹)

Rq: Possibilité d'écrire tout ça avec des moyennes

↳ Décroissance du spectre

↳ Matrice de Gram croisée entre X_{train} et X_{test}

III Ce que semble présenter "A U-Turn"

Le message principal est que pour les estimateurs linéaires (smoothers) du type

$$\hat{g}(x_0) = \langle \underbrace{\hat{\Delta}(x_0)}_{\in \mathbb{R}^n}, Y_{\text{train}} \rangle,$$

le bon paramètre de dimension effective est

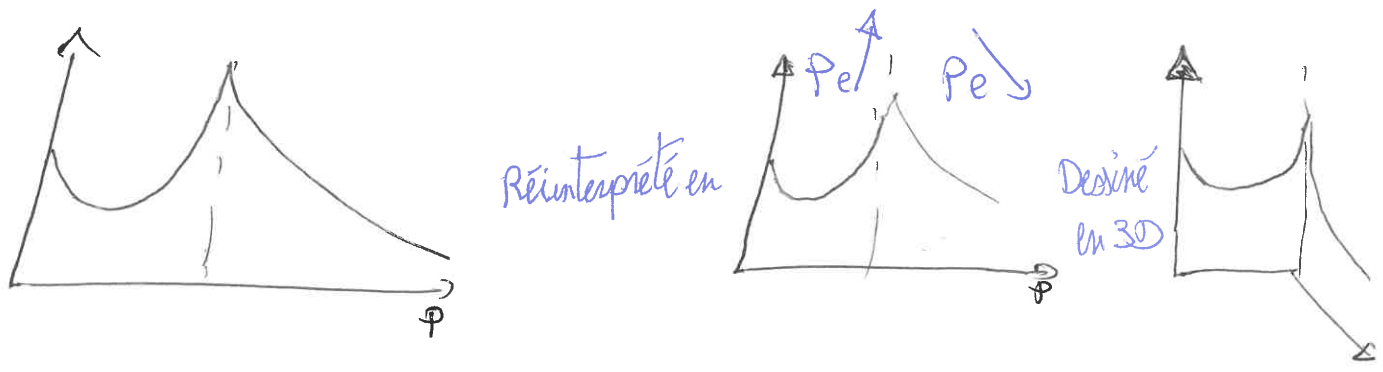
$$\begin{aligned} p_{\text{eff}} &:= \frac{1}{\sigma^2} \sum_{j \in \text{test}} \text{Var}_{\varepsilon} \hat{g}(x_j) \\ &= \frac{1}{\sigma^2} \sum_{j \in \text{test}} \|\hat{\Delta}(x_j)\|^2 \end{aligned}$$

→ Plus en lumière de la notion de transfert (DEUX designs $\begin{cases} \rightarrow X_{\text{train}} \\ \rightarrow X_{\text{test}} \end{cases}$)

→ Inclut :

- Régression à noyau
- Forêts aléatoires
- Gradient Boosting

→ Simulation / Données réelles montrant un compromis biais variance classique si l'on met p_{eff} en abscisse

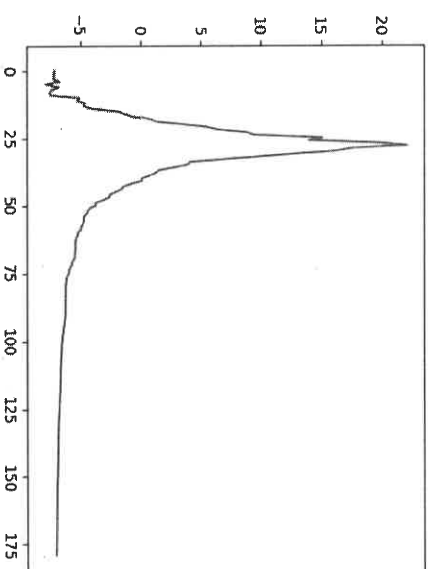


Unebychev polynomials

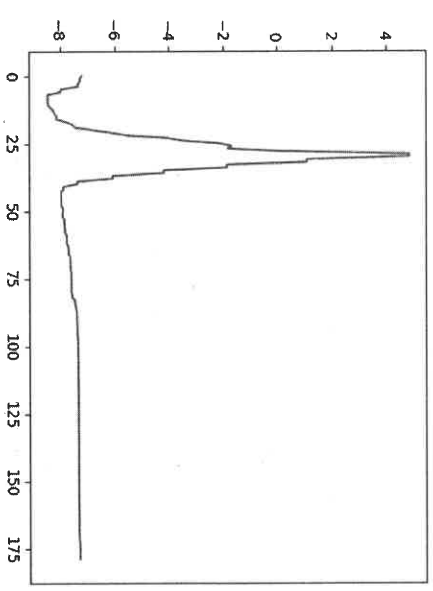
$$n_{\text{train}} = 30, n_{\text{test}} = 1000000$$

train
test

Random



Equidistant



Random

Equidistant

