

Regression à noyau, Neural Tangent Kernel et Overfitting bête

Eddie Amari
28/09/2023

I Régression et espaces à noyau reproduisant

1) Régression linéaire

Données $(X_i, Y_i) \in \mathcal{X} \times \mathcal{Y} = \mathbb{R}^d \times \mathbb{R}$

On cherche à trouver un minimiseur pour

$$\min_{\beta \in \mathbb{R}^d} \sum_{i=1}^m (Y_i - \langle X_i, \beta \rangle)^2 = \min_{\beta \in \mathbb{R}^d} \|Y - X\beta\|_2^2$$

où $X = \begin{pmatrix} X_1^T \\ \vdots \\ X_m^T \end{pmatrix} \in \mathbb{R}^{m \times d}$
 $Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix}$

Points critiques $\Leftrightarrow X^T(X\hat{\beta} - Y) = 0$

$$\Leftrightarrow X^T X \hat{\beta} = X^T Y$$

($m > d$) • Si $X^T X$ inversible, $\hat{\beta} = (X^T X)^{-1} X^T Y$

($m < d$) • Sinon, infinité de solutions à $Y = X\beta$

Rq: $XX^T = (\langle X_i, X_j \rangle)_{i,j=1}^m$
matrice de Gram

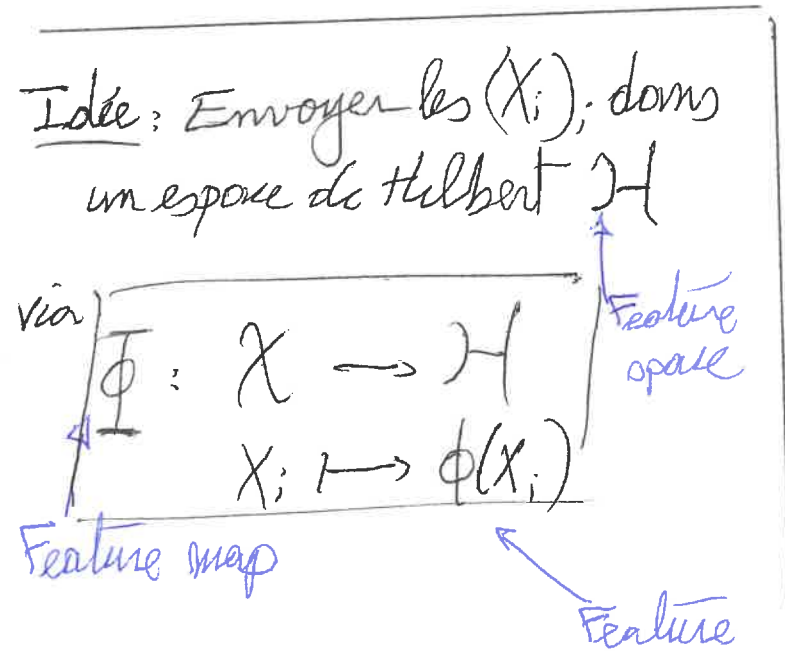
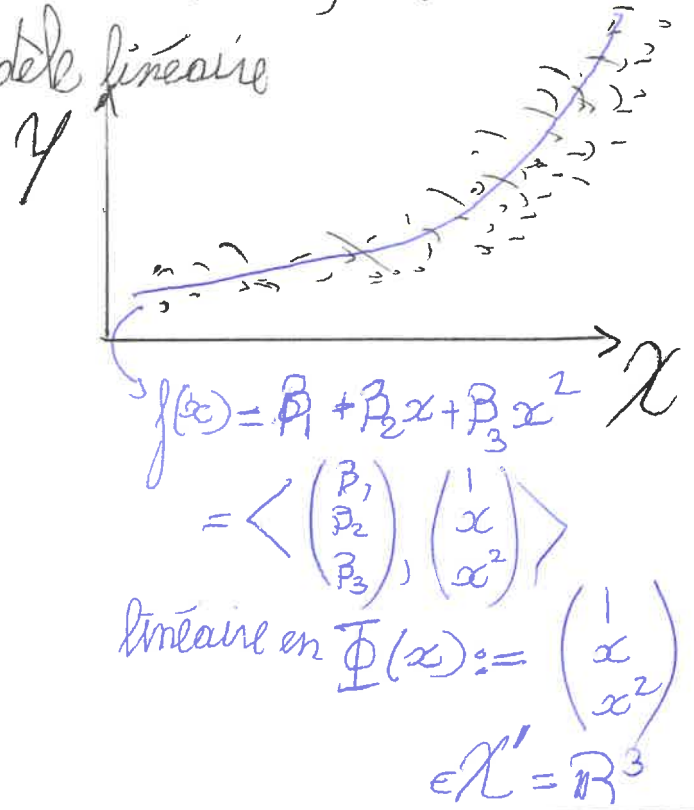
↳ Solution de norme minimale $\hat{\beta} = \arg \min_{X\beta=Y} \frac{\|\beta\|_2^2}{2}$

$$= \underbrace{X^T (X X^T)^{-1}}_{\hat{\alpha} \in \mathbb{R}^m} Y = \sum_{i=1}^m \hat{\alpha}_i X_i$$

⇒ $\hat{\beta}$ combinaison linéaire des $(X_i)_i$

2) Machines à noyau

Au lieu de travailler directement avec les $(X_i)_{i \leq m}$, on peut d'abord les transformer, puis considérer le modèle linéaire



Travailler dans le nouvel espace \mathcal{H} et résoudre le problème linéaire associé semble nécessiter de calculer $(\langle \Phi(X_i), \Phi(X_j) \rangle_{\mathcal{H}})_{i,j \leq m}$ via le calcul des $\Phi(X_i)$ et des $\langle \cdot, \cdot \rangle_{\mathcal{H}}$.

Astuce du noyau (Kernel trick)

Il est parfois plus facile de calculer le noyau

$$k: \mathcal{X} \times \mathcal{X} \longrightarrow \mathbb{R}$$

$$(x_i, x_j) \longmapsto \langle \Phi(x_i), \Phi(x_j) \rangle_{\mathcal{H}}$$

que de calculer Φ , qui vit possiblement en dimension infinie.

Inversement: On va voir que se donner un "noyau" $k: X \times X \rightarrow \mathbb{R}$ est équivalent à se donner ϕ . (3)

Def: (Noyau) $k: X \times X \rightarrow \mathbb{R}$ est appelée noyau lorsque

- k est symétrique, continue

- k est positive:

$$\forall n \geq 1, \left\{ \begin{array}{l} \forall x_1, \dots, x_n \in X \\ \forall a_1, \dots, a_n \in \mathbb{R} \end{array} \right\}, \sum_{i=1}^n \sum_{j=1}^n a_i a_j k(x_i, x_j) \geq 0$$

On appelle matrice de Gram $K := (k(x_i, x_j))_{i,j \leq n} \in \mathbb{R}^{n \times n}$

Req: • la définition de noyau est exactement la contrainte $K \succeq 0$

• $K(x_i, x_j) := g(\|x_i - x_j\|)$ est un noyau $\Leftrightarrow g \geq 0$
 Bochner
 Noon

Thm: (Moore - Aronszajn)

Pour tout noyau $k: X \times X \rightarrow \mathbb{R}$, il existe | • un hilbert $\mathcal{H} \subset \mathcal{C}(X, \mathbb{R})$
• $\Phi: X \rightarrow \mathcal{H}$

tel que $\forall x, x' \in X, k(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}}$

De plus, on a la propriété de noyau reproductant

$$\forall f \in \mathcal{H}, f(x) = \langle f, k(x, \cdot) \rangle_{\mathcal{H}}$$

Dem: On pose $\mathcal{H} := \left\{ \sum_{j=1}^N a_j K(x_j, \cdot) \mid \begin{array}{l} N \geq 1 \\ a_j \in \mathbb{R} \\ x_j \in \mathcal{X} \end{array} \right\}$

où la clôture est prise pour $\langle \sum_j a_j K(x_j, \cdot), \sum_j b_j K(y_j, \cdot) \rangle_{\mathcal{H}}$
 $= \sum_{j,l} a_j b_l K(x_j, y_l)$

Ex: Noyau polynômial $k(x, x') = (\langle x, x' \rangle + c)^p$

↳ Feature map $\Phi(x) = \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2} x_1 x_2 \\ \sqrt{2} x_1 \\ \sqrt{2} x_2 \\ 1 \end{pmatrix} \quad (p=2, d=2)$

Complexité: $k(x, x') : O(d)$

$\Phi(x) : O\left(\binom{d+2}{2}\right)$

$\uparrow \dim \mathbb{R}_{\leq 2}[X_1, \dots, X_d]$

• Radial Basis Function: $k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$

$\Phi(x) = e^{-\frac{x^2}{2\sigma^2}} \left(1, \sqrt{\frac{1}{1!\sigma^2}} x, \sqrt{\frac{1}{2!\sigma^2}} x^2, \dots \right)$

$\hookrightarrow \dim \mathcal{H} = +\infty$

Il se trouve que bien que travaillant désormais (potentiellement) en dimension infinie, on peut se limiter à la dimension n pour des problèmes d'optimisation

Thm: (Théorème du représentant)

Soient $X_1, \dots, X_m \in \mathcal{X}$, $\Phi: \mathcal{X} \rightarrow \mathcal{H}$ et $g: \mathbb{R} \rightarrow \mathbb{R}$ strictement croissante

Alors

$$\begin{aligned} \min_{h \in \mathcal{H}} L(\underbrace{h(X_1), \dots, h(X_m)}_{= \langle h, \Phi(X_1), \dots, \Phi(X_m) \rangle_{\mathcal{H}}}) + g(\|h\|_{\mathcal{H}}) \\ = \langle h, \Phi(X_1) \rangle_{\mathcal{H}} \\ = \langle h, \Phi(X_1) \rangle_{\mathcal{H}} \end{aligned}$$

on admet que des solutions de la forme

$$h^* = \sum_{i=1}^m \alpha_i \Phi(X_i) = \sum_{i=1}^m \alpha_i \Phi(X_i)$$

↳ Analogie avec la régression linéaire :

les solutions de norme minimale sont une combinaison linéaire des descripteurs $\Phi(X_i)$ (et donc vivent en $\dim \leq m$)

Application: Kernel ridge regression

$$\min_{h \in \mathcal{H}} \sum_{i=1}^m (Y_i - \langle h, \Phi(X_i) \rangle)^2 + \frac{\lambda}{2} \|h\|_{\mathcal{H}}^2$$

$$\text{En notant } \left\{ \begin{array}{l} X: \mathcal{H} \rightarrow \mathbb{R}^m \\ h \mapsto \begin{pmatrix} \langle h, \Phi(X_1) \rangle_{\mathcal{H}} \\ \vdots \\ \langle h, \Phi(X_m) \rangle_{\mathcal{H}} \end{pmatrix} \end{array} \right\} \quad \text{et} \quad \left\{ \begin{array}{l} X^T: \mathbb{R}^m \rightarrow \mathcal{H} \\ \alpha \mapsto \sum_{i=1}^m \alpha_i \Phi(X_i) = \sum_{i=1}^m \alpha_i \Phi(X_i) \end{array} \right\}$$

$$\text{on pose } \hat{\Sigma} := X^T X = \sum_{i=1}^m \Phi(X_i) \otimes \Phi(X_i)$$

⑥

$$\text{On obtient } h^* = (\hat{\Sigma} + \lambda I_{\mathcal{H}})^{-1} X^T Y$$

$$= X^T (\lambda I_{\mathbb{R}^n} + X X^T)^{-1} Y$$

le régresseur final est alors

$$\begin{aligned} \hat{f}_{KRR}^\lambda(x) &:= \langle h^*, \phi(x) \rangle_{\mathcal{H}} \\ &= \langle X^T (\lambda I_{\mathbb{R}^n} + X X^T)^{-1} Y, \phi(x) \rangle_{\mathcal{H}} \\ &= \langle (\lambda I_{\mathbb{R}^n} + X X^T)^{-1} Y, X(\phi(x)) \rangle_{\mathbb{R}^n} \\ &= \left\langle (\lambda I_{\mathbb{R}^n} + X X^T)^{-1} Y, \begin{pmatrix} k(X_1, x) \\ \vdots \\ k(X_m, x) \end{pmatrix} \right\rangle_{\mathbb{R}^n} \end{aligned}$$

qui ne dépend que de $k(\cdot, \cdot)$.

Rq: Pour $\lambda = 0$, on obtient $\hat{f}_{KR}^\lambda(x) = \langle K^{-1} Y, K(X, x) \rangle_{\mathbb{R}^n}$

$$\text{où } \begin{cases} K = (k(X_i, X_j))_{i,j \leq m} \\ K(X, x) = \begin{pmatrix} k(X_1, x) \\ \vdots \\ k(X_m, x) \end{pmatrix} \end{cases}$$

II Réseaux de neurones et Neural Tangent Kernel

1) Régression avec réseaux de neurones feedforward

Def: Un réseau de neurones (feedforward) à L couches cachées
de fonction d'activation $\Psi: \mathbb{R} \rightarrow \mathbb{R}$

et une fonction $f^{(L)}(\cdot, W): \mathbb{R}^d \rightarrow \mathbb{R}$ définie par

$$f^{(L)}(x, W) = W^{(L+1)} \frac{1}{\sqrt{k_L}} \Psi \left(W^{(L)} \frac{1}{\sqrt{k_{L-1}}} \Psi \left(\dots W^{(2)} \frac{1}{\sqrt{k_1}} \Psi(W^{(1)} x) \right) \right),$$

où $\begin{cases} W^{(i)} \in \mathbb{R}^{k_i \times k_{i-1}} \\ k_0 = d, k_L = 1 \end{cases}$

Exemple d'utilisation: la base $(f^{(L)}(\cdot, W))_{W \in \mathbb{R}^p}$ peut être insérée dans une méthode par moindres carrés, cette fois-ci non-linéaire si Ψ est non-linéaire :

$$W^* \in \operatorname{argmin}_{W \in \mathbb{R}^p} \sum_{i=1}^n (Y_i - f^{(L)}(X_i, W))^2,$$

qui peut être obtenu par descente de gradient (stochastique) autour d'une valeur initiale $W_0 \in \mathbb{R}^p$.

2) Linearization: cf. Neural Tangent Kernels (8)
 Si la SGD a des trajectoires qui restent dans le domaine de validité du développement limité

$$f^{(L)}(x, w) \approx f^{(L)}(x, w_0) + \langle \nabla_w f^{(L)}(x, w_0), x - x_0 \rangle$$

alors $w_{\text{SGD}}^* \in \underset{w \in \mathbb{R}^p}{\text{argmin}} \sum_{i=1}^m (y_i - (f^{(L)}(x_i, w_0) + \langle \nabla_w f^{(L)}(x_i, w_0), w - w_0 \rangle))^2$

ce qui revient essentiellement à faire une régression linéaire sur les features

$$\Phi(x) = \nabla_w f^{(L)}(x, w_0)$$

le noyau associé à ce descripteur est appelé Neural Tangent Kernel ou NTK

Def: (NTK) $\forall x, \tilde{x} \in \mathbb{R}^d, \forall w \in \mathbb{R}^p, K_w^{(L)}(x, \tilde{x}) := \langle \nabla_w f^{(L)}(x, w), \nabla_w f^{(L)}(\tilde{x}, w) \rangle$

↳ [Jacot et al, Neurips 2018]

3) Validité de la linéarisation en grande largeur

Dans [Liu, Zhu, Belkin - Neurips 2020] les auteurs démontrent une validité uniforme de la constante du NTK pour des initialisations w_0 gaussiennes

Thm: Soit L est fixé on a pour tout $x, \tilde{x} \in \mathbb{R}^d$

$\bullet k_1, \dots, k_L \rightarrow \infty$

$W_{a,b}^{(i)} \stackrel{\text{ind}}{\sim} N(0,1) \Rightarrow$

$$\begin{cases} \bullet \sup_{w' \in B(w,1)} |K_{w'}^{(L)}(x, \tilde{x}) - K_w^{(L)}(x, \tilde{x})| \rightarrow 0 \\ \bullet K_{w'}^{(L)}(x, x) = \Omega(1) \quad \forall w' \in B(w,1) \end{cases}$$

Autrement dit, le NTK est constant et non-trivial sur une boule de rayon constant:

$$B(w,1) \rightarrow \mathbb{R}^p$$
$$w' \mapsto \nabla_w f^{(L)}(x, w')$$

"constant" \nwarrow Abrévié car $p \rightarrow \infty$

En particulier, si on utilise une méthode d'optimisation à l'ordre 1,

$$\text{Regression } L^2 \text{ sur } \{f(x, w)\}_{w \in \mathbb{R}^p} \Leftrightarrow \text{Regression à noyau } K_w^{(L)}(\cdot, \cdot)$$

\uparrow Informel

Une fois ceci acté, on peut se demander

- Quelle tête a $K_w^{(L)}$ en grande largeur?
- ————— $K_w^{(L)}$ ————— Profondeur?
- À la limite, obtient-on une méthode statistiquement intéressante?

III Convergence du NTK et application en classification

[Radhakrishnan, Belkin, Uhler, PNAS 2023 - Wide & Deep neural networks achieve consistency for classification]

1) NTK en largeur infinie et fonction d'activation dual

a) Une couche cachée

Pour commencer, on considère $L=1$ et

$$f(x, w) = \frac{1}{\sqrt{k}} A \Psi(Bx)$$

$$\text{où } \begin{cases} A \in \mathbb{R}^{1 \times k} \\ B \in \mathbb{R}^{k \times d} \end{cases}, \quad B = \begin{pmatrix} B_{1:} \\ \vdots \\ B_{k:} \end{pmatrix} \quad \boxed{\begin{matrix} w = (A_{11}, \dots, A_{1k}, B_{11}, \dots, B_{kd}) \\ \in \mathbb{R}^{k+kd} \quad (p = k+kd) \end{matrix}}$$

On suppose que $\begin{cases} A_{1i}, B_{ij} \stackrel{\text{iid}}{\sim} N(0,1) \\ x, \tilde{x} \in \mathbb{R}^d \text{ sont normalisés } \|x\| = \|\tilde{x}\| = 1 \end{cases}$ (Initialisation standard des NN)

$$\text{Par définition, } f(x, w) = \frac{1}{\sqrt{k}} \sum_{i=1}^k A_{1i} \Psi(\underbrace{B_{i:} x}_{= \sum_{j=1}^d B_{ij} x_j})$$

$$\text{d'où } \frac{\partial f}{\partial A_{1i}}(x, w) = \frac{1}{\sqrt{k}} \Psi(B_{i:} x) \quad \text{et} \quad \frac{\partial f}{\partial B_{ij}}(x, w) = \frac{1}{\sqrt{k}} A_{1i} \Psi'(B_{i:} x) x_j$$

En particulier,

(11)

$$K^{(1)}(x, \tilde{x}) = \left\langle \nabla_w f^{(1)}(x, w), \nabla_w f^{(1)}(\tilde{x}, w) \right\rangle$$

$$= \underbrace{\frac{1}{k} \sum_{i=1}^k \psi(B_{i,:} x) \psi(B_{i,:} \tilde{x})}_{(a)} + \underbrace{\frac{1}{k} \sum_{i=1}^k \sum_{j=1}^d A_{i,j}^2 \psi'(B_{i,:} x) \psi'(B_{i,:} \tilde{x}) x_j \tilde{x}_j}_{(b)}$$

(a) et (b) font intervenir des sommes iid de variables aléatoires dépendant de $(B_{i,:} x, B_{i,:} \tilde{x})$.

Comme $B_{i,:} \sim \mathcal{N}(0, \text{Id}_{1 \times d})$, $(B_{i,:} x, B_{i,:} \tilde{x}) \in \mathbb{R}^2$ est un vecteur gaussien

De plus, $\mathbb{E}[B_{i,:} x] = \mathbb{E}[B_{i,:}] x = 0$ et

$$\text{cov}(B_{i,:} x, B_{i,:} \tilde{x}) = \mathbb{E}[x^T B_{i,:}^T B_{i,:} \tilde{x}]$$

$$= x^T \tilde{x} \quad \text{car } \mathbb{E}[B_{i,:}^T B_{i,:}] = \text{Id}_{1 \times d}$$

$$= \langle x, \tilde{x} \rangle$$

Comme $\langle x, x \rangle = \langle \tilde{x}, \tilde{x} \rangle = 1$, on a donc

$$\begin{pmatrix} B_{i,:} x \\ B_{i,:} \tilde{x} \end{pmatrix} \sim \mathcal{N}(0, \Lambda) \quad \text{où} \quad \Lambda := \begin{pmatrix} 1 & \langle x, \tilde{x} \rangle \\ \langle x, \tilde{x} \rangle & 1 \end{pmatrix}$$

Par la loi des grands nombres :

$$(a) \xrightarrow[k \rightarrow \infty]{p.d} E \left[\Psi(u) \Psi(v) \right]_{(u,v) \sim N(0,1)}$$

$$(b) \xrightarrow[k \rightarrow \infty]{p.d} \sum_{j=1}^d E \left[A^2 \Psi'(u) \Psi'(v) \right]_{\substack{A \sim N(0,1) \\ (u,v) \sim N(0,1)}} x_j \tilde{x}_j = E[\Psi'(u) \Psi'(v)] \langle x, \tilde{x} \rangle$$

$E[A^2] = 1$

la transformation $\Psi \mapsto E \left[\Psi(u) \Psi(v) \right]_{(u,v) \sim N(0,1)}$ apparaît comme centrale !

$$\Lambda = \begin{pmatrix} 1 & \langle x, \tilde{x} \rangle \\ \langle \tilde{x}, x \rangle & 1 \end{pmatrix}$$

b) Fonction d'activation duale

Def: (Activation duale)

Pour $\Psi: \mathbb{R} \rightarrow \mathbb{R}$ fonction d'activation, sa fonction duale

$$\check{\Psi}: [-1, 1] \rightarrow \mathbb{R}$$

est définie par $\check{\Psi}(\xi) := E \left[\Psi(u) \Psi(v) \right]_{(u,v) \sim N(0,1)}$, où $\Lambda = \begin{pmatrix} 1 & \xi \\ \xi & 1 \end{pmatrix}$

Rq: Pour $\Psi(\gamma) = e^{-\gamma^2/2}$, $K(x, \tilde{x}) = e^{-\frac{\|x - \tilde{x}\|^2}{2}}$

le Neural Network Gaussian Process $K(x, \tilde{x}) = E \left[\Psi(u) \Psi(v) \right]_{(u,v) \sim N(0, \begin{pmatrix} \|x\|^2 & \langle x, \tilde{x} \rangle \\ \langle \tilde{x}, x \rangle & \|\tilde{x}\|^2 \end{pmatrix})}$

donne $K(x, \tilde{x}) = \check{\Psi}(\langle x, \tilde{x} \rangle) \forall x, \tilde{x} \in \mathcal{X}^{d-1}$

les fonctions $\check{\Psi}$ ont beaucoup de bonnes propriétés, et s'étudient via les polynômes de Hermite $h'_i(x) = \sqrt{i} h_i(x)$ qui forment une base orthogonale de $L^2(N(0,1))$. (13)

Prop: Si $\Psi \in L^2(N(0,1))$, alors

- $\Psi(x) = \sum_{i=0}^{\infty} a_i h_i(x)$ a pour duale $\check{\Psi}(\cdot) = \left\{ \sum_{i=0}^{\infty} a_i^2 \right\}^{\frac{1}{2}}$
- $\check{\Psi}$ est croissante convexe sur $[0,1]$, ainsi que toutes ses dérivées
- $(\check{\Psi})' = (\check{\Psi}') \quad (\check{\cdot} \text{ commute avec la dérivation})$
- $\check{\Psi} \in \mathcal{C}^0([-1,1])$ et $\check{\Psi} \in \mathcal{C}^\infty((-1,1))$
- $\check{\Psi}$ est définie positive, au sens où $K: \mathcal{Y}^{d \times 1} \times \mathcal{Y}^{d \times 1} \rightarrow \mathbb{R}$
 $(x, \tilde{x}) \mapsto \check{\Psi}(\langle x, \tilde{x} \rangle)$
 est un noyau (positif)
- Inversement, toute fonction $\sigma: [-1,1] \rightarrow \mathbb{R}$ définie positive est la duale $\check{\Psi} = \sigma$ d'une fonction d'activation

c) Une formule de récurrence pour le NTK en largeur infinie

Pour simplifier, on normalise ψ de sorte que $\psi(1) = 1$

$$\left\{ x, \tilde{x} \in \mathbb{R}^d \longrightarrow \|x\| = \|\tilde{x}\| = 1 \right.$$

Thm: Pour des entrées $W^{(i)} \in \mathbb{R}^{k_i \times k_{i-1}}$ ont des entrées $N(0,1)$ i.i.d, alors quand $k_1 \rightarrow \infty$, puis $k_2 \rightarrow \infty$, ..., puis $k_L \rightarrow \infty$, les NTK convergent presque sûrement vers

$$\left| K^{(L)}(x, \tilde{x}) = \sum^{(L)}(x, \tilde{x}) + K^{(L-1)}(x, \tilde{x}) \psi' \left(\sum^{(L-1)}(x, \tilde{x}) \right) \right.$$

avec $K^{(0)}(x, \tilde{x}) = \langle x, \tilde{x} \rangle$

où $\left| \sum^{(L)}(x, \tilde{x}) = \psi \left(\sum^{(L-1)}(x, \tilde{x}) \right) \right.$

avec $\sum^{(0)}(x, \tilde{x}) = \langle x, \tilde{x} \rangle$

Rq: $\sum^{(L)}(x, \tilde{x})$ correspond à la dérivation par rapport à la dernière couche (linéaire)

$K^{(L-1)}(x, \tilde{x}) \psi' \left(\sum^{(L-1)}(x, \tilde{x}) \right)$ à la dérivation par rapport aux autres couches.

2) NTK en profondeur infinie

a) Énoncé du résultat

On étudie maintenant le NTK lorsque la profondeur $L \rightarrow \infty$, en laissant à chaque fois toutes les largeurs tendre vers ∞

Dans tout ce qui suit, $\{x, \tilde{x} \in \mathcal{Y}^{d-1}\}$ sont fixées

$$\left[\Psi(1) = 1 = \mathbb{E}_{y \sim N(0,1)} [\Psi(y)^2] \right]$$

L'asymptotique de $K^{(L)}$ est gouvernée par trois quantités :

$$\sqrt{\Psi(0)} = \left[\begin{aligned} A &:= \mathbb{E}_{y \sim N(0,1)} [\Psi(y)] , \quad A' := \mathbb{E}_{y \sim N(0,1)} [\Psi'(y)] = \sqrt{\Psi'(0)} \\ B' &:= \mathbb{E}_{y \sim N(0,1)} [\Psi'(y)] = \Psi'(1) \end{aligned} \right]$$

Rq: L'article étudie le classifieur limite $m_m(x) := \text{sign}(\langle K^{-1} \tilde{y}, K(x, x) \rangle)$ pour $(X_i, Y_i) \in \mathbb{R}^d \times \{-1, 1\}$ des données de classification fixées. Les résultats sur l'asymptotique de $K^{(L)}$ se répercutent sur $m_m(x)$

Thm:

$$(A=0, A' \neq 0) \lim_{L \rightarrow \infty} \frac{K^{(L)}(x, \tilde{x})}{(A')^{2L} (L+1)} = \frac{R(\|x - \tilde{x}\|)}{\|x - \tilde{x}\|^\alpha}, \quad \text{où } \begin{cases} \alpha = -2 \frac{\log(A'^2)}{\log(B')} \\ R(\cdot) \geq 0 \text{ bornée loin de zéro en } O. \end{cases}$$

→ Classifieur à noyau singulier pour $\alpha > 0$

(A=0, A'=0) Si $\|x - \tilde{x}_i\| = \min_{j \leq m} \|x - X_j\|$ et que ce minimum est unique,

$$\lim_{L \rightarrow \infty} \frac{K^{(L)}(X_j, x)}{K^{(L)}(X_i, x)} = 0 \quad \forall j \neq i$$

→ Classifieur au plus proche voisin $m(x) = \text{sign}(X_i(x))$

$$(A \neq 0) \forall x \neq \tilde{x}, \lim_{L \rightarrow \infty} \frac{K^{(L)}(x, \tilde{x})}{C(L)} = C, \neq \lim_{L \rightarrow \infty} \frac{K^{(L)}(x_0, x)}{C(L)}$$

→ Classifieur par vote majoritaire $m(x) = \text{sign}\left(\sum_{i=1}^m Y_i\right)$

→ RelU $\Psi(x) = x_+$ tombe dans ce cas

Rq: Dans le cas $A=0, A' \neq 0$, le classifieur peut être consistant quand $m \rightarrow \infty$ car il est équivalent au classifieur de Hilbert [Devroye, Györfi, Krzyżniak - 1998] pour $\alpha = \frac{d}{2}$

→ Et cet estimateur overfite! ($m(X_i) = Y_i$)

b) Idee de demonstration pour $A=0, A' \neq 0$

On note $\eta = \langle x, \tilde{x} \rangle$ et $\check{\Psi}^{(L)}(\eta) = \underbrace{\check{\Psi}_0 \dots \check{\Psi}_1}_{(L+1) \text{ fois}}(\eta)$.

La formule de récurrence sur $K^{(L)}$ donne l'expression explicite

$$K^{(L)}(\eta) = \sum_{i=0}^L \check{\Psi}^{(i)}(\eta) \prod_{j=i}^{L-1} \check{\Psi}'(\check{\Psi}^{(j)}(\eta))$$

Pour étudier cette somme, on commence d'abord par examiner l'asymptotique de $\check{\Psi}^{(L)}(\eta)$

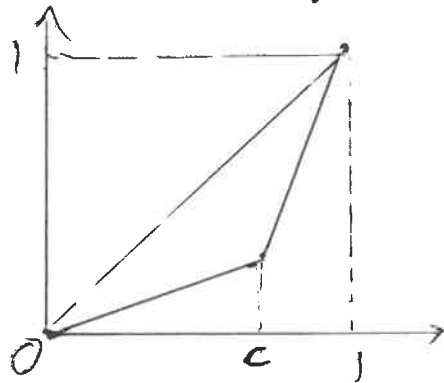
lemme: $\lim_{L \rightarrow \infty} \frac{\check{\Psi}^{(L)}(\eta)}{\check{\Psi}'(0)^L} = \frac{\tilde{R}(\eta)}{(1-\eta)^{\frac{\alpha}{2}}}$, où $\alpha = \frac{2 \log \check{\Psi}'(0)}{\log \check{\Psi}'(1)} = -\frac{2 \log(A'^2)}{\log(B')}$

$0 \neq A'^2 = \dots$

Dem: On traite le cas simple où $f = \check{\Psi}$ est telle que $f: [0, 1] \rightarrow [0, 1]$

$$f(x) = \begin{cases} ax & x \in [0, c] \\ 1 - b(1-x) & x \in [c, 1] \end{cases}$$

avec $c := \frac{b-1}{b-a}$ ($f'(0)=a, f'(1)=b$)



• Pour $x \in [0, c]$, $f^{(L)}(x) = a^L x$, d'où $\lim_{L \rightarrow \infty} \frac{f^{(L)}(x)}{a^L} = x$

• Pour $x \in [c, 1]$, il existe $L_0 \in \mathbb{N}$ (défini comme le plus petit L tel que $f^{(L)}(x) \leq c$)
 \hookrightarrow Existe car $\begin{cases} f(x) < x \\ f(0) = 0 \leq c \end{cases}$

$$\begin{aligned} \text{On écrit donc } \lim_{L \rightarrow \infty} \frac{f^{(L)}(x)}{a^L} &= \lim_{L \rightarrow \infty} \frac{f^{(L-L_0)}(f^{(L_0)}(x))}{a^{L-L_0}} a^{L_0} \\ &= \frac{f^{(L_0)}(x)}{a^{L_0}}, \text{ car } f^{(L_0)}(x) \leq c \text{ d'après la question précédente} \end{aligned}$$

Un calcul simple donne

$$\frac{1}{a^{L_0}} \in \left[\left(\frac{1-c}{1-x} \right)^{-\frac{\log a}{\log b}}, \frac{1}{a} \left(\frac{1-c}{1-x} \right)^{-\frac{\log a}{\log b}} \right]$$

$$\text{On obtient donc } \lim_{L \rightarrow \infty} \frac{f^{(L)}(x)}{a^L} = \frac{R(x)}{(1-x)^{\frac{\log a}{\log b}}},$$

$$\text{avec } R(x) := \frac{f^{(L_0)}(x) (1-x)^{-\frac{\log a}{\log b}}}{a^{L_0}}, \text{ où } \begin{cases} f^{(L_0)}(x) \in [ac, c] \\ \frac{(1-x)^{-\frac{\log a}{\log b}}}{a^{L_0}} \in \left[(1-c)^{-\frac{\log a}{\log b}}, \frac{1}{a} \right] \end{cases}$$

Pour conclure sur $K^{(L)}(\gamma) = \sum_{i=0}^L \psi^{(i)}(\gamma) \prod_{k=i}^{L-1} \psi'(\psi^{(k)}(\gamma))$, les auteurs (19)

procèdent par majoration/minoration

↳ Existence des limites laissées sous le tapis

* Majoration: long & technique

* Minoration: On utilise les propriétés de positivité de ψ .

$$\boxed{1} \quad \forall \gamma \in [0, 1], \quad \psi'(\gamma) \gamma \geq \psi(\gamma)$$

En effet, $\psi(\gamma) = \sum_{j=1}^{\infty} a_j \gamma^j$ car $\psi(0) = 0$

$$\Rightarrow \psi'(\gamma) \gamma = \gamma \sum_{j=1}^{\infty} a_j \times j \gamma^{j-1} = \sum_{j=1}^{\infty} \underbrace{j}_{\geq 1} \underbrace{a_j \gamma^j}_{\geq 0} \geq \psi(\gamma) \quad \square$$

$\boxed{2}$ On écrit

$$K^{(L)}(\gamma) = \sum_{i=0}^L \psi^{(i)}(\gamma) \prod_{k=i}^{L-1} \psi'(\underbrace{\psi^{(k)}(\gamma)}_{=\tilde{\gamma}_k})$$

$$\geq \sum_{i=0}^L \psi^{(i)}(\gamma) \prod_{k=i}^{L-1} \frac{\psi(\psi^{(k)}(\gamma))}{\psi^{(k)}(\gamma)}$$

$$= \sum_{i=0}^L \psi^{(L)}(\gamma)$$

$$= (L+1) \psi^{(L)}(\gamma)$$

Ainsi, $\lim_{L \rightarrow \infty} \frac{K^{(L)}(\gamma)}{\psi'(0)^{L(L+1)}} \geq \lim_{L \rightarrow \infty} \frac{\psi^{(L)}(\gamma)}{\psi'(0)} = \frac{\tilde{R}(\gamma)}{(1-\gamma)^{\frac{\alpha}{2}}}$

IV Quelques questions ouvertes

(20)

- Peut-on obtenir des résultats similaires avec des structures de réseaux de neurones différentes?

↳ Convolutional NN?

- La partie traitant de limite $L \rightarrow \infty$ est très artisanale et ne fait pas apparaître de structure à la limite.

↳ Formaliser la limite de noyaux généraux

([Aronszajn 1950] le fait pour des noyaux \nearrow ou \searrow)

↳ Donner une structure type RKHS pour des "noyaux" tels que

$$K(x, x) = +\infty.$$

$\Delta \int_{\mathbb{R}^d} \frac{1}{|x|} dx$ pas RKHS (Partir de Fourier pour $K(x, \tilde{x}) = g(\|x - \tilde{x}\|)$ avec $\hat{g} \geq 0$ au sens des distributions?)

↳ Comprendre plus finement le reste $R(\|x - \tilde{x}\|)$ au numérateur