Big Data

François Malgouyres

Institut de Mathématiques de Toulouse, Université Paul Sabatier

2025-2026

Plan

- - Introduction

 Classification et régression
 - Les réseaux de neurones

2025-2026



Classification et régression

- Les réseaux de neurones
 - Multi-Layer Perceptron/fully-connected neural network
 - Réseaux convolutifs
 - Réseaux Récurrents
 - Mécanismes d'attention et transformers
 - Beaucoup de variantes importantes
- L'optimisation des réseaux de neurones
 - L'optimisation des réseaux de neurones: Premières propriété
- La Rétropropagation et ses défauts
- Le paysage de la fonction objectif
 - Introduction
 - Paysage pour les réseaux larges
- La régularisation induite par la géométrie
 - Overview sur un exemple
 - Résultats formels
- Le paysage pour les réseaux linéaires
 - Introduction
 - Paysage pour les réseaux linéaires

Classification et régression

Il existe un couple de variables aléatoires (X, Y)

- On connaît \mathscr{X} et \mathscr{Y} tels que $\mathbb{P}(X \in \mathscr{X}) = \mathbb{P}(Y \in \mathscr{Y}) = 1$
- On sait que l'on va observer x d'une réalisation (x, y) de (X, Y)
- On veut construire une fonction $g: \mathcal{X} \longrightarrow \mathcal{Y}$ prédisant y

Plus précisément, pour une fonction coût $L: \mathcal{Y} \times \mathcal{Y} \longrightarrow \mathbb{R}$, on veut trouver

$$g^* \in \operatorname{argmin} R(g)$$

On appelle $R(g) = \mathbb{E}(L(g(X), Y))$ le risque de g.

Car: La personnes utilisant g observe le coût

$$R_{test}(g) \simeq \frac{1}{n'} \sum_{i=1}^{n'} L(g(x_i'), y_i')$$

pour un échantillon i.i.d $(x'_i, y'_i)_{i=1..n'}$ suivant la loi (X, Y), **indépendant de l'échantillon d'apprentissage**. La loi des grands nombres conduit à considérer R.

Classification

En classification, \mathscr{Y} est fini.

- Classification binaire : $\mathcal{Y} = \{-1, +1\}$
- Classification multi-classes : $\mathscr{Y} = \{1, \dots, C\}, C \in \mathbb{N}, C > 2$

En deep-learning, pour $C \ge 2$, on construit C fonctions

$$g_c: \mathscr{X} \longrightarrow \mathbb{R}$$
 , $c = 1..C$

on retourne

$$y \in \operatorname{argmax}_{c=1..C} g_c(x)$$

$$L(y, y') = \mathbb{1}_{y \neq y'}$$
 ou la perte logistique, etc

Exemples:

- $\mathscr{X} = \mathbb{R}^2$, $\mathscr{Y} = \{-1, +1\}$
- $\mathscr{X} = \mathbb{R}^n$: n est le nombre de pixels d'une image, $\mathscr{Y} = \{1, \dots, C\}$: chaque sortie est une "note" pour la présence d'un objet.
- $\mathcal{X} = \mathbb{R}^n$: n est la taille d'un "état" (ex: l'écran d'un jeu Atari), $\mathcal{Y} = \{1, \dots, C\}$: chaque sortie est l'espérance de gain de l'action c = 1...C.

Régression

En régression,

- \bullet $\mathscr{Y} = \mathbb{R}$
- $\mathcal{Y} = \mathbb{R}^C$

On construit C fonctions

$$g_c: \mathscr{X} \longrightarrow \mathbb{R}$$
 , $c = 1..C$

$$L(y, y') = ||y - y'||^2$$
 (Son choix est une composante de la modélisation)

Exemples:

- $\mathscr{X} = \mathbb{R}, \mathscr{Y} = \mathbb{R}$
- $\mathcal{X} = \mathbb{R}^d$: d est le nombre de mesures concernant une situation économique actuelle, $\mathcal{Y} = \mathbb{R}^C$: différentes quantités économiques futures c = 1..C.
- $\mathscr{X} = \mathbb{R}^d$: d la taille d'un vocabulaire, (0, ..., 0, 1, 0, ..., 0) est un mot $\mathscr{Y} = \mathbb{R}^C$: espace d'une représentation des mots du vocabulaire.

Classification et régression: Erreur Bayésienne

Décision Bayésienne

Sous des hypothèses faibles sur la loi de (X, Y) et si elle existe, la fonction définie pour $x \in \mathcal{X}$ par

$$g_b(x) \in \operatorname{argmin}_{y \in \mathscr{Y}} \mathbb{E}(L(y, Y) | X = x)$$

minimise le risque :

$$\forall g: \mathscr{X} \to \mathscr{Y}, \qquad R(g) \geq R(g_b)$$

On appelle g_b la **décision Bayésienne**.

Remarques:

On n'a pas forcément

$$\forall x \in \mathcal{X}, \exists a \in \mathcal{Y}, \quad \mathbb{P}(Y = a | X = x) = 1.$$

En général,

$$R(g_b) \neq 0$$
.

- On ne peut pas calculer g_b car:
 - On ne connaît pas la loi de (X, Y)
 - On devrait représenter g_b avec un ordinateur

Classification et régression : les principes

- On observe $(x_i, y_i)_{i=1..n}$ un échantillon i.i.d. de même loi que (X, Y)
- On considère une famille \mathscr{F} de fonctions g_w paramétrés par des paramètres w d'un espace Euclidien.
 - Pour nous : elle est basée sur les réseaux de neurones
- On voudrait résoudre

$$\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w}} R(g_{\mathbf{w}}).$$

• La minimisation du risque empirique consiste à résoudre

$$\widehat{\mathbf{w}} \in \operatorname{argmin}_{\mathbf{w}} \widehat{R}(g_{\mathbf{w}})$$

οù

$$\widehat{R}(g_{\mathbf{w}}) = \frac{1}{n} \sum_{i=1}^{n} L(g_{\mathbf{w}}(x_i), y_i).$$

Le contrôle du risque

ullet Pour toute fonction g, on appelle **risque en population** et **risque empirique** de g:

$$R(g) = \mathbb{E}(L(g(X), Y))$$
 et $\widehat{R}(g) = \frac{1}{n} \sum_{i=1}^{n} L(g(x_i), y_i)$

- $R^* = \inf_g R(g)$ le risque optimal
- Pour $\varepsilon > 0$, on fixe \mathbf{w}^* et $\widehat{\mathbf{w}}$ tels que :

$$R(g_{\mathbf{W}^*}) \le \inf_{\mathbf{W}} R(g_{\mathbf{W}}) + \varepsilon$$
 et $\widehat{R}(g_{\widehat{\mathbf{W}}}) \le \inf_{\mathbf{W}} \widehat{R}(g_{\mathbf{W}}) + \varepsilon$

Pour w retourné par un algorithme

On décompose

$$\begin{array}{lll} 0 \leq & R(g_{\mathbf{W}}) & - & R^* & \text{(l'excès de risque)} \\ & = & R(g_{\mathbf{W}}) & - & \widehat{R}(g_{\mathbf{W}}) & \text{(erreur de généralisation)} \\ & + & \widehat{R}(g_{\mathbf{W}}) & - & \widehat{R}(g_{\widehat{\mathbf{W}}}) & \text{(erreur d'optimisation)} \\ & + & \widehat{R}(g_{\widehat{\mathbf{W}}}) & - & \widehat{R}(g_{\mathbf{W}^*}) & \leq \varepsilon \\ & + & \widehat{R}(g_{\mathbf{W}^*}) & - & R(g_{\mathbf{W}^*}) & \text{(erreur de généralisation)} \\ & + & R(g_{\mathbf{W}^*}) & - & R^* & \text{(erreur d'approximation)} \end{array}$$

Puis on **majore** indépendamment chaque terme.

Le contrôle du risque

- Erreur d'optimisation : $|\widehat{R}(g_{\mathbf{w}}) \widehat{R}(g_{\widehat{\mathbf{w}}})|$
 - Est-ce-que l'optimisation a réussi? Pbs: Paysage de la fonction objectif, n grand, d grand.
 - ► En deep-learning : Elle est petite pour les "gros" réseaux.
- Erreur d'approximation : $|R(g_{\mathbf{w}^*}) R(g_b)|$, où $R(g_b) \simeq \inf_g R(g)$
 - ► Il suffit que $||g_{\mathbf{w}^*} g_b||$ soit petit (le choix de la norme compte).
 - Quelles fonctions peut-on approximer avec notre classe de fonctions? Mots clefs: Expressivité, théorie de l'approximation etc
 - ► En deep-learning : Elle est petite pour les "gros" réseaux.
- Erreur de généralisation : $|R(g_w) \widehat{R}(g_w)|$
 - Quelles conditions sur le réseau et l'échantillon pour avoir la convergence uniforme

$$\sup_{\mathbf{w}} |\widehat{R}(g_{\mathbf{w}}) - R(g_{\mathbf{w}})| \qquad \text{soit petit?}$$

- ► En théorie : Dimension de Vapnik-Chervonenkis ~ nombre de paramètres
- ► En pratique : Reste faible même pour les "gros" réseaux

~ régularisation implicite.

La régularisation implicite : Visualisation

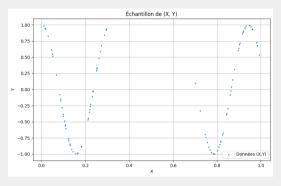


Figure: Source : demo_regularisation_implicite.ipynb

Apprentissage d'un réseau de neurones MLP, ReLU, de profondeur 6, de largeur 40. Que va-t'il se passer dans [0.3,0.7] ?

La régularisation implicite : Visualisation

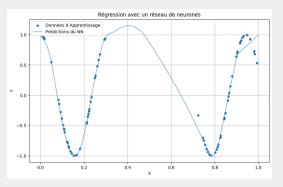


Figure: Source : demo_regularisation_implicite.ipynb

La régularisation implicite : pas de sur-apprentissage

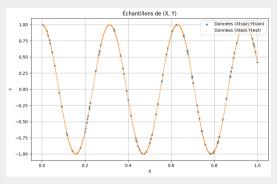


Figure: Source: demo regularisation implicite overfit.ipynb

L'échantillon d'apprentissage comprend 100 exemples.

Quelles sont les performances des réseaux de profondeur fixe (6 couches) pour différentes largeurs ?

largeur	10	20	30	40	50	60	70	80	90	100	200	400
nb. param.	581	2161	4741	8321	12901	18481	25061	32641	41221	50801	201601	803201

La régularisation implicite : pas de sur-apprentissage

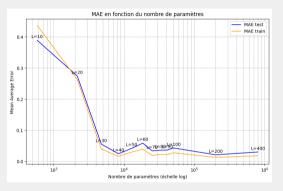


Figure: Source: demo regularisation implicite overfit.ipynb

Même très surparamétré, dans cette expérience, le réseau de neurones ne surapprend pas. Warning: Cela arrive tout de même.

Le contrôle du risque: En pratique

Pour w retourné par un algorithme

$$R(g_{\mathbf{w}}) = R(g_{\mathbf{w}}) - \widehat{R}(g_{\mathbf{w}})$$
 (erreur de généralisation)
+ $\widehat{R}(g_{\mathbf{w}})$ observable

• N'a d'intérêt que si on est capable de trouver un réseau et un w tels que

$$\widehat{R}(g_{\mathbf{w}})$$
 soit petit.

- ▶ i.e. : On a résolu les problèmes de l'expressivité et de l'optimisation.
- Erreur de généralisation : même problématique que précédemment.

Classification et régression : les principes

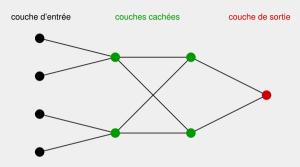
ATTENTION: Cette modélisation ne tient pas compte:

- Problèmes plus complexes : apprentissage non-supervisé, semi-supervisé, par renforcement, apprentissage de représentations ...
- du **biais** dans les données: $(x_i, y_i)_{i=1..n}$ est rarement i.i.d. selon (X, Y)
 - Les x_i peuvent ne couvrir que partiellement le domaine, être corrompus, etc
 - Les *y_i* peuvent être une décision biaisée, bruitée, etc
- du besoin de robustesse
 - ▶ On impose que y_i soit prédit pour tout $x_i + e$, où $||e|| \le \varepsilon$
- d'une régularisation du réseau
 - w est parfois quantifié
 interpretabilité, faible coût énergétique/calculatoire/mémoire
 - w est tel que g_w est 1-Lipschitz \Longrightarrow garantie de robustesse
 - Régularisations explicites : poids parcimonieux, pénalisation quadratique, dropout



- Classification et régression
- Les réseaux de neurones
 - Multi-Layer Perceptron/fully-connected neural network
 - Réseaux convolutifs
 - Réseaux Récurrents
 - Mécanismes d'attention et transformers
 - Beaucoup de variantes importantes
- L'optimisation des réseaux de neurones
- L'optimisation des réseaux de neurones: Premières propriété
- La Rétropropagation et ses défauts
- Le paysage de la fonction objectif
 - Introduction
 - Paysage pour les réseaux larges
- La régularisation induite par la géométrie
 - Overview sur un exemple
 - Résultats formels
- Le paysage pour les réseaux linéaires
 - Introduction
 - Paysage pour les réseaux linéaires

Multi-Layer Perceptron/fully-connected neural network



- Couche 0: couche d'entrée;
 Couches 1 et 2: couches cachées;
 Couche 3: couche de sortie.
- Chaque sommet (= neurone) contient un réel
- Chaque arête contient un poids

Multi-Layer Perceptron/fully-connected neural network

- → H couches (on dit aussi H 1 couches cachées)
- les tailles $n_0, n_1, ..., n_H \in \mathbb{N}$ des différentes couches
- On note $f_h(x)$: le résultat obtenu en calculant le contenu de la couche h pour l'entrée $x \in \mathbb{R}^{n_0}$
- On note $W_h \in \mathbb{R}^{n_h \times n_{h-1}}$: la matrice contenant les poids sur les arcs entre la couche h-1 et la couche h
- On note $b_h \in \mathbb{R}^{n_h}$: le biais ajouté à la couche h
- On note σ_h : la fonction d'activation appliquée à chaque couche
 - ► Typiquement, elle applique la même fonction à chaque entrée d'un vecteur

La prédiction/l'inférence : la fonction $f_H : \mathbb{R}^{n_0} \longrightarrow \mathbb{R}^{n_H}$

$$\begin{cases}
f_0(x) = x \\
f_h(x) = \sigma_h(W_h f_{h-1}(x) + b_h)
\end{cases}, \forall h = 1, \dots, H$$

Multi-Layer Perceptron/fully-connected neural network : Exemple

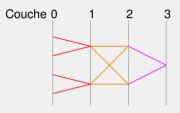


Figure: Réseau Fully-connected

Ci-dessous, chaque * est un nombre

$$W_{1} = \begin{pmatrix} * & * & 0 & 0 \\ 0 & 0 & * & * \end{pmatrix}, b_{1} = \begin{pmatrix} * \\ * \end{pmatrix}, W_{2} = \begin{pmatrix} * & * \\ * & * \end{pmatrix}, b_{2} = \begin{pmatrix} * \\ * \end{pmatrix}, W_{3} = (* * *), b_{3} = (* *),$$

$$f_{3}(x) = \sigma_{3} \left(W_{3} \sigma_{2} \left(W_{2} \sigma_{1} \left(W_{1} x + b_{1} \right) + b_{2} \right) + b_{3} \right)$$

F. Malgouyres Big Data 2025-2026 20/104

Multi-Layer Perceptron/fully-connected neural network: Vocabulaire et variantes

- Fully-connected layer : La matrice est "pleine"
- Simplifications pour l'étude théorique :
 - Fully connected **linear** network: $\sigma(t) = t$, $\forall t \in \mathbb{R}$, et $b_h = 0$, pour tout h = 1..H.
 - ► One hidden-layer Neural Network: on prend H = 2
- Fonctions d'activation :
 - Une zoologie importante : tanh, heavyside, identité, sigmoïde, etc
 - ► Des non-locales : group-sort, max pooling
 - La plus utilisée est "Rectified Linear Unit" (ReLU) : $\sigma(t) = \max(0, t)$
 - ► Souvent, celle associée à la dernière couche est l'identité ou softmax (pour la classification).
- "Batch-normalisation": Au lieu d'optimiser les biais \mathbf{b} , on règle un couple $(\operatorname{diag}(\gamma), \mathbf{b})$ permettant de centrer les données et fixer leur variance.
- "Dropout" : On minimise le risque moyen pour des sous-réseaux aléatoires.
- L'architecture du réseau: La donnée des hyperparamètres. Ex : Un réseau ReLU de largeur 100 et de profondeur 10.

Multi-Layer Perceptron/fully-connected neural network

Proposition : Propriétés de f_H , cas linéaire avec biais

Pour tout réseaux linéaire fully-connected, f_H est affine:

$$f_H(x) = W_H \cdots W_1 \ x + b'_H$$

οù

$$b'_{H} = W_{H} \cdots W_{2} b_{1} + \cdots + W_{H} b_{H-1} + b_{H}$$

Preuve:

- Une composition de fonctions affines est affine.
- La formule est facile établir avec une récurrence sur *H*. (Exercice)

Les réseaux de neurones

Proposition : Propriétés de f_H, cas ReLU

Pour tout réseau de neurones, avec la fonction d'activation ReLU et l'identité sur la dernière couche

- \bullet La fonction f_H est continue
- Il existe une partition compatible ayant moins de $2^{n_1+\cdots+n_{H-1}}$ morceaux dont les morceaux sont des polyèdres ayant au plus $n_1+\cdots+n_{H-1}$ faces.
- Réciproque dans : Arora, Raman, et al. "Understanding Deep Neural Networks with Rectified Linear Units." ICLR, 2018:
 - Toute fonction continue, affine sur un nombre fini de polyhèdres (dont l'union est \mathbb{R}^{n_0}) peut être représenté par un réseau de neurone ReLU.

Réseaux convolutifs

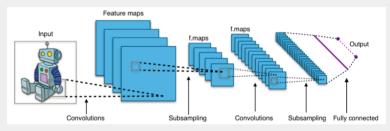
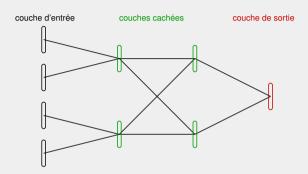


Figure: Source: wikipedia page "apprentissage profond"

Réseaux convolutifs



- Chaque sommet contient un signal/une image. On parle de canal.
- Les canaux de la couche d'entrée correspondent aux canaux (ex (r, g, b)) de x.

Réseaux convolutifs

- des signaux de taille N
- → H couches (on dit aussi H 1 couches cachées)
- les tailles $n_0, n_1, \dots, n_H \in \mathbb{N}$ des différentes couches (= $\mathbb{N} \times$ nombre de canaux)
- On note $f_h(x)$: le résultat obtenu en calculant le contenu de la couche h pour l'entrée x
- On note $W_h \in \mathbb{R}^{n_h \times n_{h-1}}$: la matrice **concaténant des matrices de convolutions** pour passer de la couche h-1 et la couche h
- On note $b_h \in \mathbb{R}^{n_h}$: le biais ajouté à la couche h
- On note σ_h : la fonction d'activation appliquée à chaque couche
 - ► Elle applique la même fonction à chaque entrée d'un vecteur

L'action du réseau : la fonction f_H

$$\begin{cases}
f_0(x) = vect(x) \\
f_h(x) = \sigma_h(W_h f_{h-1}(x) + b_h)
\end{cases}, \forall h = 1, \dots, H$$

Réseaux convolutifs: Matrice circulante, Toeplitz, bloc circulante... calculant des convolutions

On suppose $v, x \in \mathbb{R}^N$. Le signal v est supposé (N)-périodique :

$$v_{-n'} = v_{N-n'}$$
 pour tout $n' = 0..N - 1$

On a:

$$\begin{pmatrix} v_0 & v_{N-1} & \cdots & \cdots & v_1 \\ v_1 & v_0 & v_{N-1} & \cdots & \cdots & v_2 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \vdots \\ \vdots & & & \ddots & \ddots & \ddots & \vdots \\ \vdots & & & & \ddots & \ddots & v_{N-1} \\ v_{N-1} & \cdots & \cdots & v_1 & v_0 \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ \vdots \\ \vdots \\ x_{N-1} \end{pmatrix} = \begin{pmatrix} \sum_{\substack{n'=1 \\ n'=0}}^{N-1} v_{0-n'} x_{n'} \\ \sum_{\substack{n'=0 \\ n'=0}}^{N-1} v_{1-n'} x_{n'} \\ \vdots \\ \sum_{\substack{n'=1 \\ n'=0}}^{N-1} v_{1-n'} x_{n'} \end{pmatrix} = v * x$$

Réseaux convolutifs: Exemple

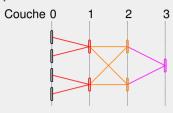


Figure: Réseau convolutif

Ci-dessous,

- chaque * est une matrice composant une convolution et un échantillonnage
- chaque est un vecteur colonne de taille *N* (cas sans sous-échantillonnage)

$$W_{1} = \begin{pmatrix} * & * & 0 & 0 \\ 0 & 0 & * & * \end{pmatrix}, b_{1} = \begin{pmatrix} \circ \\ \circ \end{pmatrix}, W_{2} = \begin{pmatrix} * & * \\ * & * \end{pmatrix}, b_{2} = \begin{pmatrix} \circ \\ \circ \end{pmatrix}, W_{3} = \begin{pmatrix} * & * \end{pmatrix}, b_{3} = \begin{pmatrix} \circ \\ \circ \end{pmatrix}, f_{3}(x) = \sigma \Big(W_{3} \sigma \big(W_{2} \sigma \big(W_{1} x + b_{1} \big) + b_{2} \big) + b_{3} \Big)$$

Réseaux convolutifs: Propriétés

Remarques

- Les formules sont les mêmes que dans le cas fully-connected. Seules les façons de construire les W_h et $f_0(x)$ diffèrent.
- Un réseau convolutif contient des couches convolutives et des couches fully-connected.
- Un réseau convolutif est un réseau fully-connected dans lequel on impose une structure.

Proposition : Propriétés de f_H , cas linéaire

Pour tout réseaux convolutif linéaire, f_H est affine:

$$f_H(x) = W_H \cdots W_1 \ x + b'_H$$

οù

$$b'_{H} = W_{H} \cdots W_{2} b_{1} + \cdots + W_{H} b_{H-1} + b_{H}$$

De plus, $x \mapsto W_H \cdots W_1 x$ est une convolution du signal de départ. (Preuve en exercice)

Réseaux Récurrents

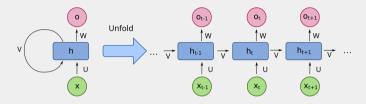


Figure: Source: wikipedia page "recurrent neural networks"

- Modèles de langage : l'entrée $x = (x_i)_{i=1..T}$ est un texte, $x_i = (0,...,0,1,0,...,0)^T \in \mathbb{R}^{d_{in}}$ est un mot encodé "one-hot" dans un dictionnaire de d_{in} mots. Le RNN prédit le prochain mot.
- Signal : les entrées $x = (x_i)_{i=1..T}$ sont les mesures de capteurs, t est le temps. Le RNN indique si le système fonctionne correctement, ou comment le maintenir dans une configuration, etc

Réseaux Récurrents

- Entrée : La série temporelle $x_1, ..., x_T \in \mathbb{R}^{d_{in}}$, pour $d_{in} \in \mathbb{N}$
- État caché : La série temporelle $h_0, \ldots, h_T \in \mathbb{R}^{d_h}$, pour un paramètre $d_h \in \mathbb{N}$, avec $h_0 = 0$ et

$$h_t = \sigma_h \Big(V h_{t-1} + U x_t + b_i \Big)$$
 pour tout $t = 1..T$

pour

- ▶ La matrice de récurrence : $V \in \mathbb{R}^{d_h \times d_h}$
- La matrice d'entrée : $U \in \mathbb{R}^{d_h \times d_{in}}$
- ▶ Le biais : $b_i \in \mathbb{R}^{d_h}$
- l'activation σ_b (souvent l'identité)
- Sortie :
 - ► Many-to-one : un vecteur

$$V\sigma(h_T) + b_O \in \mathbb{R}^{d_{out}}$$

Many-to-many : une série temporelle d'éléments de $\mathbb{R}^{d_{out}}$

$$V\sigma(h_1)+b_0,\cdots,V\sigma(h_T)+b_0$$

31/104

pour $V \in \mathbb{R}^{d_{out} \times d_h}$, $b_o \in \mathbb{R}^{d_{out}}$ et la fonction d'activation σ .

Réseaux Récurrents

Remarques

- On peut déplier les formules du RNN pour obtenir les mêmes même formules que dans le cas fully-connected (exercice). Seules les façons de construire les W_h et $f_0(x)$ diffèrent. Il y a des poids partagés entre les couches.
- Un réseau récurrent est un réseau fully-connected dans lequel on impose une structure.
- Difficultés: La profondeur du réseaux déplié est T. Quand T est grand (ex: T = 2000), cela pose des problèmes de vanishing/exploding gradient (voir chapitre 2).
 On parle de bien gérer les "dépendances à long terme", la "mémorisation".
- Principales architectures concurrentes :
 - Long short-term memory (LSTM), 1997, 114000+ citations google scholar.
 - ► Gated Recurrent Units (GRU), 2014, 31000+ citations google scholar.
 - ► State-Space Models (SSM), 2021, 1000 + citations google scholar.

Transformers: Self-attention

- Entrée : La série temporelle $x_1, ..., x_T \in \mathbb{R}^{d_{in}}$, pour $d_{in} \in \mathbb{N}$, **en lignes**
- On calcule

$$X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_T \end{pmatrix} \in \mathbb{R}^{T \times d_{in}} \quad \text{et} \quad \begin{cases} Q = XW^q \in \mathbb{R}^{T \times d_h} \\ K = XW^k \in \mathbb{R}^{T \times d_h} \\ V = XW^v \in \mathbb{R}^{T \times d_h} \end{cases}$$

puis le plongement basé sur l'attention vaut

$$A = BV \in \mathbb{R}^{T \times d_h}$$
,

pour le mécanisme d'attention

$$B = Soft_max(\frac{1}{\sqrt{d_{in}}}QK^T) \in \mathbb{R}^{T \times T}$$

où, pour tout
$$i, j$$
, $Soft_max(B)_{i,j} = \frac{\exp(B_{i,j})}{\sum_{j'=1}^{T} \exp(B_{i,j'})}$.

- Pour $d_h \in \mathbb{N}$ et les paramètres
 - ▶ Une matrice de "query" : $W^q \in \mathbb{R}^{d_{in} \times d_h}$
 - ▶ Une matrice de "key" : $W^k \in \mathbb{R}^{d_{in} \times d_h}$
 - ▶ Une matrice de "value" : $W^{V} \in \mathbb{R}^{d_{in} \times d_{h}}$

F. Malgouyres

Big Data

2025-2026

Transformers: Multi-headed self-attention

Pour $h \in \mathbb{N}$ (ex : h=8), on a h mécanismes d'attention :

$$(W_i^q, W_i^k, W_i^v)_{i=1..h} \in (\mathbb{R}^{d_{in} \times d_h} \times \mathbb{R}^{d_{in} \times d_h} \times \mathbb{R}^{d_{in} \times d_h})^h,$$

pour $d_h = \frac{d_{in}}{h}$. On concatène les résultats

$$A = (A_1, A_2, \ldots, A_h) \in \mathbb{R}^{T \times d_{in}}.$$

Le mécanisme "Multi-Headed self-attention" est utilisé dans des transformers, GPT, BERT, BART, etc.

Les transformers

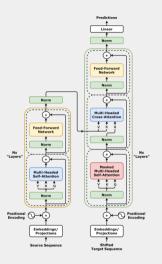


Figure: Source: wikipedia page "Transformer (deep learning architecture)"

Les transformers : Problème

Entreprise	Entreprise Nom		Nombre de	temps	coût	
			paramètres	d'entraînement	d'entraînement	
				(en jours)	(en dollars)	
OpenAl	GPT 1	2018	117 10 ⁶	quelques	< 10 000	
OpenAl	GPT 2	2019	1,5 10 ⁹	7-14	50.10 ³	
OpenAl	GPT 3	2020	175 10 ⁹	en mois	en millions	
MetaAl	BART	2019	0,4 10 ⁹	7-14	50.10 ³	

Beaucoup de variantes importantes

Des architectures pour éviter le "vanishing/exploding gradient": ex : Couches "Res-net"

$$f_h(x) = \sigma \Big(W_h \, \sigma \Big(W_{h-1} f_{h-2}(x) + b_{h-1} \Big) + b_h + f_{h-2}(x) \Big)$$

- Des architectures plus évoluées, dédiés à des applications.
 - ► En vision : U-Net (segmentation), YOLO (detection), transformer.
 - Séries temporelles et modèles de language : RNN, LSTM, GRU, SSM, transformer.
 - Robotique
 - Graphes
 - ► Etc
- etc

Modélisation avec des réseaux

- Autoencoder et VAE : Un réseau compresse, un réseau décompresse.
- Création d'embeddings : Deux réseaux envoient une image et un texte dans un même espace de caractéristiques. (Ex : Visual Query Answering (VQA))
- Generative Adversarial Networks (GAN): un réseau génère des données; un réseau discrimine les données générées de vraies données.
- Réduction de biais dans les données (FairGAN): Un réseaux envoie les données dans un espace de caractéristiques et ses poids sont optimisés pour qu'un autre réseau classifiant sur un critère pertinent fonctionne, un troisième réseau classifiant sur un critère non-pertinent échoue.
- etc

Plan

- Introductio
- 2 L'optimisation des réseaux de neurones
 - L'optimisation des réseaux de neurones: Premières propriétés
 - La Rétropropagation et ses défauts
 - Le paysage de la fonction objectif
 - La régularisation induite par la géométrie
 - Le paysage pour les réseaux linéaires



- Classification et régression
- Les réseaux de neurones
 - Multi-Layer Perceptron/fully-connected neural network
 - Réseaux convolutifs
 - Réseaux Récurrents
 - Mécanismes d'attention et transformers
 - Beaucoup de variantes importantes
- L'optimisation des réseaux de neurones
 - L'optimisation des réseaux de neurones: Premières propriétés
 - La Rétropropagation et ses défauts
- Le paysage de la fonction objectif
 - Introduction
 - Paysage pour les réseaux larges
- La régularisation induite par la géométrie
 - Overview sur un exemple
 - Résultats formels
- Le paysage pour les réseaux linéaires
 - Introduction
 - Paysage pour les réseaux linéaires

L'optimisation : Premières propriétés

- On dispose d'un échantillon $(x_i, y_i)_{i=1..n}$.
- On dispose d'un réseau de profondeur *H*, d'architecture fixée,
 - de paramètre

$$\theta = (W_H, \dots, W_1, b_H, \dots, b_1) \in \Theta$$

les fonctions d'activation sont notées σ_h

La prédiction est la fonction

$$f_{\theta}(x) = \sigma_{H} \Big(W_{H} \cdots \sigma_{2} \big(W_{2} \sigma_{1} \big(W_{1} x + b_{1} \big) + b_{2} \big) \cdots + b_{H} \Big)$$

- Les réseaux ReLU:
 - $\sigma_h = ReLU$, pour $h = 1, \dots, H-1$
 - Cas de la régression : σ_H = Id
 - Cas de la classification : σμ est softmax.
- On considère une fonction coût définie par

$$E: \Theta \longrightarrow \mathbb{R}$$

$$\theta \longmapsto E(\theta) = \sum_{i=1}^{n} L(f_{\theta}(x_i), y_i)$$

pour une fonction coût $L: \mathcal{Y} \times \mathcal{Y} \longrightarrow \mathbb{R}$.

L'optimisation : Premières propriétés

Proposition: Non-coercivité dans le cas ReLU

Pour tout réseau ReLU, pour tout échantillon d'apprentissage, la fonction E n'est pas coercive.

Preuve utilisant l'homogénéité:

On considère θ avec des biais nuls et, pour tout $\lambda > 0$, on définit θ^{λ} par

$$W_1^{\lambda} = \lambda^{H-1} W_1$$
 et $W_h^{\lambda} = \lambda^{-1} W_h, \forall h = 2, ..., H$

et des biais nuls.

Comme pour tout $\lambda > 0$, $\sigma(\lambda t) = \lambda \sigma(t)$, on a pour tout x et tout $\lambda > 0$,

$$f_{\theta^{\lambda}}(x) = \lambda^{-1} W_{H} \cdots \sigma(\lambda^{-1} W_{2} \sigma(\lambda^{H-1} W_{1} x)) \cdots$$
$$= f_{\theta}(x)$$

Donc $E(\theta) = E(\theta^{\lambda})$, pour tout $\lambda > 0$, mais $\lim_{\lambda \to +\infty} \|\theta^{\lambda}\| = +\infty$.

Donc E n'est pas coercive.

F. Malgouyres Big Data 2025-2026 42/104

L'optimisation : Premières propriétés

Proposition: Non-convexité

Pour la loss quadratique $L(y',y)=(y'-y)^2$, le réseaux ReLU fully connnected d'architecture (1,1,1) et l'échantillon constitué de l'exemple (x,y)=(1,0), le risque empirique est non convexe.

Preuve: On a pour tout $(w_1, w_2) \in \mathbb{R}^2$ avec $w_1 \ge 0$, et pour $(b_1, b_2) = (0, 0)$

$$E(w_1, w_2, b_1, b_2) = (w_2 \sigma(w_1 x + b_1) + b_2 - y)^2 = (w_1 w_2)^2.$$

Cette fonction n'est pas convexe car:

- En $(w_1, w_2) = (0, 1), E(0, 1, 0, 0) = 0$
- En $(w_1, w_2) = (1,0), E(1,0,0,0) = 0$
- En $(w_1, w_2) = (0.5, 0.5) = 0.5 \times (1, 0) + 0.5 \times (0, 1),$

$$E(0.5, 0.5, 0, 0) = 0.25^2 > 0 = 0.5 \times E(1, 0, 0, 0) + 0.5 \times E(0, 1, 0, 0).$$

F. Malgouyres Big Data 2025-2026 43/104



- Introduction
- Classification et régression
- Les réseaux de neurones
 - Multi-Layer Perceptron/fully-connected neural network
 - Réseaux convolutifs
 - Réseaux Récurrents
 - Mécanismes d'attention et transformers
 - Beaucoup de variantes importantes
- L'optimisation des réseaux de neurones
 - L'optimisation des réseaux de neurones: Premières propriété
 - La Rétropropagation et ses défauts
- Le paysage de la fonction objectif
 - Introduction
 - Paysage pour les réseaux larges
 - La régularisation induite par la géométrie
 - Overview sur un exemple
 - Résultats formels
- Le paysage pour les réseaux linéaires
 - Introduction
 - Paysage pour les réseaux linéaires

La Rétropropagation

On suppose:

La fonction coût est de la forme :

$$\theta \longmapsto E(\theta) = \sum_{i=1}^{n} L(f_{\theta}(x_i), y_i)$$

où L est C^1 .

- On suppose que, pour tout $h \in \{1, ..., H\}$, σ_h est C^1
 - On peut régulariser σ_h
 - ► Il existe un cadre formel pour faire du calcul différentiel et de l'optimisation avec des activations comme ReLU
- Sous ces hypothèses: E est différentiable.
- Les frameworks de deep learning utilisent la différentiation automatique.
- Ci-dessous, on ne considère qu'un unique exemple (x, y):
 - ► On somme si besoin plusieurs gradients $\nabla L(f_{\theta}(x_i), y_i)$

La Rétropropagation

- \mathbf{W}_h la matrice pour passer de la couche h-1 à h,
- $f_{\theta}^{h}(x)$, le contenu de la couche h,
- $\bullet \ d_{\theta}^{h}(x) = \sigma_{h}'(\mathbf{W}_{h}f_{\theta}^{h-1}(x) + \mathbf{b}_{h}).$

Proposition: Gradient pour un réseau fully-connected

On a pour h = 1..H, $i = 1..n_h$, $j = 1..n_{h-1}$

$$\frac{\partial E}{\partial (\mathbf{W}_h)_{i,j}}(\theta) = \left[f_{\theta}^{h-1}(x) \right]_j \left[d_{\theta}^h(x) \right]_i \Delta_i^h(x)$$

$$\frac{\partial E}{\partial (\mathbf{b}_h)_i}(\theta) = \left[d_{\theta}^h(x) \right]_i \Delta_i^h(x)$$

où $\Delta^h(x) \in \mathbb{R}^{n_h}$ est défini par

$$\Delta^{h}(x) = \begin{cases} \nabla_{y_{1}} L(f_{\theta}(x), y) & \text{, si } h = H \\ \mathbf{W}_{h+1}^{T}.\operatorname{diag}\left(d_{\theta}^{h+1}(x)\right).\Delta^{h+1}(x) & \text{, sinon} \end{cases}$$

Rq: \mathbf{W}_{h+1}^T remonte d'un niveau dans le réseau: **Rétropropagation**.

Rétropropagation: Problèmes connus

"Vanishing/Exploding gradient":

- ightharpoonup Si \mathbf{W}_{h+1}^T . diag $\left(d_{\theta}^{h+1}(x)\right)$ est systématiquement une contraction \Longrightarrow "Vanishing gradient"
- ightharpoonup Si \mathbf{W}_{h+1}^T . diag $\left(d_{\theta}^{h+1}(x)\right)$ augmente systématiquement la norme \Longrightarrow "Exploding gradient"

Rétropropagation: Problèmes connus

Dans certaines régions de l'espace, la fonction objectif est très irrégulière :

Pour une fonction d'activation positivement homogène¹ (notamment ReLU):

Pour $\lambda > 0$, $\lambda \sim 0$

$$\begin{aligned} \mathbf{W}_1^{\lambda} &= \lambda^{H-1} \mathbf{W}_1 & \text{et} & \mathbf{W}_h^{\lambda} &= \lambda^{-1} \mathbf{W}_h, & \forall h = 2..H \\ \mathbf{b}_1^{\lambda} &= \lambda^{H-1} \mathbf{b}_1 & \text{et} & \mathbf{b}_h^{\lambda} &= \lambda^{H-h} \mathbf{b}_h, & \forall h = 2..H \end{aligned}$$

On a, pour tout x, $f_{\theta^{\lambda}}(x) = f_{\theta}(x)$, pour h > 1

$$\frac{\partial E}{\partial \mathbf{W}_{h,i,j}}(\theta^{\lambda}) = \left[f_{\theta^{\lambda}}^{h-1}(x) \right]_{j} \left[d_{\theta^{\lambda}}^{h}(x) \right]_{i} (\Delta_{\lambda}^{h})_{i}(x)
= \lambda^{H-(h-1)} \left[f_{\theta}^{h-1}(x) \right]_{j} \left[d_{\theta}^{h}(x) \right]_{i} \lambda^{-(H-h)} \Delta_{i}^{h}(x)
= \lambda \frac{\partial E}{\partial \mathbf{W}_{h,i,i}}(\theta)$$

mais
$$\frac{\partial E}{\partial \mathbf{b}_{h,i}}(\theta^{\lambda}) = \lambda^{-(H-h)} \frac{\partial E}{\partial \mathbf{b}_{h,i}}(\theta)$$
 et $\frac{\partial E}{\partial \mathbf{W}_{1,i}}(\theta^{\lambda}) = \lambda^{-(H-1)} \frac{\partial E}{\partial \mathbf{W}_{1,i}}(\theta)$.

Les petits b_h^{λ} ont de forts gradients et seront multipliés par des forts W_h^{λ} .

¹ Attention à la non-différentiabilité en 0

Pour gagner en profondeur, faciliter l'apprentissage

Couches ResNet

$$f_h(x) = \sigma \Big(W_h \, \sigma \big(W_{h-1} f_{h-2}(x) + b_{h-1} \big) + b_h + f_{h-2}(x) \Big)$$

- Batch normalisation : On centre (au mieux) les données après chaque couche, à l'aide d'un opérateur diagonal.
- Contraindre les matrices \mathbf{W}_h à être orthogonales ou presque. Bonus: Gain de robustesse.
- Augmenter la largeur. Pb: Complexité = largeur².

Plan

- Introductio
 - L'optimisation des réseaux de neurones
 - Le paysage de la fonction objectif
 - Introduction
 - Paysage pour les réseaux larges
 - La régularisation induite par la géométrie
 - Le paysage pour les réseaux linéaires



- Le paysage de la fonction objectif
 - Introduction

Paysage de la fonction objectif: Introduction

Rappel, pour tout θ :

```
\begin{array}{lll} 0 \leq & R(f_{\theta}) & - & R^* & \text{(l'excès de risque)} \\ & = & R(f_{\theta}) & - & \widehat{R}(f_{\theta}) & \text{(erreur de généralisation)} \\ & + & \widehat{R}(f_{\theta}) & - & \widehat{R}(f_{\theta}) & \text{(erreur d'optimisation)} \\ & + & \widehat{R}(f_{\theta}) & - & \widehat{R}(f_{\theta^*}) & \leq \varepsilon \\ & + & \widehat{R}(f_{\theta^*}) & - & R(f_{\theta^*}) & \text{(erreur de généralisation)} \\ & + & R(f_{\theta^*}) & - & R^* & \text{(erreur d'approximation)} \end{array}
```

On utilise un algorithme d'optimisation pour trouver un θ , on veut que

$$\widehat{R}(f_{\theta}) - \inf_{\theta} \widehat{R}(f_{\theta})$$

soit le plus faible possible.

Idéalement, on voudrait que cette quantité soit nulle ou au moins on voudrait la borner supérieurement.

Paysage de la fonction objectif : Introduction

On suppose que $\theta \longmapsto \widehat{R}(f_{\theta})$ est C^2 partout. On a

$$\widehat{R}(f_{\theta}) = \widehat{R}(f_{\theta^*}) + \langle \nabla_{\theta} \widehat{R}(f_{\theta^*}), \theta - \theta^* \rangle + \frac{1}{2} \langle \nabla_{\theta}^2 \widehat{R}(f_{\theta^*})(\theta - \theta^*), \theta - \theta^* \rangle + o(\|\theta - \theta^*\|^2)$$

On distingue:

- θ* est un minimiseur global:
 - $\blacktriangleright \ \forall \theta, \qquad \widehat{R}(f_{\theta^*}) \leq \widehat{R}(f_{\theta})$
- $\widehat{R}(f_{\theta^*}) = \min_{\theta} \widehat{R}(f_{\theta})$ • θ^* est un minimiseur local:
- ► Il existe un voisinage ouvert \mathscr{O} de θ^* tel que

$$\forall \theta \in \mathcal{O}, \qquad \widehat{R}(f_{\theta^*}) \leq \widehat{R}(f_{\theta})$$

- θ^* est un point critique du second ordre:
 - ► On a

$$\nabla \widehat{R}(f_{\theta^*}) = 0$$
 et $\nabla^2 \widehat{R}(f_{\theta^*}) \ge 0$

- θ^* est un point critique du premier ordre:
 - ► On a

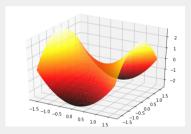
$$\nabla \widehat{R}(f_{\theta^*}) = 0$$

53/104

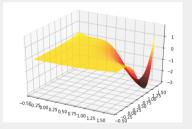
F. Malgouyres Big Data 2025-2026

Paysage de la fonction objectif : Introduction

- θ^* est un point selle si c'est un point critique qui n'est ni un minimiseur local, ni un maximiseur local
 - Un point selle θ* est strict : si ce n'est pas un point critique du second ordre (i.e., le Hessian a une v.p. négative).
 - Un point selle θ^* est non-strict: si c'est un point critique du second ordre (i.e. le Hessian est semi-defini positif et a une v.p. égale à 0. Typiquement, un terme d'ordre supérieur en fait un point selle.).



(a) Point selle strict



(b) Point selle non-strict

Paysage de la fonction objectif: Introduction

Optimisation non convexe avec les mains

Pour des fonctions non-convexes, on sait montrer que

- dans un cadre assez vaste, l'algorithme du gradient (ou gradient stochastique) converge vers un point critique du premier ordre
- dans un cadre plus restreint, l'algorithme du gradient converge vers un point critique du second ordre
- Pour le gradient stochastic, sans vitesse de convergence, que les itérés convergent vers un minimiseur local.
- S. Gadat, F. Panloup, S. Saadane. "Stochastic heavy ball." Electronic Journal of Statistics 12.1 (2018): 461-529.
- J. Lee, M. Simchowitz, M. Jordan, B. Recht." Gradient Descent Converges to Minimizers." COLT 2016.



- Classification et régression
- Les réseaux de neurones
 - Multi-Layer Perceptron/fully-connected neural network
 - Réseaux convolutifs
 - Réseaux Récurrents
 - Mécanismes d'attention et transformers
 - Beaucoup de variantes importantes
- L'optimisation des réseaux de neurones
- L'optimisation des réseaux de neurones: Premières propriété
- La Rétropropagation et ses défauts
- 3 Le paysage de la fonction objectif
 - Introduction
 - Paysage pour les réseaux larges
 - La régularisation induite par la géométrie
 - Overview sur un exemple
 - Résultats formels
- Le paysage pour les réseaux linéaires
 - Introduction
 - Paysage pour les réseaux linéaires

Paysage pour les réseaux larges

Différents énoncés décrits dans

- Gori, Tesi, "On the problem of local minima in backpropagation", IEEE PAMI, 1992
- Yu, Chen, "On the local minima free condition of backpropagation learning", IEEE Trans. Neural Networks, 1995
- Nguyen, Hein, "The loss surface of deep and wide neural networks", ICML, 2017
-

On considère:

- un problème de régression
- un réseau fully-connected de paramètre $\theta = (\mathbf{W}_1, \dots, \mathbf{W}_H, \mathbf{b}_1, \dots, \mathbf{b}_H)$
- des observations $(x^i, y^i)_{i=1..n}$:

$$E(\theta) = \widehat{R}(f_{\theta}) = \sum_{i=1}^{n} L(f_{\theta}(x^{i}) - y^{i})$$

Paysage pour les réseaux larges

Théoreme (Le Paysage pour les réseaux larges)

On suppose que σ est C^1 et que, pour tout $t \in \mathbb{R}$, $\sigma'(t) \neq 0$.

On suppose que le coût L est C^1 , à valeur dans \mathbb{R}^+ et tel que L(0) = 0. On suppose aussi que $\nabla L(y) = 0$ si et seulement si y = 0.

On note
$$X = [x^1 \cdots x^n] \in \mathbb{R}^{n_0 \times n}$$
 et $A = \begin{pmatrix} X \\ \mathbb{1}_n^T \end{pmatrix}$.

On considère un point critique du premier ordre $\theta = (\mathbf{W}_1, ..., \mathbf{W}_H, \mathbf{b}_1, ..., \mathbf{b}_H)$ de E.

On suppose que rang (A) = n et que, pour tout $h \in \{1, ..., H\}$, rang $(\mathbf{W}_h) = n_h$.

Alors, on a

$$\widehat{R}(f_{\theta})=0$$

et θ est un minimiseur global.

Les hypothèses fortes sont celles sur le rang. Elles impliquent notamment

$$n \le n_0 + 1$$
 et $n_H \le n_{H-1} \le \cdots \le n_0$

Nb: Ci-dessus l'indice 0 est arbitraire car on pourrait supposer que les x_i sont le résultat des premières couches d'un réseau.

Plan

- Introductio
 - L'optimisation des réseaux de neurones
 - Le paysage de la fonction objectif
 - La régularisation induite par la géométrie
 - Overview sur un exemple
 - Résultats formels
 - Le paysage pour les réseaux linéaires



- Classification et régression
- Les réseaux de neurones
 - Multi-Layer Perceptron/fully-connected neural network
 - Réseaux convolutifs
 - Réseaux Récurrents
 - Mécanismes d'attention et transformers
 - Beaucoup de variantes importantes
- L'optimisation des réseaux de neurones
- L'optimisation des réseaux de neurones: Premières propriété
- La Rétropropagation et ses défauts
- Le paysage de la fonction objectif
 - Introduction
 - Paysage pour les réseaux larges
 - La régularisation induite par la géométrie
 - Overview sur un exemple
 - Résultats formels
- Le paysage pour les réseaux linéaires
 - Introduction
 - Paysage pour les réseaux linéaires

Geometry-induced implicit regularization:

Deep learning context

- Implicit regularization:
 - ▶ Neural networks often have far more parameters than training examples.
 - ► Their performance is not explained by classical arguments.
 - ► Some global minimizers generalize well, some don't. The optimization finds a solution that does.
- Puzzling empirical observation:
 - ► Neuron alignment
 - Saddle-to-saddle dynamics
 - ► Flat-minimas generalize well
- NN are successful independently of the learning objective:
 - ► In classification, regression, generative models...
- For many unrelated data types:
 - Vectors, images, audio, natural langage, time-series,...

→ Need to study Neural Networks

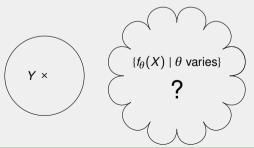
independently of the learning objective, algorithm, data type etc

Geometry of deep ReLU neural networks

- Given a deep ReLU architecture
- $X \in \mathbb{R}^{n_0 \times n}$ (resp $Y \in \mathbb{R}^{n_L \times n}$) contains n input (resp output) examples in \mathbb{R}^{n_0}
- Denoting $f_{\theta}(X) \in \mathbb{R}^{n_L \times n}$ the prediction of X by the neural network of parameter $\theta \in \mathbb{R}^p$
- Learning performs a local search solving (for instance)

$$\operatorname{argmin}_{\theta} \|f_{\theta}(X) - Y\|_{F}^{2} \iff \operatorname{argmin}_{Y' \in \{f_{\theta}(X) \mid \theta \text{ varies}\}} \|Y' - Y\|_{F}^{2}$$

- We study the local geometry of two sets:
 - ▶ the *image set* $\{f_{\theta}(X) \mid \theta \text{ varies}\};$
 - the pre-image set $\{\theta' \mid f_{\theta'}(X) = f_{\theta}(X)\}.$



Analogy with ℓ^1 regularization

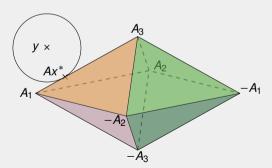
- Given $A \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$
- ullet Write ℓ^1 regularization under the form

$$| \operatorname{argmin}_{x} ||Ax - y||^{2}$$

$$||x||_{1} \le \tau$$

• The geometrical object is:

$${Ax \mid ||x||_1 \le \tau} = \tau \operatorname{conv}(A_1, -A_1, \cdots, A_p, -A_p).$$



Geometry of Neural Networks: In parameter space

- We study the local geometry of two sets:
 - the *image set* $\{f_{\theta}(X) \mid \theta \text{ varies}\};$
 - the pre-image set $\{\theta' \mid f_{\theta'}(X) = f_{\theta}(X)\}$.



Notation for Deep ReLU architecture

- Depth L, widths $(n_0, ..., n_L)$,
- $\sigma_I = ReLU$, $\forall I \neq L$, and σ_L is analytic (e.g. the identity of softmax)
- Parameter $\theta = (W_L, ..., W_1, b_L, ..., b_1) \in \mathbb{R}^p$
- For all l = 0, ..., L, and all $x \in \mathbb{R}^{n_0}$

$$f_l(x) = \begin{cases} x & \text{if } l = 0\\ \sigma_l(W_l f_{l-1}(x) + b_l) & \text{otherwise.} \end{cases}$$

- For all $x \in \mathbb{R}^{n_0}$, $f_{\theta}(x) = f_L(x)$.
- When differentiable, we denote $\mathbf{D} f_{\theta}(X) : \mathbb{R}^p \longrightarrow \mathbb{R}^{n_L \times n}$ such that

$$f_{\theta+\theta'}(X) = f_{\theta}(X) + \mathbf{D}f_{\theta}(X)(\theta') + o(\|\theta'\|)$$

$$\{f_{\theta}(X) + \mathbf{D}f_{\theta}(X)(\theta') \mid \|\theta' - \theta\| < \varepsilon\}$$

$$\{f_{\theta'}(X) \mid \|\theta' - \theta\| < \varepsilon\}$$

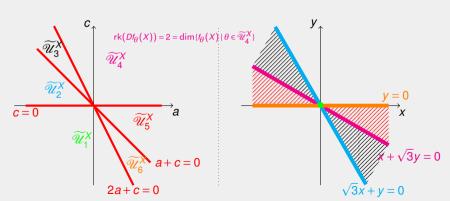
• A key quantity : $\operatorname{rk}(\mathbf{D}f_{\theta}(X)) = \dim\{f_{\theta'}(X) \mid \|\theta' - \theta\| < \varepsilon\}$

F. Malgouyres Big Data 2025-2026

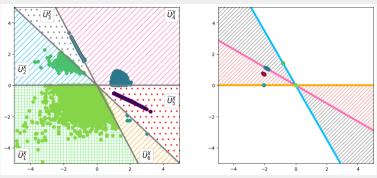
Why does geometry induce regularization? (1/2)

- A ReLU network of parameter $\theta = (a, b, c, d) \in \mathbb{R}^4$: $f_{\theta}(x) = b\sigma(ax + c) + d$
- A sample X = (0, 1, 2)
- We have $f_{\theta}(X) = (b\sigma(c) + d, b\sigma(a+c) + d, b\sigma(2a+c) + d) \in \mathbb{R}^{1\times 3}$

$$\theta = (a, b, c, d) \in \mathbb{R}^4$$
 $\xrightarrow{f_{\theta}(X)} \mathbb{R}^{1 \times 3} \longrightarrow \mathbb{R}^2 \sim (1, 1, 1)^{\perp}$



Why does geometry induce regularization? (2/2)



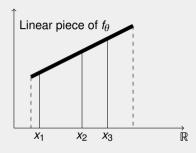
Denote

 $d(\theta)$: lower-semicontinuous envelop of $\operatorname{rk}(Df_{\theta}(X))$

For a learning objective L

 \Rightarrow Implicit regularization with $d(\theta)$

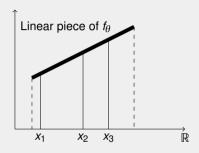
What is the effect of this regularization? (1/1)



 $Df_{\theta}(x_1), Df_{\theta}(x_2), Df_{\theta}(x_3)$ are linearly dependent $\Rightarrow \operatorname{rk}(Df_{\theta}(X))$ is reduced

 \Rightarrow Regularization reduces the number of linear pieces 'perceived' by X, fills holes linearly

What is the effect of this regularization? (1/1)



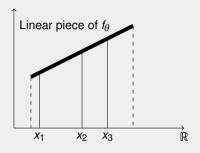
 \Rightarrow rk($Df_{\theta}(X)$) is reduced \Rightarrow Regularization reduces the number of linear pieces 'perceived' by X, fills holes linearly

 $Df_{\theta}(x_1), Df_{\theta}(x_2), Df_{\theta}(x_3)$ are linearly dependent

Experiment: Recovery of y = |x|

X contains 10 random numbers
ReLU MLP of Width = 20 and depth = 6

What is the effect of this regularization? (1/1)

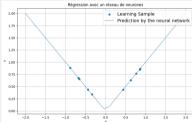


 $Df_{\theta}(x_1), Df_{\theta}(x_2), Df_{\theta}(x_3)$ are linearly dependent $\Rightarrow \text{rk}(Df_{\theta}(X))$ is reduced $\Rightarrow \text{Regularization reduces the number of linear pieces 'perceived' by <math>X$, fills holes linearly

Experiment: Recovery of y = |x|

X contains 10 random numbers
ReLU MLP of Width = 20 and depth = 6
init = he_normal, Optimizer: Adam

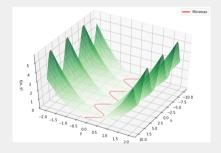
Learning rate= 0.002, batch size= 5, 500 epochs



Regularization and minima flatness

Define $L(\theta) \triangleq R(f_{\theta}(X))$

If θ is a local minimizer of L, and $\theta' \longmapsto f_{\theta'}(X)$ is smooth with constant local dimension in the vicinity of θ , \Rightarrow there exists $\varepsilon > 0$ s.t. $\{\theta' \mid f_{\theta'}(X) = f_{\theta}(X), \|\theta' - \theta\| < \varepsilon\}$ are local minimizers of L.



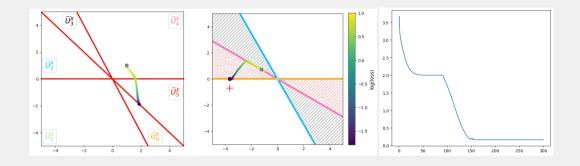
Constant rank theorem: $\dim\{\theta' \mid f_{\theta'}(X) = f_{\theta}(X), \|\theta' - \theta\| < \varepsilon\} = nb_{param} - d(\theta)$

- ⇒ More regularization
 - ⇒ 'valleys' of larger dimension
 - ⇒ escape is difficult for stoch. algo.

Near a critical point $\theta^{k+1} = \theta^k - \eta(\nabla L(\theta^k) + noise)$.

F. Malgouyres Big Data 2025-2026 69/104

Saddle-to-saddle dynamic on the example

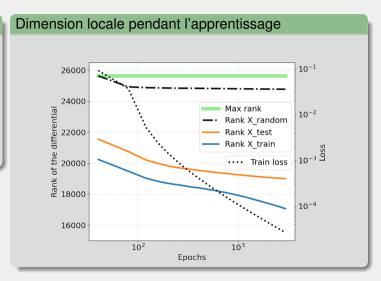


The saddle-to-saddle phenomenon occurs when the sequence $(d(\theta^k))_{k \in \mathbb{N}}$ is increasing.

Experience 1: Comportement pendant l'apprentissage

Description de l'expérience

- architecture (784,30,30,30,10);
- nombre of parametres: 25720;
- Données d'apprentissage MNIST X_{train} de taille 4000;
- Données de test MNIST X_{test} de taille 20000;
- Données aléatoires X_{random} (Bruit Gaussien) de taille 20000.



Experience 2: Comportement quand la largeur varie

Description de l'expérience

- architecture (784, w, w, w, 10), pour plusieurs w;
- Données d'apprentissage MNIST X_{train} de taille 4000;
- Données de test MNIST X_{test} de taille 10000;
- Données aléatoires X_{random} (Bruit Gaussien) de taille 40000.

Functional dimensions as the width increases





- Classification et régression
- Les réseaux de neurones
 - Multi-Layer Perceptron/fully-connected neural network
 - Réseaux convolutifs
 - Réseaux Récurrents
 - Mécanismes d'attention et transformers
 - Beaucoup de variantes importantes
- L'optimisation des réseaux de neurones
- L'optimisation des réseaux de neurones: Premières propriété
- La Rétropropagation et ses défauts
- Le paysage de la fonction objectif
 - Introduction
 - Paysage pour les réseaux larges
- La régularisation induite par la géométrie
 - Overview sur un exemple
 - Résultats formels
- Le paysage pour les réseaux linéaires
 - Introduction
 - Paysage pour les réseaux linéaires

On note E les arcs du réseaux et B les neurones des couches 1 jusqu'à H.

On note $\theta \in \mathbb{R}^E \times \mathbb{R}^B$ que l'on identifie, si besoin, à

$$(W_1,\ldots,W_H,b^1,\ldots,b^H)\in (\mathbb{R}^{n_1\times n_0}\times\cdots\times\mathbb{R}^{n_H\times n_{H-1}})\times (\mathbb{R}^{n_1}\times\cdots\times\mathbb{R}^{n_H}).$$

On rappelle

$$\begin{cases}
f_0(x) = x \\
f_h(x) = \sigma_h(W_h f_{h-1}(x) + b_h)
\end{cases}, \forall h = 1, \dots, H$$

Pour $n \in \mathbb{N}^*$, on définit la fonction pattern d'activation

$$a: \mathbb{R}^{n_0 \times n} \times (\mathbb{R}^E \times \mathbb{R}^B) \longrightarrow \{0,1\}^{(n_1 + \dots + n_{H-1}) \times n}$$

 $(X,\theta) \longmapsto a(X,\theta)$

οù

$$a(X,\theta)_{i,j} = \begin{cases} 1 & \text{si } [W_h f_{h-1}(x^j) + b_h]_v \ge 0 \\ 0 & \text{sinon,} \end{cases}$$

pour v le neurone correspondant à $i, h \in \{1, ..., H-1\}$ la couche de v, et la j^{ieme} colonne x^j de X.

Pour $X \in \mathbb{R}^{n_0 \times n}$ fixé, elle atteint un nombre fini de valeur que l'on note $\Delta_1^X, \dots, \Delta_{q_X}^X$.

F. Malgouyres Big Data 2025-2026

On note, pour $X \in \mathbb{R}^{n_0 \times n}$ fixé et pour $j \in \{1, ..., q_X\}$,

$$\widetilde{\mathscr{U}}_{j}^{X} = \operatorname{Int}\left\{\theta \in \mathbb{R}^{E} \times \mathbb{R}^{B} \mid a(X,\theta) = \Delta_{j}^{X}\right\},\,$$

et $m_X = \left| \left\{ j \in \{1, ..., q\} \mid \widetilde{\mathcal{U}}_j^X \neq \emptyset \right\} \right|.$

Quitte à changer l'ordre, on suppose que les $\widetilde{\mathscr{U}}_1^X,\ldots,\widetilde{\mathscr{U}}_{m_X}^X$ sont non-vides.

Lemme (Prédiction polynomiale par morceaux en θ)

Pour tout $n \in \mathbb{N}^*$, $X \in \mathbb{R}^{n_0 \times n}$, les ensembles $\widetilde{\mathcal{U}}_1^X, \dots, \widetilde{\mathcal{U}}_{m_X}^X$ sont non-vides, ouverts et disjoints

- Pour tout $j = 1, ..., m_X, \theta \longmapsto f_{\theta}(X)$ est une fonction polynomiale de degré inférieur à H sur $\widetilde{\mathcal{U}}_j^X$.
- Le complémentaire $\left(\cup_{j=1}^{m_X}\widetilde{\mathscr{U}_j^X}\right)^c$ est fermé et de mesure de Lebesgue nulle.

Rq: La fonction $\theta \mapsto f_{\theta}(X)$ est une composition de fonctions continues, elle est continue.

Pour $n \in \mathbb{N}^*$, $X \in \mathbb{R}^{n_0 \times n}$, et $j \in \{1, ..., m_X\}$, on note

$$r_j^X = \max_{\theta \in \widetilde{\mathcal{U}}_i^X} \operatorname{rk} \left(\mathbf{D} f_{\theta}(X) \right) \qquad \text{et} \qquad \mathcal{U}_j^X = \{ \theta \in \widetilde{\mathcal{U}}_j^X \mid \operatorname{rk} \left(\mathbf{D} f_{\theta}(X) \right) = r_j^X \}.$$

Théoreme (J. Bona-Pellissier, F. Bachoc, F. Malgouyres 2024)

Pour tout réseau ReLU, $n \in \mathbb{N}^*$, $X \in \mathbb{R}^{n_0 \times n}$, par définition

- $\mathcal{U}_1^X, \dots, \mathcal{U}_{m_Y}^X$ sont non-vides et disjoints;
- Pour tout $j \in \{1,...,m_X\}$, la fonction $\theta \mapsto a(X,\theta)$ est constante sur \mathcal{U}_j^X et prend des valeurs distinctes pour $j' \neq j$;
- Pour tout $j \in \{1,...,m_X\}$, $\theta \mapsto \operatorname{rk}(\mathbf{D}f_{\theta}(X))$ est constante sur \mathcal{U}_i^n et vaut r_i^X .

De plus,

- Les ensembles $\mathcal{U}_1^X, \dots, \mathcal{U}_{m_Y}^X$ sont ouverts;
- $\left(\bigcup_{i=1}^{m_X} \mathcal{U}_i^X\right)^c$ est fermé et de mesure de Lebesgue nulle;
- Pour tout $j \in \{1,...,m_X\}$, $\theta \mapsto f_{\theta}(X)$ est une fonction polynomiale de degré inférieur à H sur \mathcal{U}_j^X .

Corollaire

Pour tout réseaux ReLU fully-connected, profond.

Pour tout $n \in \mathbb{N}^*$, $X \in \mathbb{R}^{n_0 \times n}$, $j \in \{1, ..., m_X\}$ et $\theta \in \mathcal{U}_i^X$, il existe $\varepsilon_{X,\theta} > 0$ tel que

L'image locale

$$\{f_{\theta'}(X) \in \mathbb{R}^{n_L \times n} \mid \|\theta' - \theta\| < \varepsilon_{X,\theta}\}$$

est une variété régulière de dimension $rk(\mathbf{D}f_{\theta}(X))$;

la pre-image

$$\{\theta' \in \mathbb{R}^E \times \mathbb{R}^B \mid f_{\theta'}(X) = f_{\theta}(X) \text{ and } \|\theta' - \theta\| < \varepsilon_{X,\theta}\}$$

est une variété régulière de dimension $|E| + |B| - \operatorname{rk}(\mathbf{D}f_{\theta}(X))$.

C'est une conséquence directe du théorème précédent et du théorème du rang constant.

Théorème du rang constant.

Theorem (Théorème du rang constant)

Soient $U \subset \mathbb{R}^n$ un ouvert, $a \in U$, et $f: U \to \mathbb{R}^p$ une application de classe \mathscr{C}^1 . Si la différentielle de f est de rang constant f sur G, alors il existe :

- $un \mathscr{C}^1$ -difféomorphisme φ d'un ouvert $V \subset \mathbb{R}^n$ contenant 0 sur un ouvert de U, avec $\varphi(0) = a$, et
- $un \mathscr{C}^1$ -difféomorphisme ψ d'un ouvert de \mathbb{R}^p contenant $f(\varphi(V))$ sur un ouvert de \mathbb{R}^p , avec $\psi(f(a)) = 0$, tels que :

$$\forall x = (x_1, \ldots, x_n) \in V, \quad (\psi \circ f \circ \varphi)(x) = (x_1, \ldots, x_r, 0, \ldots, 0).$$

On note $r = \operatorname{rk}(\mathbf{D}f_{\theta}(X))$, n = |E| + |B|, $p = nn_L$, $a = \theta$, $f : \theta' \mapsto f_{\theta'}(X)$. On prend $\varepsilon_{X,\theta}$ tel que $B(\theta, \varepsilon_{X,\theta}) \subset \varphi(V)$ et on note $V' = \varphi^{-1}(B(\theta, \varepsilon_{X,\theta}))$.

$$\varphi \quad \text{ est une carte de } \quad V' \cap \left(\{0\} \times \mathbb{R}^{n-r}\right) \quad \text{ sur } \quad \left\{\theta' \in \mathbb{R}^E \times \mathbb{R}^B \mid f_{\theta'}(X) = f_{\theta}(X) \quad \text{and} \quad \|\theta' - \theta\| < \varepsilon_{X,\theta}\right\},$$

et

$$\psi^{-1} \quad \text{ est une carte de } \quad \psi \circ f \circ \varphi \left(V' \cap \left(\mathbb{R}^r \times \{0\} \right) \right) \quad \text{ sur } \quad \left\{ f_{\theta'} \left(X \right) \in \mathbb{R}^{n_L \times n} \mid \|\theta' - \theta\| < \varepsilon_{X,\theta} \right\}.$$

F. Malgouyres Big Data 2025-2026

Regularization induite par la géométrie

$$\dim^{-}(\theta, X) = \lim_{\epsilon \to 0} \min_{j: B(\theta, \epsilon) \cap \mathcal{U}_{j}^{X} \neq \emptyset} r_{j}^{X},$$

$$\dim^{+}(\theta, X) = \lim_{\epsilon \to 0} \max_{j: B(\theta, \epsilon) \cap \mathcal{U}_{i}^{X} \neq \emptyset} r_{j}^{X}.$$

$$(P)$$
: Minimise $R(f_{\theta}(X))$.

$$(P_k)$$
: $\underset{\theta:\dim^-(\theta,X)\leq k}{\text{Minimise}} R(f_{\theta}(X)),$

pour $k \in \mathbb{N}$ donné.

Corollaire

Pour tout $n \in \mathbb{N}^*$, $X \in \mathbb{R}^{n_0 \times n}$, et toute fonction objectif régulière $R : \mathbb{R}^{n_L \times n} \longrightarrow \mathbb{R}$,

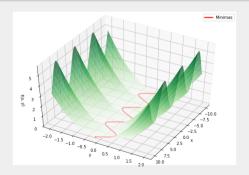
- Pour tout $\theta^* \in \bigcup_{i=1}^{p_X} \mathscr{U}_i^X$. On note $k = \operatorname{rk}(\mathbf{D}f_{\theta^*}(X))$.
 - θ^* est un point critique de (P) \iff $(\theta^*, 1)$ satisfait les conditions KKT de (P_k) .
- Pour tout $\theta^* \in \mathbb{R}^E \times \mathbb{R}^B$. On note $k = \dim^+(\theta^*, X)$. Nous avons
 - θ^* est un minimiseur local de (P) \iff θ^* est un minimiseur local de (P_k) .
 - θ^* est un point selle de (P) \iff θ^* est un point selle de (P_k) .

Minima flatness

Any local minimizer $\theta \in \bigcup_{i=1}^{p_X} \mathscr{U}_i^X$ of (P) is dimension $(|E| + |B| - \operatorname{rk}(\mathbf{D}f_{\theta}(X)))$ flat.

Définition

A local minimizer θ of (P) is said to be **dimension** k **flat**, for $k \in \mathbb{N}$, if and only if there exist $\varepsilon > 0$ and a smooth manifold $\mathcal{M} \subset \mathbb{R}^R \times \mathbb{R}^B$ of dimension k such that $\theta \in \mathcal{M}$, and every $\theta' \in \mathcal{M} \cap B(\theta, \varepsilon)$ is also a local minimizer of (P).



Lien avec le nombre de zones affines perçues par X: Cas une couche cachée On considère un réseau ReLU de largeurs $(1, n_1, 1)$, avec $\sigma_2 = Id$.

On rappelle que $a(X,\theta) \in \{0,1\}^{n_1 \times n}$, et pour $k \in [1,n_1], i \in [1,n]$

$$a(X,\theta)_{k,i} = \begin{cases} 1 & \text{if } [W_1 x^{(i)} + b_1]_k \ge 0, \\ 0 & \text{sinon.} \end{cases}$$

On note, pour tout X et θ ,

$$\mathscr{A}(X,\theta) = \{\delta \in \{0,1\}^{n_1} \mid \text{il existe } i \in [1,n] \text{ satisfaisant } a(x^{(i)},\theta) = \delta\}.$$

Pour $i \in [1, n]$,

$$\begin{cases}
\mathbf{e}_{i} = (0, \dots, 0, x^{(i+1)} - x^{(i)}, \dots, x^{(n)} - x^{(i)}), \\
\mathbf{e}_{n+i} = (x^{(i)} - x^{(1)}, \dots, x^{(i)} - x^{(i-1)}, 0, \dots, 0).
\end{cases}$$

Ainsi que $\mathbf{e}_0 = \mathbf{e}_{2n}$. On définit aussi, pour tout $i \in [1, n]$,

$$\mathbb{1}_{i} = (0, \dots, 0, 1, \dots, 1) \in \mathbb{R}^{1 \times n}$$
 and $\mathbb{1}_{n+i} = (1, \dots, 1, 0, \dots, 0) \in \mathbb{R}^{1 \times n}$.

Lien avec le nombre de zones affines perçues par X: Cas une couche cachée

Théoreme

On considère une architecture ReLU (E,V,Id), avec L=2 et $n_0=n_2=1$. On considère $n \in \mathbb{N}^*$, et un échantillon $X=(x^{(1)},x^{(2)},\ldots,x^{(n)}) \in \mathbb{R}^{1\times n}$ tel que

$$x^{(1)} < x^{(2)} < \cdots < x^{(n)}$$

Pour tout $j \in [1, p_X]$, il existe $\alpha \in \{1, ..., 2n\}^{n_1}$ tel que pour tout $\theta \in \mathcal{U}_j^X$ et tout $k \in [1, n_1]$, $a(X, \theta)_{k,:} = \mathbb{1}_{\alpha_k}$, et

$$\mathsf{rk}\left(\mathbf{D}f_{\theta}(X)\right) = \mathsf{rk}\left(\mathbb{1}, \mathbf{e}_{\alpha_{1}-1}, \mathbf{e}_{\alpha_{1}}, \dots, \mathbf{e}_{\alpha_{n_{1}}-1}, \mathbf{e}_{\alpha_{n_{1}}}\right).$$

On a donc $\mathcal{U}_{j}^{X} = \widetilde{\mathcal{U}}_{j}^{X}$ et, pour tout $\theta \in \mathcal{U}_{j}^{X}$,

$$\frac{1}{2}|\mathscr{A}(X,\theta)| \le \operatorname{rk}(\mathbf{D} f_{\theta}(X)) \le 2|\mathscr{A}(X,\theta)|.$$

F. Malgouyres Big Data 2025-2026 82/104

Plan

- 1 Introduction
 - L'optimisation des réseaux de neurones
 - Le paysage de la fonction objectit
 - La régularisation induite par la géométrie
 - Le paysage pour les réseaux linéaires
 - Introduction
 - Paysage pour les réseaux linéaires
 - Paysage pour les réseaux linéaires (cas à 1 couche cachée)



- Le paysage pour les réseaux linéaires
 - Introduction

Paysage de la fonction objectif: Introduction

Rappel, pour tout θ :

```
\begin{array}{lll} 0 \leq & R(f_{\theta}) & - & R^* & \text{(l'excès de risque)} \\ & = & R(f_{\theta}) & - & \widehat{R}(f_{\theta}) & \text{(erreur de généralisation)} \\ & + & \widehat{R}(f_{\theta}) & - & \widehat{R}(f_{\theta}) & \text{(erreur d'optimisation)} \\ & + & \widehat{R}(f_{\theta}) & - & \widehat{R}(f_{\theta^*}) & \leq \varepsilon \\ & + & \widehat{R}(f_{\theta^*}) & - & R(f_{\theta^*}) & \text{(erreur de généralisation)} \\ & + & R(f_{\theta^*}) & - & R^* & \text{(erreur d'approximation)} \end{array}
```

On utilise un algorithme d'optimisation pour trouver un θ , on veut que

$$\widehat{R}(f_{\theta}) - \inf_{\theta} \widehat{R}(f_{\theta})$$

soit le plus faible possible.

Idéalement, on voudrait que cette quantité soit nulle ou au moins on voudrait la borner supérieurement.

Paysage de la fonction objectif : Introduction

On suppose que $\theta \longmapsto \widehat{R}(f_{\theta})$ est C^2 partout. On a

$$\widehat{R}(f_{\theta}) = \widehat{R}(f_{\theta^*}) + \langle \nabla_{\theta} \widehat{R}(f_{\theta^*}), \theta - \theta^* \rangle + \frac{1}{2} \langle \nabla_{\theta}^2 \widehat{R}(f_{\theta^*})(\theta - \theta^*), \theta - \theta^* \rangle + o(\|\theta - \theta^*\|^2)$$

On distingue:

- θ* est un minimiseur global:
 - $\blacktriangleright \ \forall \theta, \qquad \widehat{R}(f_{\theta^*}) \leq \widehat{R}(f_{\theta})$
- $\widehat{R}(f_{\theta^*}) = \min_{\theta} \widehat{R}(f_{\theta})$ • θ^* est un minimiseur local:
- ▶ Il existe un voisinage ouvert \emptyset de θ^* tel que

$$\forall \theta \in \mathcal{O}, \qquad \widehat{R}(f_{\theta^*}) \leq \widehat{R}(f_{\theta})$$

- θ^* est un point critique du second ordre:
 - ► On a

$$\nabla \widehat{R}(f_{\theta^*}) = 0$$
 et $\nabla^2 \widehat{R}(f_{\theta^*}) \ge 0$

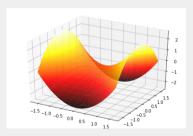
- θ^* est un point critique du premier ordre:
 - ► On a

$$\nabla \widehat{R}(f_{\theta^*}) = 0$$

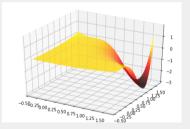
F. Malgouyres Big Data 2025-2026 86/104

Paysage de la fonction objectif : Introduction

- \bullet θ^* est un point selle si c'est un point critique qui n'est ni un minimiseur local, ni un maximiseur local
 - Un point selle θ* est strict : si ce n'est pas un point critique du second ordre (i.e., le Hessian a une v.p. négative).
 - Un point selle θ^* est non-strict: si c'est un point critique du second ordre (i.e. le Hessian est semi-defini positif et a une v.p. égale à 0. Typiquement, un terme d'ordre supérieur en fait un point selle.).



(a) Point selle strict



(b) Point selle non-strict

Paysage de la fonction objectif : Introduction

Optimisation non convexe avec les mains

Pour des fonctions non-convexes, on sait montrer que

- dans un cadre assez vaste, l'algorithme du gradient (ou gradient stochastique) converge vers un point critique du premier ordre
- dans un cadre plus restreint, l'algorithme du gradient converge vers un point critique du second ordre
- Pour le gradient stochastic, sans vitesse de convergence, que les itérés convergent vers un minimiseur local.
- S. Gadat, F. Panloup, S. Saadane. "Stochastic heavy ball." Electronic Journal of Statistics 12.1 (2018): 461-529.
- J. Lee, M. Simchowitz, M. Jordan, B. Recht." Gradient Descent Converges to Minimizers." COLT 2016.



- Le paysage pour les réseaux linéaires

 - Paysage pour les réseaux linéaires

Pointeurs bibliographiques

- Baldi, Hornik, "Neural networks and principal component analysis: Learning from examples without local minima", Neural networks, 1989.
- Baldi, Hornik, "Learning in linear neural networks: A survey", IEEE Transactions on neural networks, 1995.
- Kawaguchi, "Deep learning without poor local minima", NeurIPS 2016
- (Notamment dans les groupes de S. Arora à Princeton, de F. Bach à l'ENS)
- E. Achour, S. Gerchinovitz, F. Malgouyres, "The loss landscape of deep linear neural networks: a second-order analysis", Journal of Machine Learning Research, 2024.

On note

• $X \in \mathbb{R}^{d_X \times n}$ et $Y \in \mathbb{R}^{d_Y \times n}$ les matrices contenant les données.

•
$$\widehat{R}(\mathbf{W}) = \sum_{i=1}^{m} \|W_H W_{H-1} \cdots W_2 W_1 x_i - y_i\|_2^2 = \|W_H \cdots W_1 X - Y\|^2$$

•

•

$$\Sigma_{XX} = \sum_{i=1}^{n} x_i x_i^T = XX^T \in \mathbb{R}^{d_X \times d_X} \quad , \quad \Sigma_{YY} = \sum_{i=1}^{n} y_i y_i^T = YY^T \in \mathbb{R}^{d_Y \times d_Y},$$

$$\Sigma_{XY} = \sum_{i=1}^n x_i y_i^T = XY^T \in \mathbb{R}^{d_X \times d_Y} \quad , \quad \Sigma_{YX} = \sum_{i=1}^n y_i x_i^T = YX^T \in \mathbb{R}^{d_Y \times d_X},$$

 $\Sigma^{1/2} = \Sigma_{YX} \Sigma_{XX}^{-1} X \in \mathbb{R}^{d_y \times n},$

sa SVD

$$\Sigma^{1/2} = U \Delta V^T,$$

avec $U \in \mathbb{R}^{d_y \times d_y}$, $\Delta = \text{diag}((\delta_i)_{i=1..d_y}) \in \mathbb{R}^{d_y \times n}$, $V \in \mathbb{R}^{n \times n}$.

•
$$r_{max} = \min(d_x, n_1, ..., n_{H-1}, d_y)$$

On suppose

- $d_y \le d_x \le n$
- Σ_{XX} est inversible et Σ_{XY} est de rang plein
- les valeurs singulières de $\Sigma^{1/2}$ sont distinctes

Lemme (Inférence et valeurs critiques)

Soit $\mathbf{W} = (W_1, ..., W_H)$ un point critique du premier ordre de \widehat{R} et $r = rk(W_H \cdots W_1)$. Il existe un unique sous-ensemble $\mathscr{S} \subset [\![1,d_y]\!]$ de taille r tel que:

$$W_H \cdots W_1 = U_{\mathcal{S}} U_{\mathcal{S}}^T \Sigma_{YX} \Sigma_{XX}^{-1}$$
.

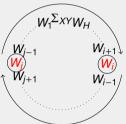
La valeur critique correspondante vaut $\widehat{R}(\mathbf{W}) = \operatorname{tr}(\Sigma_{YY}) - \sum_{i \in \mathscr{S}} \delta_i^2$. On dit que le point critique \mathbf{W} est associé à \mathscr{S} .

Proposition

Pour tout $\mathcal{S} \subset [1, d_v]$ de taille $r \in [0, r_{max}]$, il existe un point critique **W** associé à \mathcal{S} .

- Pivot $(i,j) \in [1,H]^2$, avec i > j
- Blocs complémentaires du pivot (i,j):

Premier bloc complémentaire : $W_{j-1} \cdots W_1 \Sigma_{XY} W_H \cdots W_{i+1}$



Second bloc complémentaire : $W_{i-1} \cdots W_{j+1}$

- Pivot tenu pour W : L'un des blocs complémentaires est de rang $rk(W_H...W_1)$
- W est point critique tenu: W est un point critique et tous les pivots sont tenus pour W

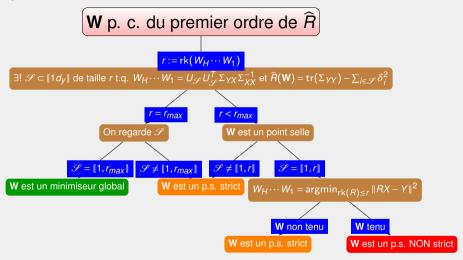


Figure: Theorem [E. Achour et al.]

Proposition (E. Achour et al.)

- Pour H = 2, il n'existe pas de point selle non-strict.
- Pour $H \ge 3$, pour tout $r < r_{max}$, il existe des points critiques associés à [1, r] tenus et non-tenus.

Proposition (E. Achour et al.)

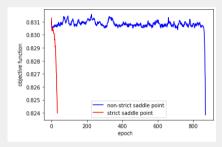
On note $\mathscr{S}_{max} = [1, r_{max}]$ et $Q_{max} = [1, d_y] \setminus \mathscr{S}_{max} = [r_{max} + 1, d_y]$.

W est un minimiseur global de \widehat{R} si et seulement si il existe des matrices inversibles

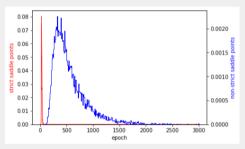
$$D_{H-1} \in \mathbb{R}^{d_{H-1} \times d_{H-1}}, \dots, D_1 \in \mathbb{R}^{d_1 \times d_1}, \text{ et des matrices } A_R \in \mathbb{R}^{(d_y - r_{max}) \times (d_{H-1} - r_{max})}, \\ (W_h)_{DR} \in \mathbb{R}^{(d_h - r_{max}) \times (d_{h-1} - r_{max})} \text{ pour } h \in [2, H-1], \text{ et } M_D \in \mathbb{R}^{(d_1 - r_{max}) \times d_x} \text{ tels que:}$$

$$\begin{aligned} W_{H} &= \begin{bmatrix} U_{\mathcal{S}_{max}}, U_{Q_{max}} A_{R} \end{bmatrix} D_{H-1}^{-1} \\ W_{h} &= D_{h} \begin{bmatrix} I_{r_{max}} & 0 \\ 0 & (W_{h})_{DR} \end{bmatrix} D_{h-1}^{-1} \qquad \forall \, h \in [2, H-1] \\ W_{1} &= D_{1} \begin{bmatrix} U_{\mathcal{S}_{max}}^{T} \sum_{YX} \sum_{XX}^{-1} \\ M_{D} \end{bmatrix}. \end{aligned}$$

Empiriquement, on trouve:



(a) Initialisation au voisinage d'un point selle strict vs non-strict



(b) Histogramme des epoch d'échappement

Conclusion [E. Achour et al.]

- Classification de l'ensemble des points critiques en: minimiseurs globaux; points selles stricts; point selles non-stricts.
- Tout point critique du second ordre qui n'est pas un minimiseur global conduit à une solution de la régression linéaire sous contrainte de rang.
- Les points selles non-stricts sont associés avec r_{max} valeurs plateau pour le risque empirique
- Paramétrisation des minimiseurs globaux



- Classification et régression
- Les réseaux de neurones
 - Multi-Layer Perceptron/fully-connected neural network
 - Réseaux convolutifs
 - Réseaux Récurrents
 - Mécanismes d'attention et transformers
 - Beaucoup de variantes importantes
- L'optimisation des réseaux de neurones
 - L'optimisation des réseaux de neurones: Premières propriété
 - La Rétropropagation et ses défauts
- Le paysage de la fonction objectif
 - Introduction
 - Paysage pour les réseaux larges
- La régularisation induite par la géométrie
 - Overview sur un exemple
 - Résultats formels
- 5 Le paysage pour les réseaux linéaires
 - Introduction
 - Paysage pour les réseaux linéaires

On simplifie

- Cas H=2
- On note pour $A \in \mathbb{R}^{n_2 \times n_1}$ et $B \in \mathbb{R}^{n_1 \times n_0}$

$$E(A, B) = \sum_{i=1}^{L} ||y_i - ABx_i||^2$$

• On note pour $D \in \mathbb{R}^{n_2 \times n_0}$

$$F(D) = \sum_{i=1}^{L} \|y_i - Dx_i\|^2$$

On note

$$\begin{split} \Sigma_{XX} &= \sum_{i=1}^{L} x_i x_i^T \in \mathbb{R}^{n_0 \times n_0} \qquad , \qquad \Sigma_{XY} = \sum_{i=1}^{L} x_i y_i^T \in \mathbb{R}^{n_0 \times n_2} \\ \Sigma_{YX} &= \sum_{i=1}^{L} y_i x_i^T \in \mathbb{R}^{n_2 \times n_0} \qquad , \qquad \Sigma_{YY} = \sum_{i=1}^{L} y_i y_i^T \in \mathbb{R}^{n_2 \times n_2} \end{split}$$

Remarques

On a

- Pour tout $C \in \mathbb{R}^{n_1 \times n_1}$ inversible, $AB = (AC)(C^{-1}B) = A'B'$

$$\Sigma_{YX}\Sigma_{XX}^{-1} \in \operatorname{argmin}_D F(D)$$

③ Soit $M \in \mathbb{R}^{n \times p}$ avec $p \le n$ de rang p. Pour tout $x \in \mathbb{R}^n$, la projection $P_M(x)$ de x sur l'espace vectoriel généré par les colonnes de M vaut

$$P_M(x) = M(M^T M)^{-1} M^T x$$

Lemme 1

Si Σ_{XX} est inversible. Soient $A \in \mathbb{R}^{n_2 \times n_1}$ telle que rang $(A) = n_1$ et $B \in \mathbb{R}^{n_1 \times n_0}$. Alors, (A, B) est un point critique du premier ordre de E si et seulement si

$$AB\Sigma_{XX}B^T = \Sigma_{YX}B^T$$
 et $B = (A^TA)^{-1}A^T\Sigma_{YX}\Sigma_{XX}^{-1}$

Lemme 2

Si Σ_{XX} est inversible. Soient $A \in \mathbb{R}^{n_2 \times n_1}$ telle que rang $(A) = n_1$ et $B \in \mathbb{R}^{n_1 \times n_0}$.

Les deux assertions suivantes sont équivalentes:

(A, B) est un point critique du premier ordre de E



$$AB = P_A \Sigma_{YX} \Sigma_{XX}^{-1}$$

et

$$P_A \Sigma = P_A \Sigma P_A = \Sigma P_A$$

pour

$$\Sigma = \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY} \in \mathbb{R}^{n_2 \times n_2}$$

On diagonalise (Σ est symétrique)

$$\Sigma = U\Lambda U^T$$

avec $\Lambda \in \mathbb{R}^{n_2 \times n_2}$ diagonale et $U \in \mathbb{R}^{n_2 \times n_2}$ unitaire.

Pour $\mathscr{S} \subset \{1, \dots, n_2\}$, on note $U_{\mathscr{S}}$ la matrice extraite de U en prenant les colonnes d'indice dans \mathscr{S} .

Proposition 1: Paramétrisation des points critiques

Si Σ_{XX} est inversible. Soient $A \in \mathbb{R}^{n_2 \times n_1}$ et $B \in \mathbb{R}^{n_1 \times n_0}$ avec A telle que rang $(A) = n_1$.

On suppose que les valeurs propres de Σ sont distinctes.

Alors, (A, B) est un point critique du premier ordre de E si et seulement si il existe $C \in \mathbb{R}^{n_1 \times n_1}$ inversible et $\mathcal{S} \subset \{1, \dots, n_2\}$ de taille n_1 tels que

$$A = U_{\mathcal{S}}C$$
 et $B = C^{-1}U_{\mathcal{S}}^T\Sigma_{YX}\Sigma_{XX}^{-1}$

Pour simplifier les notations, on suppose les valeurs propres de Σ ordonnées:

$$\lambda_1 > \lambda_2 > \cdots > \lambda_{n_2}$$

Proposition 2: minimiseur global ⇔ minimiseur local

Si Σ_{XX} est inversible. Soient $A \in \mathbb{R}^{n_2 \times n_1}$ et $B \in \mathbb{R}^{n_1 \times n_0}$ avec A telle que rang $(A) = n_1$.

On suppose que les valeurs propres de Σ sont distinctes.

• pour tout point critique du premier ordre (A, B) de E et pour $\mathcal S$ définissant A et B (voir la proposition précédente), on a

$$AB = P_{U_{\mathcal{S}}} \Sigma_{YX} \Sigma_{XX}^{-1}$$

$$E(A,B) = \operatorname{trace}(\Sigma_{YY}) - \sum_{i \in \mathscr{S}} \lambda_i$$

- ② Si (A, B) est un minimiseur global alors c'est un point critique du premier ordre associé à $\mathcal{S} = \{1, \dots, n_1\}$.
- (A, B) est minimiseur local si et seulement si c'est un minimiseur global.

Merci pour votre attention!