

Introduction aux mathématiques pour l'apprentissage de réseaux de neurones

François Malgouyres

Institut de Mathématiques de Toulouse, Université Paul Sabatier, ANITI



Plan

- 1 Introduction
- 2 Les réseaux de neurones
- 3 Propriétés de la fonction objectif

Introduction

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

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

Pour une fonction coût $L : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$, on cherche

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

où

$$R(g) = \mathbb{E}(L(g(X), Y))$$

Introduction

- **Classification** : $\mathcal{Y} = \{1, \dots, n_{out}\}$, pour $n_{out} \geq 2$

On construit

$$f: \mathcal{X} \longrightarrow \mathbb{R}^{n_{out}}$$

on retourne

$$g(x) \in \operatorname{argmax}_{i=1..n_{out}} f_i(x)$$

Typiquement :

$$L(y, y') = \mathbb{1}_{y \neq y'}$$

- **Régression** : $\mathcal{Y} = \mathbb{R}^{n_{out}}$, pour $n_{out} \geq 1$

On construit

$$g: \mathcal{X} \longrightarrow \mathbb{R}^{n_{out}}$$

Typiquement :

$$L(y, y') = \|y - y'\|^2$$

Introduction

- On observe $(x_i, y_i)_{i=1..n}$ un échantillon i.i.d. de même loi que (X, Y)
- On considère une famille \mathcal{F} de fonctions $g_{\mathbf{w}}$ paramétrées par des paramètres \mathbf{w} d'un espace Euclidien.
 - ▶ Dans cet exposé : les **réseaux de neurones**
- On cherche à résoudre

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

Une approche (naïve) consiste à résoudre

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

où

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

Introduction

ATTENTION : Cette modélisation **ne** tient **pas** compte :

- du **biais** dans les données: $(x_i, y_i)_{i=1..n}$ est rarement i.i.d. selon (X, Y)
 - ▶ x_i peut ne couvrir que partiellement le domaine, être corrompue, etc
 - ▶ y_i peut être une décision biaisée, bruitée, etc
- du besoin de **robustesse**
 - ▶ On impose que y_i soit prédit pour tout $x_i + e$, où $\|e\| \leq \varepsilon$
- d'une **régularisation** du réseau
 - ▶ w est parfois **quantifié** \implies **interprétabilité**, faible **coût calculatoire**
 - ▶ w est tel que g_w est **Lipschitz** \implies garantie de **robustesse**
 - ▶ **Régularisations explicites** : poids parcimonieux, pénalisation quadratique, dropout
 - ▶ Il peut y avoir une **régularisation implicite**
- Problèmes plus complexes. Ex: "embedding de mots"

Introduction

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

$$R(g_{\mathbf{w}^*}) \leq \inf_{\mathbf{w}} R(g_{\mathbf{w}}) + \varepsilon \quad \text{et} \quad \hat{R}(g_{\hat{\mathbf{w}}}) \leq \inf_{\mathbf{w}} \hat{R}(g_{\mathbf{w}}) + \varepsilon$$

- Pour \mathbf{w} retourné par un algorithme

On **décompose**

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

Introduction : Erreur de généralisation

- **Convergence uniforme** : Quelles conditions sur le réseau et n garantissent que

$$\widehat{R}(g_{\mathbf{w}}) \simeq R(g_{\mathbf{w}}),$$

pour tout \mathbf{w} ?

- ▶ N. Harvey, C. Liaw, A. Mehrabian. "Nearly-tight VC-dimension and Pseudodimension Bounds for Piecewise Linear Neural Networks." COLT 2017.
 - ★ VC-dimension des réseaux de neurones ReLU feedforward à H couches, W paramètres est $O(WH \log(W))$.
- ▶ P. Bartlett, D. Foster, M. Telgarsky. "Spectrally-normalized margin bounds for neural networks." NeurIPS 2017.
- ▶ N. Golowich, A. Rakhlin, O. Shamir. "Size-independent sample complexity of neural networks." COLT 2018.
- ▶ etc

N'expliquent pas les performances des réseaux de neurones lorsque $n \leq W$

Introduction : Erreur d'approximation

- Est-ce que

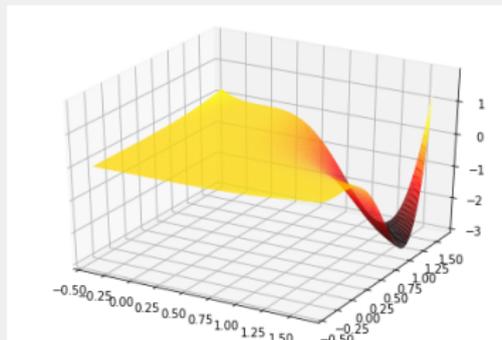
$$\inf_{\mathbf{w}} R(g_{\mathbf{w}}) \simeq \inf_g R(g)?$$

- Quelle architecture de réseau permet d'approximer une certaine classe de fonctions, avec une précision ε ?

Voir l'exposé de **Rémi Gribonval**, à 14h.

Introduction : Erreurs d'optimisation

- le calcul de la fonction objectif dépend de **beaucoup d'entrées** (n grand)
- La fonction objectif est souvent **non-différentiable**
- **w** vit dans un espace de **grande dimension**
- les **propriétés de la fonction objectif** ne se prêtent pas forcément à l'optimisation. On parle de **paysage de la fonction objectif**.



Plan

1

Introduction

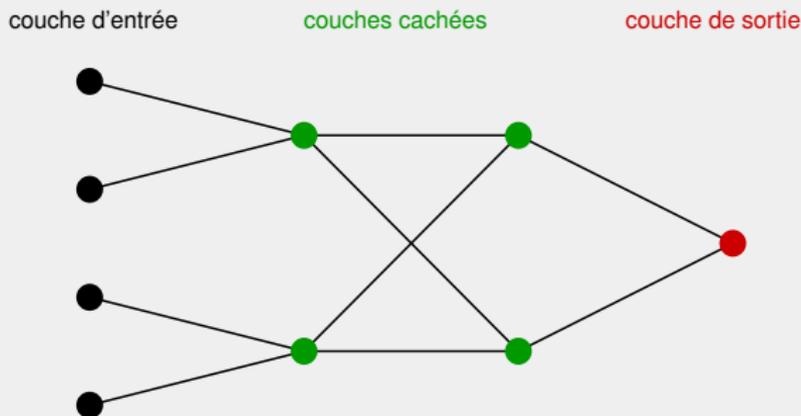
2

Les réseaux de neurones

3

Propriétés de la fonction objectif

La structure Feedforward/Le perceptron



- 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

La structure Feedforward/Le perceptron

- H couches (on dit aussi $H - 1$ couches cachées)
- les tailles $n_0, n_1, \dots, n_H \in \mathbb{N}$ 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 : la fonction d'activation appliquée à chaque couche
 - ▶ elle applique la même fonction à chaque entrée d'un vecteur
 - ▶ dépend de h , en général

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

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

La structure Feedforward/Le perceptron: Exemple

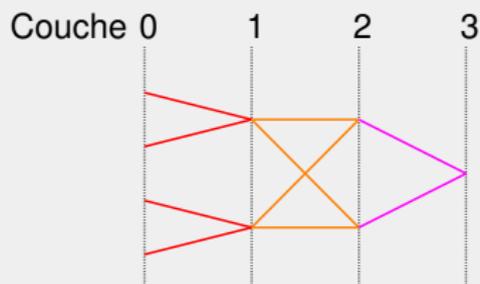


Figure: Feedforward

Ci-dessous, chaque * est un nombre

$$W_1 = \begin{pmatrix} * & * & 0 & 0 \\ 0 & 0 & * & * \end{pmatrix}, b_1 = \begin{pmatrix} * \\ * \end{pmatrix}, W_2 = \begin{pmatrix} * & * \\ * & * \end{pmatrix}, b_2 = \begin{pmatrix} * \\ * \end{pmatrix}, W_3 = \begin{pmatrix} * & * \end{pmatrix},$$
$$b_3 = \begin{pmatrix} * \end{pmatrix},$$

$$f_3(x) = \sigma(W_3 \sigma(W_2 \sigma(W_1 x + b_1) + b_2) + b_3)$$

La structure Feedforward/Le perceptron: Vocabulaire et variantes

- **Fully-connected** : Tous les facteurs sont “pleins”
- **Fonctions d’activation** :
 - ▶ Dépendent de la couche
 - ▶ Une zoologie importante (tanh, heavyside, identité, sigmoïde, etc)
 - ▶ La plus utilisée est ReLU : $\sigma(t) = \max(0, t)$
 - ▶ Des non-locales : group-sort
 - ▶ Souvent, celle associée à la dernière couche est l’identité.
- **“Batch-normalisation”** : Au lieu d’optimiser les biais \mathbf{b} , on règle un couple $(\text{diag}(\gamma), \mathbf{b})$ permettant centrer les données et fixer leur variance.

Réseaux convolutifs

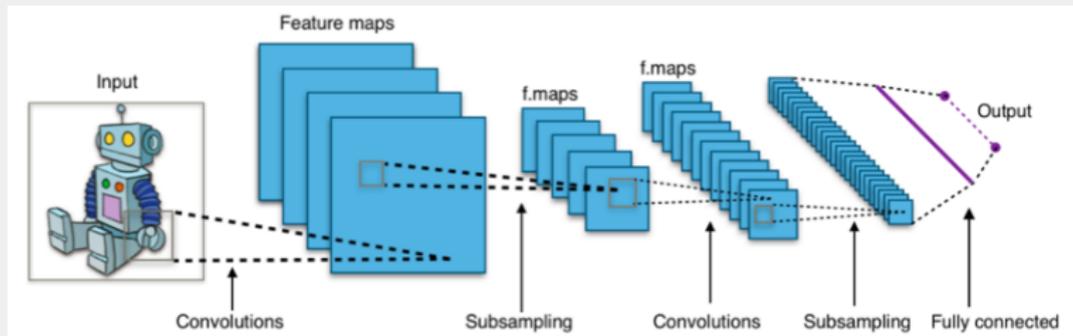
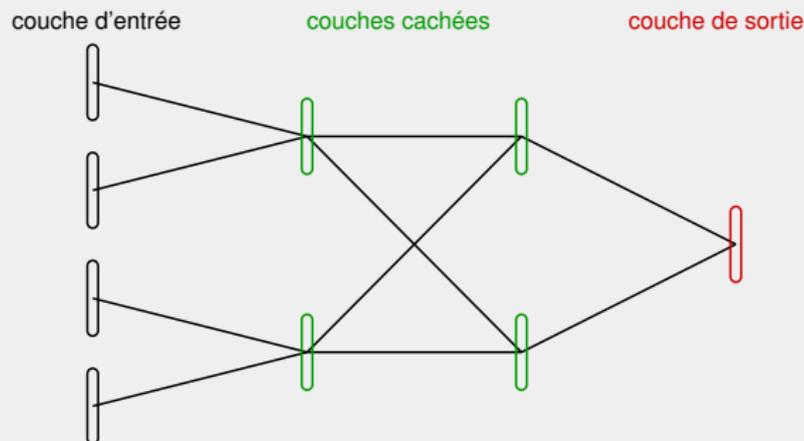


Figure: Source : wikipedia page "apprentissage profond"

Réseaux convolutifs



- Chaque sommet contient un signal/une image. On parle de **canal**.
- Les canaux de la couche d'entrée correspondent aux canaux 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 (= $N \times$ nombre de canal)
- 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 : la fonction d'activation appliquée à chaque couche
 - ▶ elle applique la même fonction à chaque entrée d'un vecteur
 - ▶ dépend de h , en général

L'action du réseau : la fonction f_H

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

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

On suppose $v, x \in \mathbb{R}^{N+1}$. Le signal v est supposé $(N+1)$ -périodique de période:

$$v_{-n'} = v_{N+1-n'} \quad \text{pour tout } n' = 0..N$$

On a:

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

Réseaux convolutifs: Exemple

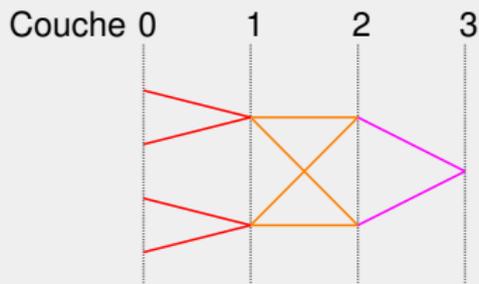


Figure: Réseau convolutif

Ci-dessous,

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

$$W_1 = \begin{pmatrix} * & * & 0 & 0 \\ 0 & 0 & * & * \end{pmatrix}, b_1 = \begin{pmatrix} \circ \\ \circ \end{pmatrix}, W_2 = \begin{pmatrix} * & * \\ * & * \end{pmatrix}, b_2 = \begin{pmatrix} \circ \\ \circ \end{pmatrix}, W_3 = \begin{pmatrix} * & * \end{pmatrix}, b_3 = \begin{pmatrix} \circ \end{pmatrix},$$

$$f_3(x) = \sigma \left(W_3 \sigma \left(W_2 \sigma \left(W_1 x + b_1 \right) + b_2 \right) + b_3 \right)$$

Réseaux convolutifs: Propriétés

Remarque

Les formules sont les mêmes que dans le cas feedforward. Seules les façons de construire les W_h et $f_0(x)$ diffèrent.

Un réseau convolutif est un réseau feedforward dans lequel on impose une structure.

Il existe beaucoup de variantes:

- Une topologie contenant plusieurs type de structures (convolutif, feed-forward ...)
- Couches "Res-net":

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

- "Max pooling"
- Recurrent neural networks
- Recursive neural networks
- etc

Les combinaisons de réseaux

- **Autoencoder** : Un réseau compresse, un réseau décompresse
- **Création d'embeddings** : Deux réseaux envoient une image et un texte dans un même espace de caractéristiques. Le mapping permet d'interroger le contenu d'images.
- **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 réseau classifiant sur un critère pertinent fonctionne, un réseau classifiant sur un critère non-pertinent échoue.
- **etc**

Plan

1 Introduction

2 Les réseaux de neurones

3 Propriétés de la fonction objectif

- Introduction

- Paysage de la fonction objectif

Outline

- 1 Introduction
- 2 Les réseaux de neurones
- 3 Propriétés de la fonction objectif
 - Introduction
 - Paysage de la fonction objectif
 - Paysage pour les réseaux larges
 - Paysage pour les réseaux linéaires

Premières propriétés

- On dispose d'un échantillon $(x_i, y_i)_{i=1..n}$.
- On dispose d'un réseau de profondeur H , de topologie fixé, dont
 - ▶ les arcs sont paramétrés par $\mathbf{w} \in \mathbb{R}^{H \times S}$
 - ▶ les nœuds sont paramétrés par $\mathbf{b} \in \mathbb{R}^{n_1 + \dots + n_H}$
 - ▶ la fonction d'activation est notée σ
- On connaît : $W_h : \mathbb{R}^S \longrightarrow \mathbb{R}^{n_h \times n_{h-1}}$
- La prédiction est la fonction

$$f_{\mathbf{w}, \mathbf{b}}(x) = \sigma \left(W_H(\mathbf{w}_H) \cdots \sigma \left(W_2(\mathbf{w}_2) \sigma \left(W_1(\mathbf{w}_1)x + b_1 \right) + b_2 \right) \cdots + b_H \right)$$

ou, dans le cas feedforward,

$$f_{\mathbf{w}, \mathbf{b}}(x) = \sigma \left(\mathbf{W}_H \cdots \sigma \left(\mathbf{W}_2 \sigma \left(\mathbf{W}_1 x + b_1 \right) + b_2 \right) \cdots + b_H \right)$$

- Cas de la **régression**
- On considère une fonction objectif définie par

$$E : \mathbb{R}^{H \times S} \times \mathbb{R}^{n_1 + \dots + n_H} \longrightarrow \mathbb{R}$$
$$(\mathbf{w}, \mathbf{b}) \longmapsto E(\mathbf{w}, \mathbf{b}) = \sum_{i=1}^n L(f_{\mathbf{w}, \mathbf{b}}(x_i) - y_i)$$

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

Premières propriétés

Proposition: Non-coercivité dans le cas ReLU

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

Preuve utilisant l'homogénéité:

On considère $\mathbf{w} \in \mathbb{R}^{H \times S}$ et, pour tout $\lambda > 0$, on définit $\mathbf{w}^\lambda \in \mathbb{R}^{H \times S}$ par

$$\mathbf{w}_1^\lambda = \lambda^{H-1} \mathbf{w}_1 \quad \text{et} \quad \mathbf{w}_h^\lambda = \lambda^{-1} \mathbf{w}_h, \forall h = 2..H$$

Comme pour $\sigma(\lambda t) = \lambda \sigma(t)$ et, pour tout h , $W_h(\lambda \mathbf{w}_h) = \lambda W_h(\mathbf{w}_h)$ on a

$$f_{\mathbf{w},0}(x) = f_{\mathbf{w}^\lambda,0}(x) \quad , \text{ pour tout } x \text{ et tout } \lambda > 0.$$

Donc $E(\mathbf{w}, 0) = E(\mathbf{w}^\lambda, 0)$, pour tout $\lambda > 0$.

Donc E n'est pas coercive.

Premières propriétés

Proposition: Structure de la fonction objectif: le cas ReLU

Pour tout réseau, avec le fonction d'activation ReLU, pour tout $x \in \mathbb{R}^{n_0}$, pour tout $j = 1 \cdots n_H$, la fonction

$$\begin{aligned} \mathbb{R}^{H \times S} \times \mathbb{R}^{n_1 + \cdots + n_H} &\longrightarrow \mathbb{R} \\ (\mathbf{w}, \mathbf{b}) &\longmapsto [f_{\mathbf{w}, \mathbf{b}}(x)]_j \end{aligned}$$

est **continue**. C'est un **polynôme de degré H par morceaux**.

Il y a au plus $2^{n_1 + \cdots + n_H}$ morceaux.

Si L est quadratique, la fonction E est **continue et c'est un polynôme de degré $2H$ par morceaux**.

Ex.: Le bord $\partial C = \overline{C} - \overset{\circ}{C}$ d'un morceau C est toujours de mesure nulle.

La Rétropropagation

On suppose:

- L est C^1 .
- On suppose que σ est C^1

Sous ces hypothèses: E **est différentiable**.

- On ne considère qu'un unique exemple (x, y) :
 - ▶ On somme si besoin les gradients
 - ▶ Cas de l'algorithme du gradient stochastique

Rétropropagation: cas feedforward complètement connecté

On note $f_{\mathbf{w},\mathbf{b}}^h(x)$, le contenu de la couche h , et

$$d_{\mathbf{w},\mathbf{b}}^h(x) = \sigma'(\mathbf{W}_h f_{\mathbf{w},\mathbf{b}}^{h-1}(x) + \mathbf{b}_h)$$

Proposition: Gradient pour un réseau feedforward

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

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

$$\frac{\partial E}{\partial \mathbf{b}_{h,i}}(\mathbf{w}, \mathbf{b}) = \left[d_{\mathbf{w},\mathbf{b}}^h(x) \right]_i \Delta_i^h(x)$$

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

$$\Delta^h(x) = \begin{cases} \nabla L(f_{\mathbf{w},\mathbf{b}}(x) - y) & , \text{ si } h = H \\ \mathbf{W}_{h+1}^T \cdot \text{diag} \left(d_{\mathbf{w},\mathbf{b}}^{h+1}(x) \right) \cdot \Delta^{h+1}(x) & , \text{ sinon} \end{cases}$$

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

Rétropropagation: Problèmes connus

- “Vanishing/Exploding gradient” :

- ▶ Si $\mathbf{W}_{h+1}^T \cdot \text{diag}(d_{\mathbf{w},\mathbf{b}}^{h+1}(x))$ est systématiquement une contraction \implies “Vanishing gradient”
- ▶ Si $\mathbf{W}_{h+1}^T \cdot \text{diag}(d_{\mathbf{w},\mathbf{b}}^{h+1}(x))$ augmente systématiquement la norme \implies “Exploding gradient”

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

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

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

$$\mathbf{W}_1^\lambda = \lambda^{H-1} \mathbf{W}_1 \quad \text{et} \quad \mathbf{W}_h^\lambda = \lambda^{-1} \mathbf{W}_h, \forall h = 2..H$$

$$\mathbf{b}_1^\lambda = \lambda^{H-1} \mathbf{b}_1 \quad \text{et} \quad \mathbf{b}_h^\lambda = \lambda^{-h} \mathbf{b}_h, \forall h = 2..H$$

On a, pour tout x ,

$$f_{\mathbf{w}^\lambda, \mathbf{b}^\lambda}(x) = f_{\mathbf{w}, \mathbf{b}}(x)$$

mais (génériquement) $\frac{\partial E}{\partial \mathbf{W}_2}$ est grand, la constante Lipschitz de $\frac{\partial E}{\partial \mathbf{W}_1}$ est très faible en \mathbf{W}^λ .

¹ Attention à la non-différentiabilité

Outline

- 1 Introduction
- 2 Les réseaux de neurones
- 3 **Propriétés de la fonction objectif**
 - Introduction
 - **Paysage de la fonction objectif**
 - Paysage pour les réseaux larges
 - Paysage pour les réseaux linéaires

Paysage de la fonction objectif : Introduction

Rappel, pour tout² \mathbf{w} :

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

On utilise un algorithme d'optimisation pour trouver un \mathbf{w}^{calc} , on veut que

$$\widehat{R}(f_{\mathbf{w}^{calc}}) - \inf_{\mathbf{w}} \widehat{R}(f_{\mathbf{w}})$$

soit le plus faible possible.

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

²Pour alléger les notations, on oublie \mathbf{b} .

Paysage de la fonction objectif : Introduction

On suppose que $\mathbf{w} \mapsto \widehat{R}(f_{\mathbf{w}})$ est C^2 partout. On a

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

On distingue:

- **\mathbf{w}^* est un minimiseur global:**

- ▶ $\forall \mathbf{w}, \quad \widehat{R}(f_{\mathbf{w}^*}) \leq \widehat{R}(f_{\mathbf{w}})$
- ▶ $\widehat{R}(f_{\mathbf{w}^*}) = \min_{\mathbf{w}} \widehat{R}(f_{\mathbf{w}})$

- **\mathbf{w}^* est un minimiseur local:**

- ▶ Il existe un voisinage ouvert \mathcal{O} de \mathbf{w}^* tel que

$$\forall \mathbf{w} \in \mathcal{O}, \quad \widehat{R}(f_{\mathbf{w}^*}) \leq \widehat{R}(f_{\mathbf{w}})$$

- **\mathbf{w}^* est un point critique du second ordre:**

- ▶ On a

$$\nabla \widehat{R}(f_{\mathbf{w}^*}) = 0 \quad \text{et} \quad \nabla^2 \widehat{R}(f_{\mathbf{w}^*}) \geq 0$$

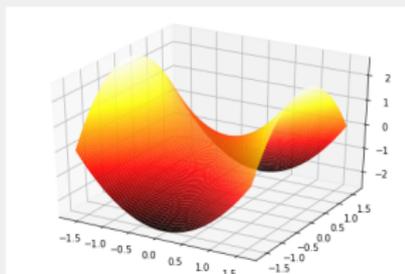
- **\mathbf{w}^* est un point critique du premier ordre:**

- ▶ On a

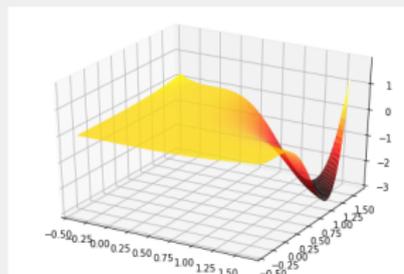
$$\nabla \widehat{R}(f_{\mathbf{w}^*}) = 0$$

Paysage de la fonction objectif : Introduction

- **w^* est un point selle** si c'est un point critique qui n'est ni un minimiseur local, ni un maximiseur local
 - ▶ **Un point selle w^* 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 w^* 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



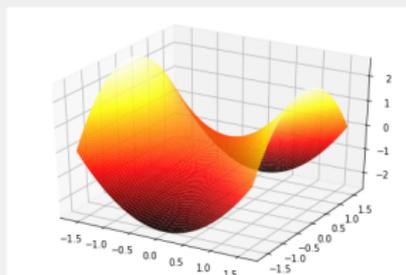
(b) Point selle non-strict

Paysage de la fonction objectif : Introduction

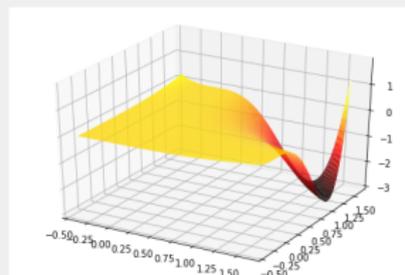
Optimisation non convexe avec les mains

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**



(a) Point selle strict



(b) Point selle non-strict

- 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.

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 feedforward complètement connecté
- des observations $(x_i, y_i)_{i=1..n}$

Paysage pour les réseaux larges

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 (\mathbf{w}, \mathbf{b}) de \hat{R} .

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

Alors, on a

$$\hat{R}(\mathbf{f}_{\mathbf{w}, \mathbf{b}}) = 0$$

et (\mathbf{w}, \mathbf{b}) est un minimum global.

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

$$n \leq n_0 + 1 \quad \text{et} \quad n_H \leq n_{H-1} \leq \cdots \leq n_0$$

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.

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", NIPS 2016
 - (Notamment dans les groupes de S. Arora à Princeton, de F. Bach à l'ENS)
 - En collaboration avec **E. Achour** (Polytechnique, M2 MVA, actuellement en 2^{de} année de thèse), **S. Gerchinovitz** (IRT Saint Exupéry)
-
- un problème de régression
 - un réseau feedforward complètement connecté, linéaire (notamment sans biais)
 - des observations $(x_i, y_i)_{i=1..n}$

Paysage pour les réseaux linéaires

On note

- $X \in \mathbb{R}^{d_x \times n}$ et $Y \in \mathbb{R}^{d_y \times n}$ les matrices contenant les données.

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

-

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

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

-

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

sa SVD

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

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

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

Paysage pour les réseaux linéaires

On suppose

- $d_y \leq d_x \leq n$
- 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} un point critique du premier ordre de \hat{R} . On pose $r = rk(W_H \cdots W_1)$.
Il existe un unique sous-ensemble $\mathcal{S} \subset \llbracket 1, d_y \rrbracket$ 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 $\hat{R}(\mathbf{W}) = \text{tr}(\Sigma_{YY}) - \sum_{i \in \mathcal{S}} \delta_i^2$.
On dit que le point critique \mathbf{W} est associé à \mathcal{S} .

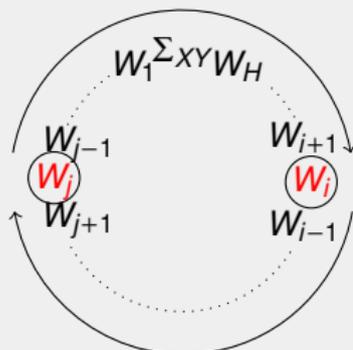
Proposition

Pour tout $\mathcal{S} \subset \llbracket 1, d_y \rrbracket$ de taille $r \in \llbracket 0, r_{max} \rrbracket$, il existe un point critique \mathbf{W} associé à \mathcal{S} .

Paysage pour les réseaux linéaires

- **Pivot** $(i, j) \in \llbracket 1, H \rrbracket^2$, 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 étroit pour W** : L'un des blocs complémentaires est de rang $\text{rk}(W_H \cdots W_1)$
- **W est point critique étroit**: W est un point critique et tous les pivots sont étroits pour W

Paysage pour les réseaux linéaires

W p. c. du premier ordre de \hat{R}

$$r := \text{rk}(W_H \cdots W_1)$$

$\exists! \mathcal{S} \subset \llbracket 1, d_Y \rrbracket$ de taille r t.q. $W_H \cdots W_1 = U_{\mathcal{S}} U_{\mathcal{S}}^T \Sigma_{YX} \Sigma_{XX}^{-1}$ et $\hat{R}(W) = \text{tr}(\Sigma_{YY}) - \sum_{i \in \mathcal{S}} \delta_i^2$

$$r = r_{\max}$$

On regarde \mathcal{S}

$$\mathcal{S} = \llbracket 1, r_{\max} \rrbracket$$

W est un minimiseur global

$$\mathcal{S} \neq \llbracket 1, r_{\max} \rrbracket$$

W est un p.s. strict

$$r < r_{\max}$$

W est un point selle

$$\mathcal{S} \neq \llbracket 1, r \rrbracket$$

W est un p.s. strict

$$\mathcal{S} = \llbracket 1, r \rrbracket$$

$$W_H \cdots W_1 = \text{argmin}_{\text{rk}(R) \leq r} \|RX - Y\|^2$$

W non étreint

W est un p.s. strict

W étreint

W est un p.s. NON strict

Figure: Theorem [E. Achour et al.]

Paysage pour les réseaux linéaires

Proposition (E. Achour et al.)

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

Proposition (E. Achour et al.)

On note $\mathcal{S}_{max} = \llbracket 1, r_{max} \rrbracket$ et $Q_{max} = \llbracket 1, d_y \rrbracket \setminus \mathcal{S}_{max} = \llbracket r_{max} + 1, d_y \rrbracket$.

W est un minimiseur global de \widehat{R} si et seulement si il existe des matrices inversibles $D_{H-1} \in \mathbb{R}^{d_{H-1} \times d_{H-1}}, \dots, D_1 \in \mathbb{R}^{d_1 \times d_1}$, et des matrices $A_R \in \mathbb{R}^{(d_y - r_{max}) \times (d_{H-1} - r_{max})}$, $(W_h)_{DR} \in \mathbb{R}^{(d_h - r_{max}) \times (d_{h-1} - r_{max})}$ pour $h \in \llbracket 2, H-1 \rrbracket$, et $M_D \in \mathbb{R}^{(d_1 - r_{max}) \times d_x}$ tels que:

$$W_H = [U_{\mathcal{S}_{max}}, U_{Q_{max}} A_R] D_{H-1}^{-1}$$

$$W_h = D_h \begin{bmatrix} I_{r_{max}} & 0 \\ 0 & (W_h)_{DR} \end{bmatrix} D_{h-1}^{-1} \quad \forall h \in \llbracket 2, H-1 \rrbracket$$

$$W_1 = D_1 \begin{bmatrix} U_{\mathcal{S}_{max}}^T & \Sigma_{YX} \Sigma_{XX}^{-1} \\ & M_D \end{bmatrix}.$$

Paysage pour les réseaux linéaires

Conclusion [E. Achour et al.]

- **Caractérisation** de l'ensemble: des 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**

Autres travaux sur les réseaux de neurones :

- Condition d'identifiabilité et de stabilité de \mathbf{w} : SIAM data science, MSML, ITW, SPARS et thèse de **J. Bona-Pellissier** en cours
- Apprentissage certifié de fonction monotone : AAI, workshop safeAI

Merci pour votre attention !