

17.4 Programmation différentiable et apprentissage d'EDO

Ici on va introduire une alternative pour l'apprentissage des réseaux de type OdeNet. Cela va nous permettre d'introduire la programmation différentiable qui va bien au delà de l'apprentissage d'EDO.

17.4.1 OpD vs DpO

L'approche OdeNet consiste à écrire le gradient d'une fonction coût qui dépend de la solution au temps final ou a un ensemble de temps par rapport aux paramètres de l'EDO. En pratique cela fait intervenir le gradient de la solution par rapport aux paramètres. Sachant que la fonction qui relie les paramètres à la solution n'étant pas explicite le calcul de gradient est complexe. Pour effectuer ce calcul les OdeNet propose de calculer le gradient au niveau continu avec ce qu'on appelle la **La rétro-propagation continue** (celle introduite au chapitre précédent) ou **la méthode d'adjoint**. Plus globalement on parle de méthode de type **Optimiser puis Discrétiser** car le calcul de gradient se fait au niveau continu de l'EDO avant qu'on utilise un flot numérique. Une autre approche consiste à **Discrétiser puis Optimiser**. Il s'agit de discrétiser l'EDO puis de calculer son gradient. Une fois discrétisée et codée la résolution d'une EDO revient à un programme informatique enchaînant les fonctions simples. Dans ce cadre pour calculer le gradient de la solution par rapport aux paramètres il faut être capable de différencier automatiquement à travers notre programme. Pour cela l'outil essentiel est la **différentiation automatique** qui est l'outil caché derrière les logiciels d'apprentissage comme Pytorch, Tensorflow ou Jax. Cela permet de calculer la dérivée automatiquement une succession de fonction dans un programme. Coder de façon à ce que cette différentiation automatique soit utilisable dans l'ensemble du code s'appelle la **Programmation différentiable**. Cela permet notamment d'optimiser n'importe quelle partie du code ou d'étudier sa sensibilité aux entrées.

On va dans un premier introduire les concepts importants pour la différentiation automatique de programme puis on discutera plus précisément le cas des EDO. Ces notes sont principalement basés sur le superbe document [\[1.45\]](#).

17.4.2 Rappels d'algèbre et d'analyse

On va commencer par des rappels d'outils sur la différentiation et d'algèbre linéaires fondamentaux pour ce chapitre.

Définition 17.37. Dérivées directionnelle et partielle d'une fonction scalaire. On se donne une fonction $f : \mathbb{R}^d \rightarrow \mathbb{R}$, la dérivée directionnelle de f en \mathbf{x} dans la direction \mathbf{v} est donnée par

$$\partial f(\mathbf{x})[\mathbf{v}] := \lim_{\delta \rightarrow 0} \frac{f(\mathbf{x} + \delta \mathbf{v}) - f(\mathbf{x})}{\delta}$$

sous condition que la limite existe. On parle de dérivée partielle notée $\partial_i f(\mathbf{x})$ si $\mathbf{v} = \mathbf{e}_i$ le i ème vecteur de la base canonique.

Définition 17.38. Gradient. On se donne une fonction $f : \mathbb{R}^d \rightarrow \mathbb{R}$ différentiable en chaque point \mathbf{x} . Le **gradient** est le vecteur des dérivées partielles

$$\nabla f(\mathbf{x}) := \begin{pmatrix} \partial_1 f(\mathbf{x}) \\ \vdots \\ \partial_d f(\mathbf{x}) \end{pmatrix} = \begin{pmatrix} \partial f(\mathbf{x})[\mathbf{e}_1] \\ \vdots \\ \partial f(\mathbf{x})[\mathbf{e}_d] \end{pmatrix}$$

On peut relier la dérivée directionnelle et le gradient par la relation:

$$\partial f(\mathbf{x})[\mathbf{v}] = \sum_{i=1}^d v_i \partial f(\mathbf{x})[\mathbf{e}_i] = \langle \mathbf{v}, \nabla f(\mathbf{x}) \rangle$$

Maintenant on va définir les notions équivalentes pour les fonctions à valeurs vectorielles.

Définition 17.39. Dérivées directionnelle et partielle d'une fonction vectorielle. On se donne une fonction $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^p$ définie par $\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_p(\mathbf{x}))$ la dérivée directionnelle de \mathbf{f} en \mathbf{x} dans la direction \mathbf{v} est donnée par

$$\partial \mathbf{f}(\mathbf{x})[\mathbf{v}] := \lim_{\delta \rightarrow 0} \begin{pmatrix} \frac{f_1(\mathbf{x} + \delta \mathbf{v}) - f_1(\mathbf{x})}{\delta} \\ \vdots \\ \frac{f_p(\mathbf{x} + \delta \mathbf{v}) - f_p(\mathbf{x})}{\delta} \end{pmatrix}$$

sous condition que la limite existe. On parle de dérivée partielle notée $\partial_i \mathbf{f}(\mathbf{x})$ si $\mathbf{v} = \mathbf{e}_i$ le i ème vecteur de la base canonique.

Définition 17.40. Jacobienne. La **Jacobienne** d'une fonction vectorielle

différentiable $f : \mathbb{R}^d \rightarrow \mathbb{R}^p$ en \mathbf{x} est donnée par la matrice des dérivées partielles de chaque fonction scalaire:

$$\partial \mathbf{f}(\mathbf{x}) := \begin{pmatrix} \partial_1 f_1(\mathbf{x}) & \dots & \partial_d f_1(\mathbf{x}) \\ \vdots & \ddots & \vdots \\ \partial_1 f_p(\mathbf{x}) & \dots & \partial_d f_p(\mathbf{x}) \end{pmatrix} \in \mathbb{R}^{p \times d}$$

On a aussi

$$\partial \mathbf{f}(\mathbf{x})[\mathbf{v}] = \sum_{i=1}^q v_i \partial_i f(\mathbf{x}) = \partial \mathbf{f}(\mathbf{x}) \mathbf{v} \in \mathbb{R}^p$$

Le produit entre la Jacobienne et un vecteur \mathbf{v} permet de mesurer la variation de la fonction dans la direction $\mathbf{v} \in \mathbb{R}^d$ d'entrée. On peut aussi mesurer la variation de la fonction dans une direction de sortie $\mathbf{u} \in \mathbb{R}^p$. Cela revient à calculer le gradient de la fonction

$$\langle \mathbf{u}, \mathbf{f} \rangle(\mathbf{x}) := \langle \mathbf{u}, \mathbf{f}(\mathbf{x}) \rangle \in \mathbb{R}$$

En développant \mathbf{u} sur une base $\mathbf{u} = \sum_{j=1}^p u_j \mathbf{e}_j \in \mathbb{R}^p$ on peut, en utilisant la linéarité du gradient par

$$\nabla \langle \mathbf{u}, \mathbf{f} \rangle(\mathbf{x}) = \sum_{j=1}^p u_j \nabla f_j(\mathbf{x}) = \partial \mathbf{f}(\mathbf{x})^\top \mathbf{u} \in \mathbb{R}^d$$

Si on détaille le calcul composante par composante on obtient

$$\nabla_i \langle \mathbf{u}, \mathbf{f} \rangle(\mathbf{x}) = [\partial \mathbf{f}(\mathbf{x})^\top \mathbf{u}]_i = \lim_{\delta \rightarrow 0} \frac{\langle \mathbf{u}, \mathbf{f}(\mathbf{x} + \delta \mathbf{e}_i) - \mathbf{f}(\mathbf{x}) \rangle}{\delta}$$

On voit que la i ème composante $\nabla_i \langle \mathbf{u}, \mathbf{f} \rangle(\mathbf{x})$ revient à la variation dans la direction \mathbf{e}_i projeté dans la direction \mathbf{u} .

La règle de la différentiation en chaîne est un des outils les plus fondamentaux en différentiation automatique. On a déjà vu son utilisation dans le calcul du gradient des réseaux de neurones.

Pour construire nos outils de programmation différentiable on a besoin de pouvoir mesurer les variations infinitésimales le long des directions d'entrées ou de sortie. On a vu rapidement que ces notions étaient encodées par des applications linéaires reliées à la Jacobienne. On va maintenant formaliser ses

notions dans les espaces Euclidiens

Pour les vecteurs on va se servir du produit scalaire usuel. Pour les matrices on se muni du produit scalaire suivant:

$$\langle A, B \rangle_{\mathcal{M}_{n,m}(\mathbb{R})} = \text{Tr}(A^t B).$$

On va maintenant introduire une notion fondamentale pour la suite.

Définition 17.41. Opérateur adjoint d'une application linéaire. On se donne deux espaces Euclidiens $(\mathcal{E}, \langle \cdot, \cdot \rangle_{\mathcal{E}})$ et $(\mathcal{F}, \langle \cdot, \cdot \rangle_{\mathcal{F}})$ l'adjoint d'une application linéaire $l : \mathcal{E} \rightarrow \mathcal{F}$ est l'unique application linéaire $l^* : \mathcal{F} \rightarrow \mathcal{E}$ tel que $\mathbf{v} \in \mathcal{E}$ et $\mathbf{u} \in \mathcal{F}$ satisfasse:

$$\langle l(\mathbf{v}), \mathbf{u} \rangle_{\mathcal{F}} = \langle \mathbf{v}, l^*(\mathbf{u}) \rangle_{\mathcal{E}}$$

Dans le cas $l(\mathbf{v}) = A\mathbf{v}$ on obtient $l^*(\mathbf{u}) = A^t \mathbf{u}$ par dualité du produit scalaire.

On va pouvoir introduire les deux outils qui seront très utilisés dans la suite.

Définition 17.42. Produit Jacobienne-Vecteur (JVP). Pour $\mathbf{f} : \mathcal{E} \rightarrow \mathcal{F}$, l'application linéaire $\partial \mathbf{f}(\mathbf{x}) : \mathcal{E} \rightarrow \mathcal{F}$, qui relie \mathbf{v} à

$$\partial \mathbf{f}(\mathbf{x})[\mathbf{v}] = \partial \mathbf{f}(\mathbf{x}) \mathbf{v}$$

est appelé **Produit Jacobienne Vecteur (JVP)**. La fonction:

$$\partial \mathbf{f} : \mathcal{E} \rightarrow (\mathcal{E} \rightarrow \mathcal{F})$$

est une fonction qui va \mathcal{E} à l'espace des applications linéaire $\mathcal{L}(\mathcal{E}, \mathcal{F})$.

Définition 17.43. Produit Vecteur-Jacobienne (VJP). Pour $\mathbf{f} : \mathcal{E} \rightarrow \mathcal{F}$, l'adjoint $\partial \mathbf{f}^*(\mathbf{x}) : \mathcal{F} \rightarrow \mathcal{E}$, du JVP est appelé **Produit Vecteur-Jacobienne (VJP)**. Il satisfait:

$$\nabla \langle \mathbf{u}, \mathbf{f} \rangle_{\mathcal{F}}(\mathbf{w}) = \partial \mathbf{f}^*(\mathbf{w}) \mathbf{u}$$

avec $\langle \mathbf{u}, \mathbf{f} \rangle_{\mathcal{F}}(\mathbf{x}) := \langle \mathbf{u}, \mathbf{f}(\mathbf{x}) \rangle_{\mathcal{F}}$ La fonction:

$$(\partial \mathbf{f})^* : \mathcal{E} \rightarrow (\mathcal{F} \rightarrow \mathcal{E})$$

est une fonction qui va \mathcal{E} à l'espace des applications linéaire $\mathcal{L}(\mathcal{F}, \mathcal{E})$.

Maintenant faut étendre les précédentes définitions aux cas où les entrées et

sorties sont multiples. Par exemple pour les réseaux de neurones ce sont des fonctions qui dépendent de leurs entrées et des poids du réseaux.

Définition 17.44. JVP et VJP pour les fonctions à plusieurs entrées.

On considère une fonction différentiable $\mathbf{f}(\mathbf{x}) = \mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_M)$ avec $\mathbf{f} : \mathcal{E} = \mathcal{E}_1 \times \dots \times \mathcal{E}_M \rightarrow \mathcal{F}$, L'application JVP dans la direction $\mathbf{v} = (\mathbf{v}_1, \dots, \mathbf{v}_M) \in \mathcal{E}$ est donnée par

$$\begin{aligned}\partial \mathbf{f}(\mathbf{x})[\mathbf{v}] &= \partial \mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_M) [\mathbf{v}_1, \dots, \mathbf{v}_M] \in \mathcal{F} \\ &= \sum_{i=1}^M \partial_i \mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_M) [\mathbf{v}_i]\end{aligned}$$

L'application VJP dans la direction \mathbf{u} est donnée par

$$\begin{aligned}\partial \mathbf{f}(\mathbf{x})^*[\mathbf{u}] &= \partial \mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_M)^*[\mathbf{u}] \in \mathcal{E} \\ &= (\partial_1 \mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_M)^*[\mathbf{u}], \dots, \partial_M \mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_M)^*[\mathbf{u}])\end{aligned}$$

Définition 17.45. JVP et VJP pour des fonctions a plusieurs sorties.

On considère une fonction différentiable $\mathbf{f}(\mathbf{x}) = (\mathbf{f}_1(\mathbf{w}), \dots, \mathbf{f}_T(\mathbf{w}))$, de la forme $\mathcal{E} \rightarrow \mathcal{F}$ et $\mathbf{f}_i : \mathcal{E} \rightarrow \mathcal{F}_i$ avec $\mathcal{F} := \mathcal{F}_1 \times \dots \times \mathcal{F}_T$. L'application JVP dans la direction $\mathbf{v} \in \mathcal{E}$ est donnée par

$$\partial \mathbf{f}(\mathbf{x})[\mathbf{v}] = (\partial \mathbf{f}_1(\mathbf{x})[\mathbf{v}], \dots, \partial \mathbf{f}_T(\mathbf{x})[\mathbf{v}]) \in \mathcal{F}$$

L'application VJP dans la direction $\mathbf{u} = (\mathbf{u}_1, \dots, \mathbf{u}_T)$ est donnée par

$$\begin{aligned}\partial \mathbf{f}(\mathbf{x})^*[\mathbf{u}] &= \partial \mathbf{f}(\mathbf{x})^*[\mathbf{u}_1, \dots, \mathbf{u}_T] \in \mathcal{E} \\ &= \sum_{i=1}^T \partial \mathbf{f}_i(\mathbf{x})^*[\mathbf{u}_i]\end{aligned}$$

Définition 17.46. Dérivée en chaine. On se donne deux fonctions $\mathbf{f} : \mathcal{E} \rightarrow \mathcal{F}$ et $\mathbf{g} : \mathcal{F} \rightarrow \mathcal{G}$ globalement différentiable. La composition $\mathbf{g} \circ \mathbf{f}$ est différentiable. Le **produit Jacobienne-Vecteur (JVP)** de la composition est donnée par:

$$\partial(\mathbf{g} \circ \mathbf{f})(\mathbf{x}) = \partial \mathbf{g}(\mathbf{f}(\mathbf{x})) \partial \mathbf{f}(\mathbf{x}) \quad (17.13)$$

Le **produit Vecteur-Jacobienne (VJP)** de la composition est donnée par:

$$\partial(\mathbf{g} \circ \mathbf{f})^*(\mathbf{x}) = \partial \mathbf{f}^*(\mathbf{x}) \partial \mathbf{g}^*(\mathbf{f}(\mathbf{x})) \quad (17.14)$$

Si $\mathcal{G} = \mathbb{R}$ on obtient

$$\nabla(\mathbf{g} \circ \mathbf{f})(\mathbf{x}) = \partial \mathbf{f}(\mathbf{x})^\top \nabla \mathbf{g}(\mathbf{f}(\mathbf{x})) \quad (17.15)$$

La dérivée à la chaine est l'outil central de la rétro-propagation et de la différentiation automatique. Cette définition générale sera très utile par la suite.

17.4.3 Programmes paramétrisés et différentiation

17.4.3.1 Programme paramétrisés et représentation

Comme dit précédemment, l'outil central pour la différentiation automatique va être la dérivée en chaine qui permet de dériver une composition de fonction. Avant de l'introduire son application a un programme informatique on va regarder comme on représente un programme mathématiquement.

Si on considère un programme qui est une succession de fonction informatique et donc une composition de fonction mathématique comme sur l'image

Figure 17.47



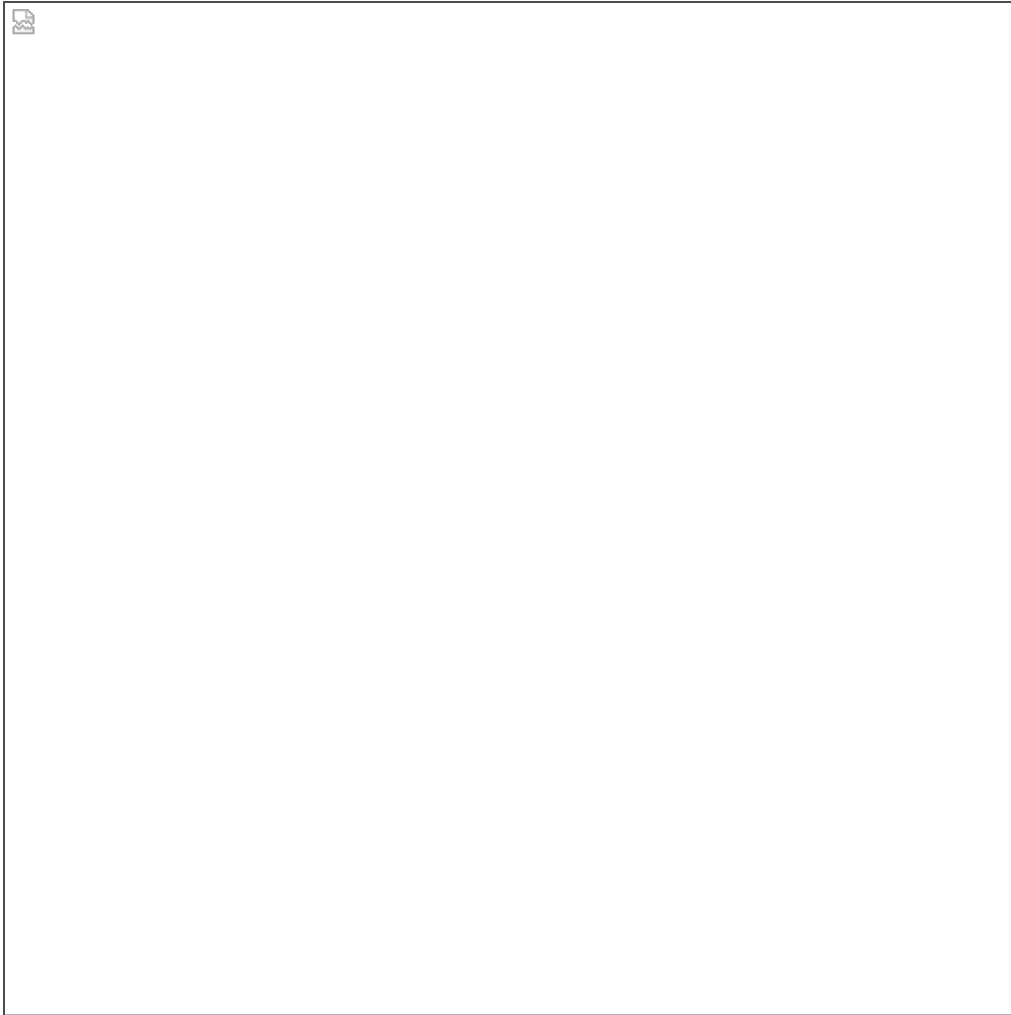


Figure 17.47. Exemple de programme séquentiel. Image issue de [1.45].

En pratique, un programme ne suit pas une structure aussi simple. En effet on peut avoir une fonction a plusieurs entrées dont les données sont calculées par plusieurs fonctions antérieures. On peut donner l'exemple [Figure 17.48](#).

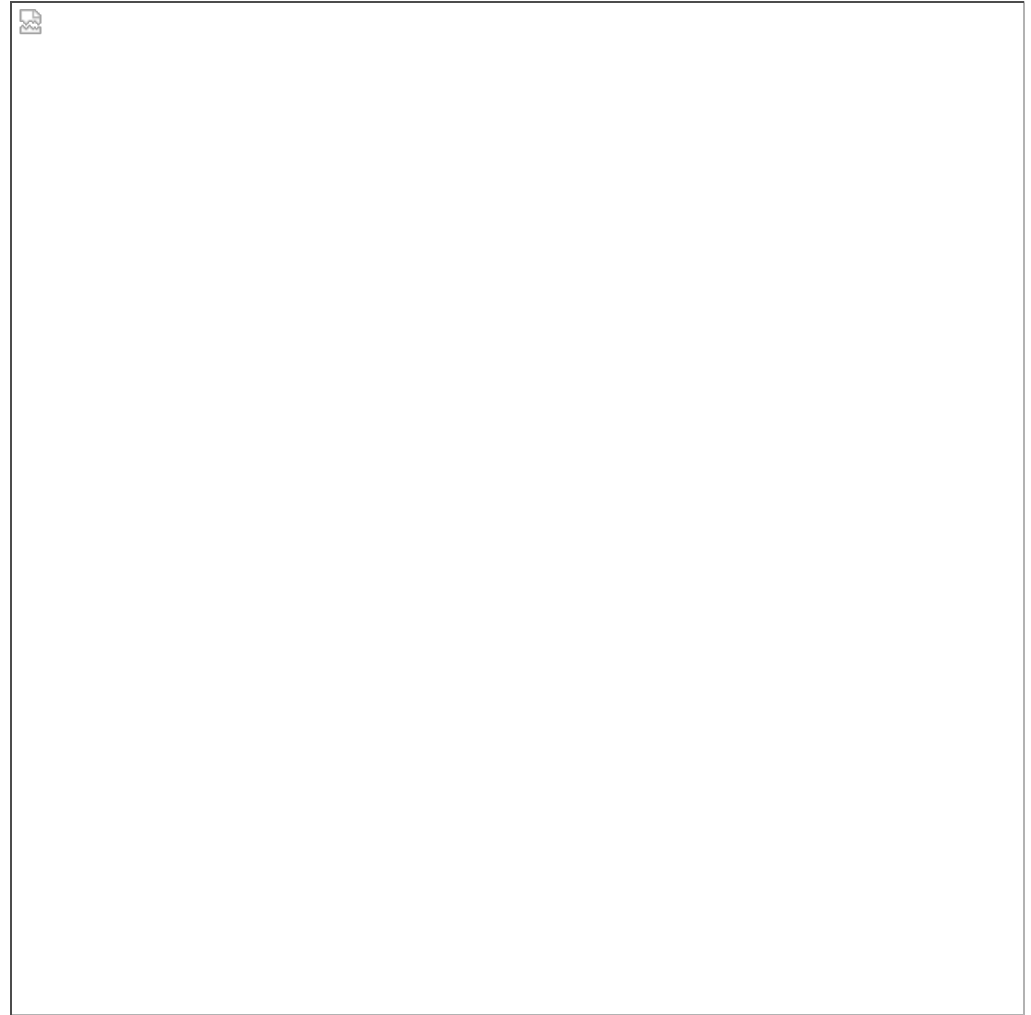


Figure 17.48. Exemple de programme non séquentiel. Image issue de [1.45].

En pratique pour représenter un programme on va utiliser un **graphe orienté acyclique**. Cette structure a été introduite dans le cours de rappel [Section 2.4](#). Chaque noeud représentera une fonction du programme chaque arête (i, j) représentera la dépendance entre les sorties de la fonction i et les entrées de la fonction j . Afin de définir l'ordre dans lequel les fonctions doivent être évaluées on a besoin choisir une relation d'ordre sur les noeuds du graphe. En général on utilise l'ordre **topologique**. Il est défini par $i \leq j$ si il existe aucun chemin de j à i . Cela revient à dire que **chaque sommet apparaît bien avant ses successeurs**. Cela permet d'assurer qu'aucune fonction n'est appelée avant que ses données d'entrées soit calculées. En prenant cette façon de représenter un programme, l'exécution d'un programme revient à l'algorithme suivant.

Algorithm 17.49. Execution d'un programme.

- **Données** : les fonctions f_1, \dots, f_K dans l'ordre topologique du graphe
- **Entrées** : état \mathbf{x}_0 .
- $\forall k = 1 < K$:
 - On récupère la liste des noeuds parents du noeud k :

$$i_1, \dots, i_{p_k} = \text{Parent}(k)$$

- On récupère les données associées $\mathbf{x}_{i_1}, \dots, \mathbf{x}_{p_k}$.
 - On calcul $\mathbf{x}_k = f_k(\mathbf{x}_{i_1}, \dots, \mathbf{x}_{p_k})$.
- On renvoie \mathbf{x}_K .

Ce formalisme peut aussi fonctionner si on a des fonctions qui ont des multiples entrées ou sorties il suffit de les concatener toutes ensemble et ensuite les fonctions filtreront utiles ou non.

Cette description des programmes nous permettra de mettre en place la différentiation automatique. Avant cela on va montrer comme différentier les briques élémentaires des programmes informatiques.

17.4.3.2 Différentiation des structure des contrôles

Les structure de contrôles (instruction conditionnelle et boucle) sont la bases des programmes informatiques, on va donc commencer par regarder comme différentier ses outils.

On va commencer par les opérateurs de comparaison comme $<$. On peut l'écrire comme des opérateur binaire de la forme

$$y = O(x_1, x_2)$$

avec $y \in \{0, 1\}$. Pour donner une définition mathématique de ses opérateurs on va réintroduire la fonction de Heaviside.

Définition 17.50. Fonction de Heaviside. La fonction de Heaviside est donnée par

$$H(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases}$$

La fonction de Heaviside nous permet de redéfinir les opérateurs de

comparaison.

- Opérateur \geq donné par $Ge(x_1, x_2) = H(u_1 - u_2)$
- Opérateur \leq donné par $Le(x_1, x_2) = H(u_2 - u_1)$
- Opérateur $=$ donné par $Eq(x_1, x_2) = Ge(x_1, x_2)Ge(x_2, x_1) = H(u_1 - u_2)H(u_2 - u_1)$
- Opérateur $!$ donné par $Neq(x_1, x_2) = 1.0 - H(u_1 - u_2)H(u_2 - u_1)$

La fonction de Heaviside est différentiable partout sans en zero et sa dérivée est nulle partout. Par conséquent la dérivée ne va pas apporter des informations. Pour palier a ses deux défauts l'idée de remplacer la fonction de heaviside par une fonction régularisée.

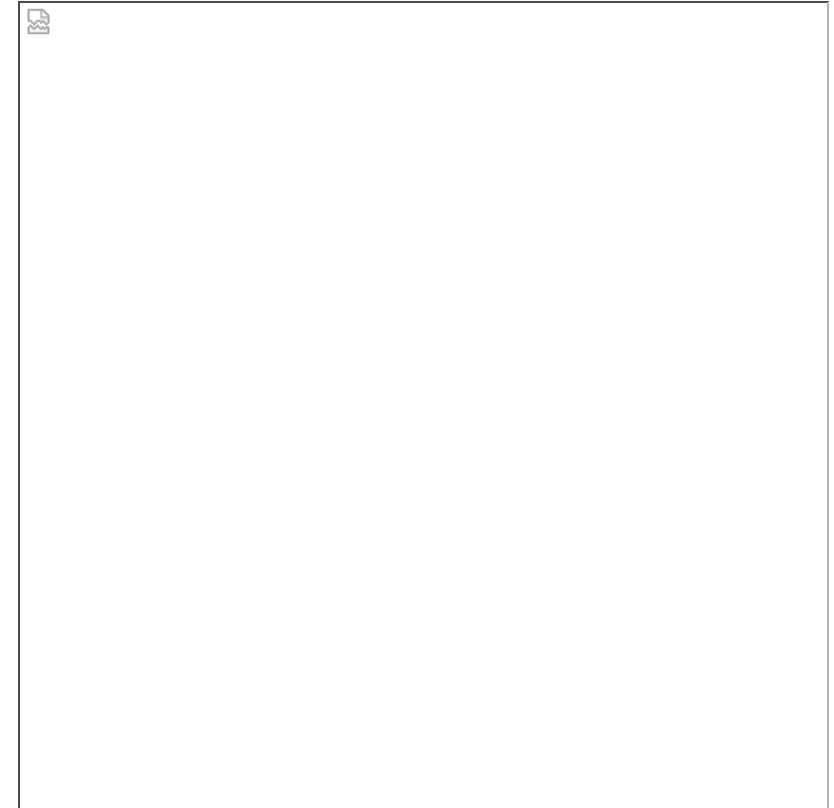


Figure 17.51. fonction sigmoïde pour plusieurs σ

La fonction la plus naturelle est la fonction **sigmoïde** très utilisée en apprentissage. Elle est donnée par

$$\text{sig}_\sigma(x) = \frac{1}{1 + e^{-\frac{x}{\sigma}}}$$

Sur la figure [Figure 17.51](#) on voit que lorsque $\sigma \rightarrow 0$ la fonction sigmoïde tend vers la fonction de Heaviside. En pratique on peut donc définir des opérateurs de comparaison régularisés en remplaçant dans les définitions précédentes par une **fonction sigmoïde**.

On va maintenant introduire l'opérateur d'égalité. On voit que la stratégie pour obtenir un opérateur différentiable on relaxe l'opérateur en lui demandant de renvoyer $y \in [0, 1]$ au lieu de $y \in \{0, 1\}$. La première question est donc de savoir comment relaxer l'opération d'égalité. L'égalité peut être vue comme la limite d'une mesure de **similarité**.

Définition 17.52. Fonction de similarité. une fonction de similarité S entre deux données x_1 et x_2 est une fonction maximale en $x_1 - x_2$. On peut définir une fonction de similarité assez générale à partir de fonction noyau [Définition 3.68](#) par la formule

$$S(x_1, x_2) = \frac{k(x_1, x_2)}{\sqrt{k(x_1, x_1)}\sqrt{k(x_2, x_2)}} \quad (17.16)$$

Si on prend un noyau invariant par translation $k(x - y)$ on obtient

$$S(x_1, x_2) = \frac{k(x_1 - x_2)}{k(0)} \quad (17.17)$$

Afin de régulariser l'opérateur d'égalité on va donc le remplacer par une fonction de similarité S_σ qui va tendre vers l'opérateur égal lorsque $\sigma \rightarrow 0$

► **Exemple 17.53.**

On donne ici deux noyaux qui peuvent définir des fonctions de similarités:

$$k_\sigma(t) = \exp\left(-\frac{t^2}{\sigma}\right)$$

et

$$k_\sigma(t) = \operatorname{sech}\left(\frac{t}{\sigma}\right)$$

avec $\operatorname{sech}(x) = \frac{2}{e^x + e^{-x}}$. On les retrouve sur les figures [Figure 17.54](#) et [Figure 17.55](#).

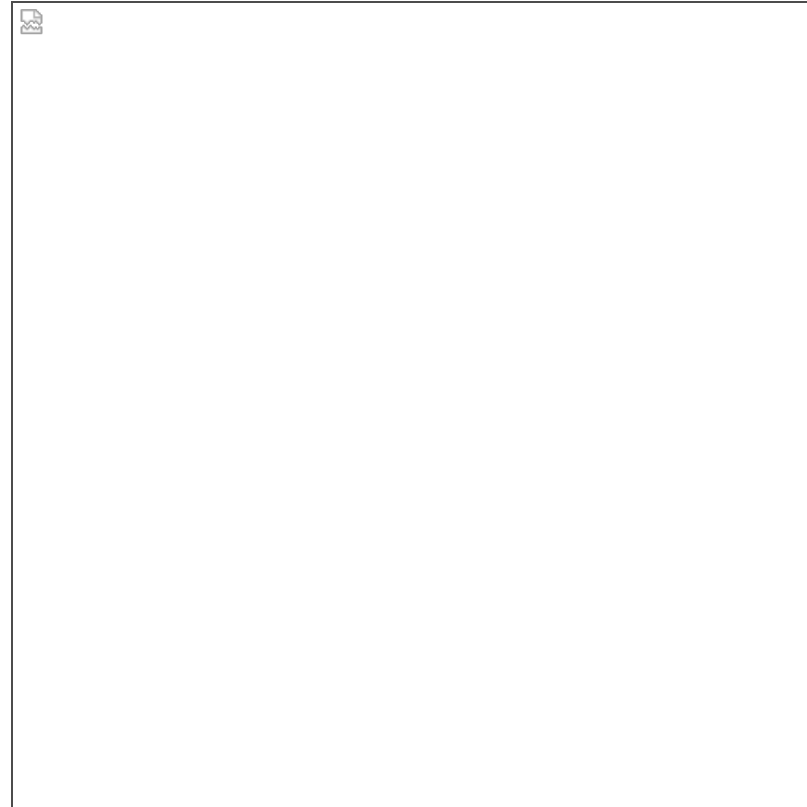


Figure 17.54. fonction de similarité Gaussienne.

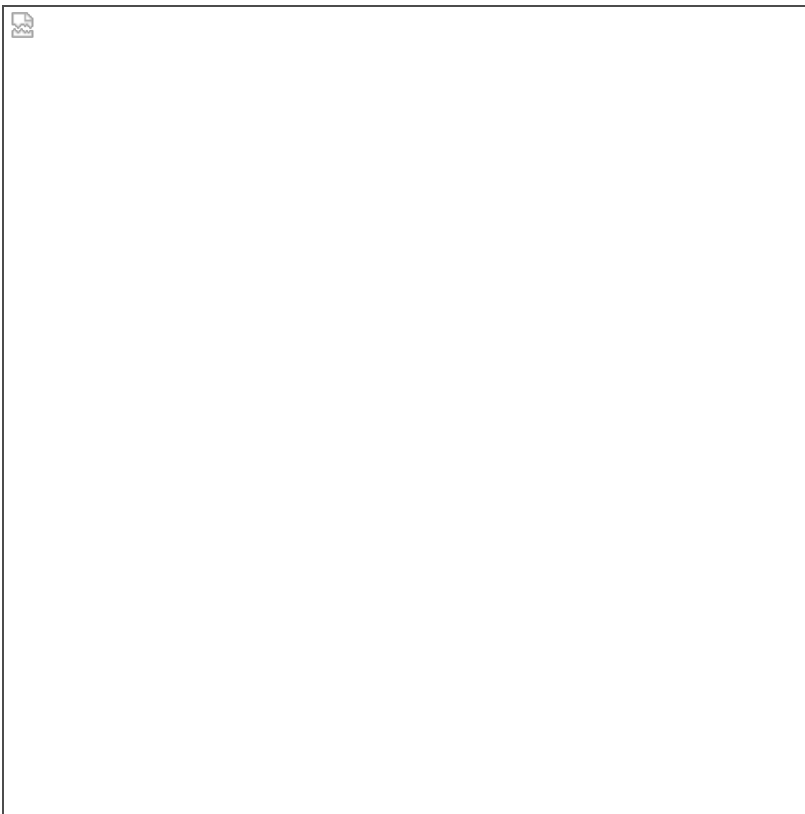


Figure 17.55. fonction de similarité Sech.

Evidemment une fois qu'on a défini une version régularisée de l'opérateur égal par une fonction de similarité de la forme (17.17) on peut définir aussi l'opérateur **différent !=**.

Après les opérateurs de comparaison on va maintenant considérer les opérateurs logiques: **et, ou, non** ainsi que les généralisations au cas non binaire. Il s'agit d'opérateur travaillant sur des booléens $y \in \{0, 1\}$. On va commencer par les rappeler

$$\begin{aligned} \text{and}(y_1, y_2) &:= \begin{cases} 1 & \text{si } y_1 = y_2 = 1 \\ 0 & \text{sinon} \end{cases} \\ \text{or}(y_1, y_2) &:= \begin{cases} 1 & \text{si } 1 \in \{y_1, y_2\} \\ 0 & \text{sinon} \end{cases} \\ \text{not}(y) &:= \begin{cases} 0 & \text{si } y = 1 \\ 1 & \text{si } y = 0 \end{cases} \end{aligned}$$

les opérateurs de comparaison retourne normalement des booléens donc dans

$\{0, 1\}$. Pour dériver les opérateurs de comparaisons on utilise des versions régularisées qui renvoient des doubles dans $[0, 1]$. Cela revient donc à dire que un boolean devient un réel entre $[0, 1]$ très probablement proche de zéro ou de 1. Cela permet d'obtenir facilement des versions régulières de nos opérateurs logiques

Définition 17.56. Opérateurs logiques régularisés. Soit $y_1, y_2 \in [0, 1]$. alors les opérateurs logiques deviennent des opérateurs de la forme:

$$[0, 1] \times [0, 1] \rightarrow [0, 1]$$

définis par

$$\begin{aligned} \text{and}(y_1, y_2) &= y_1 y_2, \quad \text{or}(y_1, y_2) = y_1 + y_2 - y_1 y_2 \\ \text{not}(y) &= 1 - y \end{aligned}$$

et

$$\text{any}(y_1, \dots, y_n) = \prod_{i=1}^n y_i, \quad \text{any} = 1 - \prod_{i=1}^n (1 - y_i)$$

Ses opérateurs de comparaisons et logiques ont des interprétations probabilistes. En effet on remplace nos booléens par des lois de probabilités logistiques. Maintenant qu'on a introduit ses outils fondamentaux de comparaison et logiques on va pouvoir introduire la dérivation des instructions conditionnelles comme mes "if".

une structure conditionnelle "ifelse" est définie par la fonction mathématique $\text{ifelse}(y, \mathbf{x}_1, \mathbf{x}_2) : \{0, 1\} \times E \times E \rightarrow E$ donnée par

$$\begin{aligned} \text{ifelse}(y, \mathbf{x}_1, \mathbf{x}_2) &:= \begin{cases} \mathbf{x}_1 & \text{si } y = 1 \\ \mathbf{x}_2 & \text{si } y = 0 \end{cases} \\ &= y\mathbf{x}_1 + (1 - y)\mathbf{x}_2 \end{aligned}$$

Cette fonction n'est pas différentiable car le comportement dépend d'un booléen. Par contre la fonction est dérivable dans chaque branche. Elle est donc facile de calculer les dérivées par rapport à \mathbf{x}_1 et \mathbf{x}_2 . Elles sont données par

$$\partial_{\mathbf{x}_1} \text{ifelse}(y, \mathbf{x}_1, \mathbf{x}_2) = yI_d, \quad \partial_{\mathbf{x}_2} \text{ifelse}(y, \mathbf{x}_1, \mathbf{x}_2) = (1 - y)I_d$$

Si $\mathbf{x}_1, \mathbf{x}_2$ sont donnés par des fonctions de $\mathbf{x}_1 = \mathbf{f}_1(\mathbf{z}_1), \mathbf{x}_2 = \mathbf{f}_2(\mathbf{z}_2)$ on remplace juste l'identité par les jacobiniennes de ses fonctions. Le problème vient de la différentiation par rapport à y . On repart d'un cas particulier

$$\text{ifelse}(g(z_1), z_1, z_2) := \begin{cases} f_1(z_1) & \text{si } g(z_1) \geq 0 \\ f_2(z_2) & \text{sinon} \end{cases}$$

$$= H(g(z_1))f_1(z_1) + (1 - H(g(z_1)))f_2(z_2)$$

Maintenant calculons la dérivée de cette fonction par rapport à z_1 :

$$\partial_{x_1} \text{ifelse} = (\partial_{z_1} g)(z_1)H'(g(z_1))f_1(z_1) + H(g(z_1))\partial_{z_1} f_1(z_1) + (1 - (\partial_{z_1} g)(z_1))H$$

Puisque la dérivée de la fonction de Heaviside H est nulle on obtient

$$\partial_{x_1} \text{ifelse} = H(g(z_1))\partial_{z_1} f_1(z_1) + f_2(z_2)$$

et donc on voit que la dérivée de g n'influe nullement la dérivée $\partial_{x_1} \text{ifelse}$ ce qui pose problème. Pour cette raison il faut aussi régulariser le "**ifelse**". En pratique il suffit d'utiliser des opérateurs de comparaison régularisés (utilisant la fonction sigmoïde), des opérateurs logiques régularisés [Définition 17.56](#) permettra d'obtenir des "**ifelse**" différentiable si on utilise

$$yx_1 + (1 - y)x_2$$

avec $y \in [0, 1]$. Par exemple dans l'exemple précédent il suffit de remplacer $H(g(z_1))$ par $\text{sig}_\theta(g(z_1))$.

Notamment dans le cadre des opérateurs logique il peut être assez fréquent d'utiliser les fonctions **Argmax** et **Argmin**. De la même façon on utilise des versions régularisées appelées **SoftArgmax** et **SoftArgmin**

Définition 17.57. Fonction SoftArgmax et SoftArgmin. La fonction Softargmax d'un vecteur $x \in \mathbb{R}^d$ est donnée par

$$\text{softargmax}(x) := \frac{\exp(x)}{\sum_{i=1}^d \exp(x_i)}$$

et le Softargmin est donné par

$$\text{softargmin}(x) = \text{softargmax}(-x)$$

Résumé.

On a principalement étudié des opérateurs qui travaillent sur des booléens $\{0, 1\}$ ce qui les rend non différentiables. Globalement la stratégie consiste à remplacer les booléens par des réels dans $[0, 1]$ avec des fonctions qui ne produisent pas (comme le \geq) et ne produisent pas des booléens mais

des réels dans $[0, 1]$ qui se concentrent vers 0 ou vers 1.

Maintenant qu'on a traité les opérateurs de comparaison et logique et les structures conditionnelles on va introduire le traitement des boucles. Pour les boucles **for** la situation est assez simple. En effet une boucle "for" peut se voir comme une composition de fonctions. Si chaque fonction est différentiable la boucle sera différentiable. Différencier une boucle "for" est donc similaire à différencier un programme. Ce sera introduit dans la sous section suivante. Pour les boucles **while** la situation est plus compliquée car le nombre d'itération dépend d'une structure conditionnelle. On expliquera pas en détails les difficultés mais pour différencier à travers une boucle "while" on va régulariser la structure conditionnelle de test et dans le test on va mettre une condition d'arrêt supplémentaire avec un nombre maximal d'opérateur. On peut le faire avec opérateur **égal** régularisé.

17.4.3.3 Différentiation automatique

Maintenant qu'on a vu comment on modélisait un programme avec un graphe puis comment on différenciait les structures de bases comment les opérateurs binaires, les structures conditionnelles et les boucles. Maintenant on va introduire les outils fondamentaux de différentiation automatique qui vont permettre de différencier à travers un programme donc un ensemble de fonction.

On part d'une séquence simple de fonctions:

$$\begin{aligned} x_0 &\in \mathcal{E}_0 \\ x_1 &= f_1(x_0) \in \mathcal{E}_1 \\ &\vdots \\ x_K &= f_K(x_{K-1}) \in \mathcal{E}_K \\ f(x_0) &= x_K \end{aligned}$$

En utilisant la règle de la dérivée en chaîne [\(17.13\)](#) on obtient:

$$\partial f(x_0) = \partial f_K(x_{K-1}) \partial f_{K-1}(x_{K-2}) \dots \partial f_2(x_1) \partial f_1(x_0)$$

On voit donc que la jacobienne de la transformation totale est le produit des jacobienes de chacune des transformations. Ce calcul peut être très coûteux car on va effectuer K produits entre des matrices. En pratique on utilise rarement la jacobienne. On veut souvent calculer la jacobienne dans une direction donc multiplier par un vecteur d'entrée v . La formule va nous donner une approche où on calcule, au fur et à mesure, nos fonctions f_i on va aussi calculer leurs jacobienes et appliquer l'application linéaire associée au résultat

de l'application linéaire de la fonction précédente.

Algorithm 17.58. Différentiation automatique en mode avant (forward).

- **Données** : les fonctions f_1, \dots, f_K .
- **Entrées** : état $\mathbf{x}_0 \in \mathcal{E}_0$ et direction $\mathbf{v} \in \mathcal{E}_0$
- On initialise $\mathbf{t}_0 = \mathbf{v}$
- $\forall k = 1 < K$
 - On calcul $\mathbf{x}_k = \mathbf{f}_k(\mathbf{x}_{k-1}) \in \mathcal{E}_k$
 - On calcul $\mathbf{t}_k = \partial \mathbf{f}_k(\mathbf{x}_{k-1})[\mathbf{t}_{k-1}] \in \mathcal{E}_k$
- On renvoie $\mathbf{f}(\mathbf{x}_0) = \mathbf{x}_K$ et $\partial \mathbf{f}(\mathbf{x}_0)[\mathbf{v}] = \mathbf{t}_K$.

On retrouve sur la figure [Figure 17.59](#) un schéma expliquant l'algorithme.

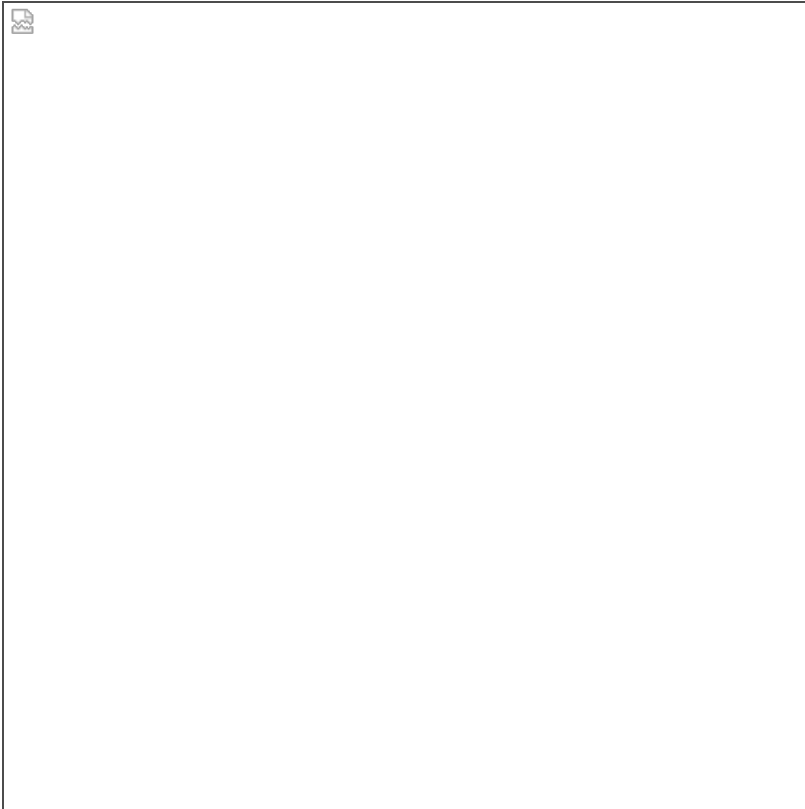


Figure 17.59. principe de l'autodifférentiation avant. Image issue du livre [\[1.45\]](#).

On souhaite maintenant essayer d'estimer l'empreinte mémoire de cet algorithme c'est à dire la quantité d'informations maximal qui sera stocker lors

de l'auto-différentiation avant. Si on regarde juste l'évaluation de l'empreinte mémoire de la fonction \mathbf{f} donnée par la composition de nos K fonctions elle est juste celle du plus grand \mathbf{x}_k car chaque fonction utilise uniquement le résultat de la précédente. Maintenant on va regarder ce qui se passe pour notre algorithme. En pratique on voit qu'on stocke uniquement deux objets \mathbf{x}_k et \mathbf{t}_k donc la l'empreinte maximal en mémoire est celle des plus gros objets de la séquence.

- Maintenant on va regarder l'exemple ou on veut calculer le gradient d'une fonction de cout appliquée à la sortie de notre fonction composée.

$$\nabla_{\mathbf{x}_0} L(\mathbf{f}(\mathbf{x}_0))$$

- avec $L : \mathcal{E}_K \rightarrow \mathbb{R}$. Ce cadre apparait quand lorsqu'on veut calculer le gradient d'une fonction de coût appliquée au résultat de \mathbf{f} . La formule de la dérivée en chaine nous donne:

$$\nabla_{\mathbf{x}_0} L(\mathbf{f}(\mathbf{x}_0)) = (\partial \mathbf{f}^*(\mathbf{x}_0)) \nabla_{\mathbf{x}} L(\mathbf{f}(\mathbf{x}_0))$$

- En pratique on voit donc qu'on cherche à calculer le plus souvent un terme de type VJP:

$$\partial \mathbf{f}^*(\mathbf{x})[\mathbf{u}] = \partial \mathbf{f}^*(\mathbf{x})\mathbf{u}$$

- avec $\mathbf{v} \in \mathcal{E}$ un vecteur. En utilisant la formule de la dérivée en chaine on va donc chercher à calculer

$$\partial \mathbf{f}^*(\mathbf{x}_0) = \partial \mathbf{f}_1^*(\mathbf{x}_0) \partial \mathbf{f}_2^*(\mathbf{x}_1) \dots \partial \mathbf{f}_{K-1}^*(\mathbf{x}_{K-2}) \partial \mathbf{f}_K^*(\mathbf{x}_{K-1})$$

- Ce calcul va donner lieu a un deuxième algorithme appelé l'auto-différentiation rétrograde.

Algorithm 17.60. Différentiation automatique en mode rétrograde (backward).

- **Données** : les fonctions f_1, \dots, f_K .
- **Entrées** : état $\mathbf{x}_0 \in \mathcal{E}_0$ et direction $\mathbf{u} \in \mathcal{E}_K$
- $\forall k = 1 \text{ à } K$
 - On calcul $\mathbf{x}_k = \mathbf{f}_k(\mathbf{x}_{k-1}) \in \mathcal{E}_k$
- On initialise $\mathbf{r}_K = \mathbf{u}$
- $\forall k = K \text{ vers } 1$
 - On calcul $\mathbf{r}_{k-1} = \partial \mathbf{f}_k^*(\mathbf{x}_{k-1})[\mathbf{r}_k] \in \mathcal{E}_{k-1}$

- On renvoie $\mathbf{f}(\mathbf{x}_0) = \mathbf{x}_K$ et $\partial \mathbf{f}(\mathbf{x}_0)[\mathbf{u}] = \mathbf{r}_0$.

La quantité \mathbf{r}_k qui décrit la variation de notre fonction \mathbf{f}_k par rapport à une direction de sortie est appelé ***l'état adjoint*** de \mathbf{x}_k . On retrouve sur la figure [Figure 17.61](#) un schéma expliquant l'algorithme.

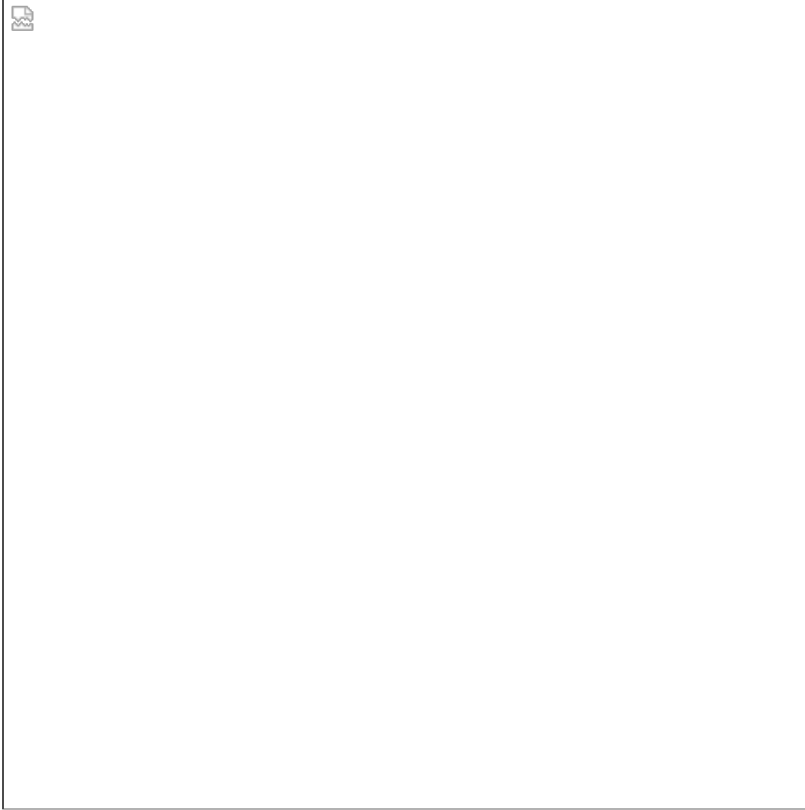


Figure 17.61. principe de l'auto-différentiation rétrograde. Image issue du livre [\[1.45\]](#).

L'algorithme rétrograde est celui qui est populairement utilisé pour les réseaux de neurones et pour l'apprentissage en général car il permet de calculer naturellement le gradient d'une fonction coût [\(17.15\)](#). Cependant la consommation mémoire est différente et plus importante que dans le cas en avant. On voit la différence sur la figure [Figure 17.62](#)

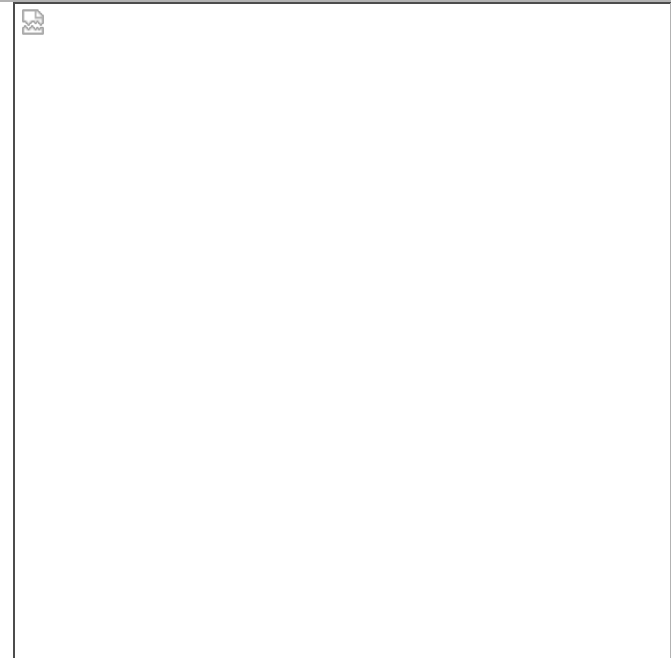
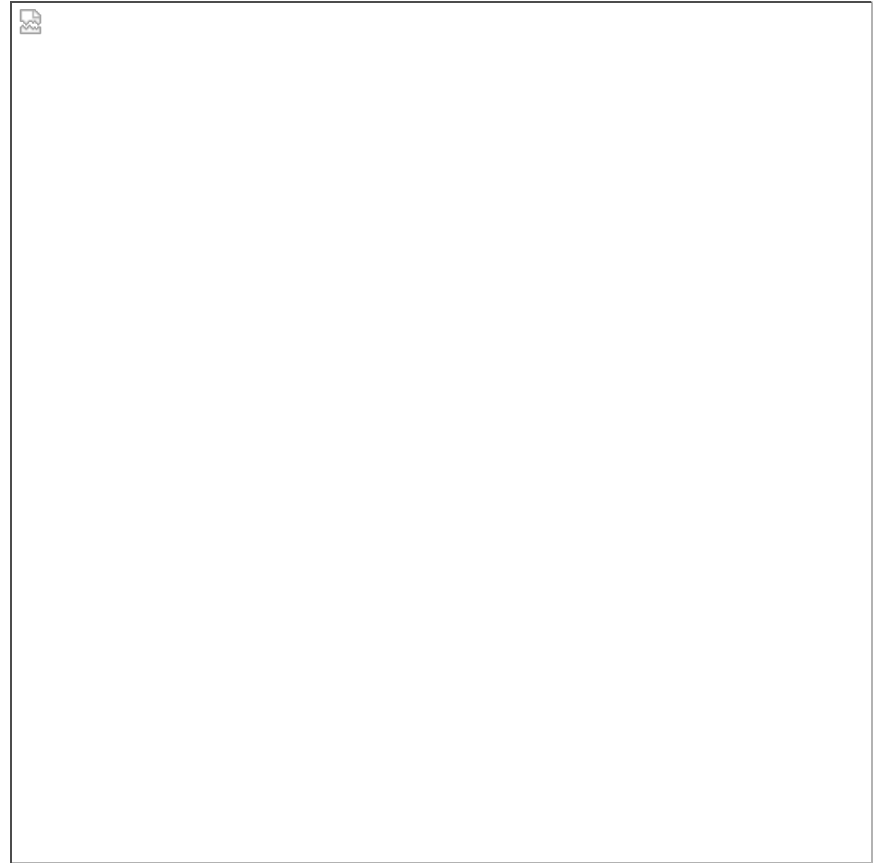


Figure 17.62. Comparaison de l'empreinte mémoire entre les deux approches

d'auto-différentiation. Gauche: en avant. Droite: rétrograde. Image issue du livre [1.45].

Maintenant on va comparer l'empreinte mémoire des deux approches pour le calcul de la jacobienne complète et pas seulement d'un produit jacobienne vecteur. On suppose qu'on a $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^m$. On suppose que les dimension intermédiaire sont tous de taille d . Dans le cas avant pour calculer la jacobienne on va calculer le produit $\partial \mathbf{f} \mathbf{v}$ en prenant comme valeur de \mathbf{v} tout les vecteurs de base \mathbf{e}_i de \mathbb{R}^d . On va donc devoir faire d calculs de type JVP. Chaque appel de type JVP a comme empreinte mémoire $O(\max(d, m))$ et le cout de calcul est donné par le cout d'un produit matrice vecteur donc $\sum_{i=1}^{d-1} d^2 + dm$ ce qui donne donc un cout total par la jacobienne entière de de l'ordre $O(md^2 + Kd^3)$. Dans le cas arrière il faudra aussi multiplier par tous les vecteurs de base de \mathbb{R}^m . On devrait faire donc m calculs de type VJP. Pour la complexité de calcul elle est la même par VJP que par JVP donc on obtient $O(m^2d + Kmd^2)$. Par contre le cout mémoire est de l'ordre $O(Kd + m)$.

► Exemple 17.63.

Si on prend un réseau avec 10 entrées et 10 neurones par couche avec 6 couches. Le mode avant nous donne une complexité mémoire en $O(10)$ et de calcul en $O(6100)$ la ou le mode rétrograde nous donnera une complexité mémoire en $O(60)$ et de calcul en $O(610)$.

Utilisation pratique.

En général on conseille le mode rétrograde si $m < d$ (fréquent en apprentissage) et le mode en avant si $m \geq d$. En effet pour le cas $d = m = 1$ les deux ont un coût CPU en $O(K)$ mais la mémoire de la méthode en avant est en $O(1)$.

Le cas d'un programme séquentiel introduit ici est suffisant pour décrire précisément l'algorithme de retro-propagation du gradient. Il suffit s'utiliser la définition de VJP pour le cas à plusieurs entrées. On va juste écrire l'algorithme de la rétropropagation du gradient d'un modèle paramétrique \mathbf{f}_θ à plusieurs couches. Il s'agit donc d'une application du mode rétrograde.

Algorithm 17.64. Rétropropagation du gradient.

- **Données** : les fonctions f_1, \dots, f_K et une fonction de L
- **Entrées** : Données entrée-sortie (\mathbf{x}, \mathbf{y}) et vecteur de paramètre

$$\theta = (\theta_1, \dots, \theta_K)$$

- on initialise $\mathbf{x}_0 = \mathbf{x}$
- $\forall k = 1 \text{ à } K$
 - On calcul $\mathbf{x}_k = \mathbf{f}_k(\mathbf{x}_{k-1}, \theta_k) \in \mathcal{E}_k$
- On calcule $L(\mathbf{x}_K, \mathbf{y})$
- On initialise $\mathbf{r}_K = \nabla_1 L(\mathbf{x}_K, \mathbf{y})$
- $\forall k = K \text{ vers } 1$
 - On calcul $(\mathbf{r}_{k-1}, \mathbf{g}_k) = \partial \mathbf{f}_k^*(\mathbf{x}_{k-1}, \theta_k)[\mathbf{r}_k] \in \mathcal{E}_{k-1} \times \Theta_k$
- On renvoie la fonction de cout total et son gradient $L(\mathbf{x}_K, \mathbf{y})$ et $\partial L(\mathbf{x}_K, \mathbf{y}) = (\mathbf{g}_1, \dots, \mathbf{g}_K)$.

Dan cet algorithme le terme \mathbf{g}_k correspond à l'adjoint des poids donc à la variation de la fonction \mathbf{f}_k par rapport aux variations de θ_k .

Jusqu'à présent on a traité le cas d'une séquence simple. En pratique on a vu que les programmes étaient plutôt définis par des graphes que par des séquences linéaires. On va donc juste introduire les modes avant et rétrograde dans le cas plus général d'un programme défini par un graphe.

Algorithm 17.65. Différentiation automatique en mode avant pour un graphe.

- **Données** : les fonctions f_1, \dots, f_K dans l'ordre topologique du graphe
- **Entrées** : état $\mathbf{x}_0 \in \mathcal{E}_0$ et direction $\mathbf{v} \in \mathcal{E}_0$
- On initialise $\mathbf{t}_0 = \mathbf{v}$
- $\forall k = 1 < K$:
 - On récupère les noeuds parents $i_1, \dots, i_{p_k} = \text{pa}(k)$
 - On calcul $\mathbf{x}_k = \mathbf{f}_k(\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_{p_k}})$
 - On calcul $\mathbf{t}_k = \sum_{j=1}^{p_k} \partial_j \mathbf{f}_k(\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_{p_k}})[\mathbf{t}_{i_j}]$
- On renvoie $\mathbf{f}(\mathbf{x}_0) = \mathbf{x}_K$ et $\partial \mathbf{f}(\mathbf{x}_0)[\mathbf{v}] = \mathbf{t}_K$.

Algorithm 17.66. Différentiation automatique en mode rétrograde pour un graphe.

- **Données** : les fonctions f_1, \dots, f_K donnée dans l'ordre topologique de graphe.
- **Entrées** : état $\mathbf{x}_0 \in \mathcal{E}_0$ et direction $\mathbf{u} \in \mathcal{E}_K$

- $\forall k = 1 \text{ à } K$
 - On récupère les noeuds parents $i_1, \dots, i_{p_k} = pa(k)$
 - On calcul $\mathbf{x}_k = \mathbf{f}_k(\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_{p_k}})$
 - On instancie $\mathbf{l}_k = \partial \mathbf{f}_k^*(\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_{p_k}})$
- On initialise $\mathbf{r}_K = \mathbf{u}$ et $\mathbf{r}_k = 0, \quad \forall k \in \{0, K-1\}$
- $\forall k = K \text{ vers } 1$
 - On récupère les noeuds parents $i_1, \dots, i_{p_k} = pa(k)$
 - On calcule

$$(\delta_{i_1}, \dots, \delta_{i_{p_k}}) = \mathbf{l}_k[\mathbf{r}_k]$$

avec $(\delta_{i_1}, \dots, \delta_{i_{p_k}})$ les sorties de VJP avec plusieurs entrées définies dans [Définition 17.44](#).

- On calcule $\mathbf{r}_{i_j} \leftarrow \mathbf{r}_{i_j} + \delta_{i_j}, \quad \forall j \in \{1, \dots, p_k\}$
- On renvoie $\mathbf{f}(\mathbf{x}_0) = \mathbf{x}_K$ et $\partial \mathbf{f}(\mathbf{x}_0)[\mathbf{u}] = \mathbf{r}_0$.

Si une valeur intermédiaire \mathbf{s}_k est utilisée par les fonctions $f_{j_1}, \dots, f_{j_{c_k}}$ ou $j_1, \dots, j_{c_k} = \text{enfant}(k)$, les dérivées par rapport à \mathbf{x}_k va donc utiliser la somme des variations de toutes les fonctions. Cela nous donne

$$\mathbf{r}_k := \sum_{j \in \{j_1, \dots, j_{c_k}\}} \delta_{k,j}$$

avec $\delta_{k,j}$ Dans l'algorithme lorsqu'on a la variation d'un noeud on va la distribuer dans toutes les variations \mathbf{r} des enfants de ce noeud.

Une fois que vous avez une mode différentiation automatique il suffit de le coupler avec des fonctions de bases différentiable (on peut parler de fonctions primitives) pour obtenir un code différentiable. Les fonctions primitives sont les fonctions mathématiques usuelles ou les fonctions régularisées introduites précédemment comme le "et", le "=" etc. Il est important que pour ces fonctions les modes JVP ou VJP soient définis. On a besoin aussi que les compositions de ces fonctions primitives restent des fonctions dérivables.

Le mode rétrograde qui est principalement utilisé en apprentissage et en optimisation est assez consommateur en mémoire. En effet la complexité mémoire est de l'ordre de K le nombre de fonction (noeuds du graphe de calcul). Pour limiter il existe plusieurs approches. La première approche est simple mais limitée en pratique elle consiste à construire des fonctions

inversible. Cela permet de calculer l'état dont on a besoin pour le mode-rétrograde juste avec l'évaluation de l'inverse.

Algorithm 17.67. Différentiation automatique en mode rétrograde avec fonctions inversibles.

- **Données** : les fonctions f_1, \dots, f_K .
- **Entrées** : état $\mathbf{x}_0 \in \mathcal{E}_0$ et direction $\mathbf{u} \in \mathcal{E}_K$
- On initialise $\mathbf{r}_K = \mathbf{u}$
- On calcul $\mathbf{x}_K = \mathbf{f}_K \circ \dots \circ \mathbf{f}_1(\mathbf{x}_0)$
- $\forall k = K \text{ vers } 1$
 - On calcul $\mathbf{x}_{k-1} = \mathbf{f}_{k-1}^{-1}(\mathbf{x}_k)$
 - On calcul $\mathbf{r}_{k-1} = \partial \mathbf{f}_k^*(\mathbf{x}_{k-1})[\mathbf{r}_k]$
- On renvoie $\partial \mathbf{f}(\mathbf{x}_0)[\mathbf{u}] = \mathbf{r}_0$.

Bien que des architectures de réseaux peuvent être inversible cette stratégie reste limitée en pratique. L'approche la plus utilisée est celle des points de contrôle.

Définition 17.68. Méthode de point de contrôle pour le mode rétrograde. La méthode des points de contrôle fonctionne en ne stockant sélectivement qu'un sous-ensemble de nœuds intermédiaires, appelés points de contrôle, et en recalculant les autres.

On va commencer par introduire le cas extrême où l'on ne stocke aucune valeur dans le cadre d'un programme séquentiel.

Algorithm 17.69. Différentiation automatique en mode rétrograde sans stockage mémoire.

- **Données** : les fonctions f_1, \dots, f_K .
- **Entrées** : état $\mathbf{x}_0 \in \mathcal{E}_0$ et direction $\mathbf{u} \in \mathcal{E}_K$
- On initialise $\mathbf{r}_K = \mathbf{u}$
- $\forall k = K \text{ vers } 1$
 - On calcul $\mathbf{x}_{k-1} = \mathbf{f}_{k-1} \circ \dots \circ \mathbf{f}_1(\mathbf{x}_0)$
 - On calcul $\mathbf{r}_{k-1} = \partial \mathbf{f}_k^*(\mathbf{x}_{k-1})[\mathbf{r}_k]$
- On renvoie $\partial \mathbf{f}(\mathbf{x}_0)[\mathbf{u}] = \mathbf{r}_0$.

A ajouter: méthode de checkpointing

17.4.4 Différentier à travers des fonctions implicites

Jusqu'à présent, nous avons appris à différentier à travers des fonctions explicites. En pratique, il est possible de se retrouver avec des fonctions implicites à différentier. Par exemple, si on prend un schéma d'Euler implicite en temps, la détermination du temps $n + 1$ va nécessiter la résolution d'un système linéaire ou nonlinéaire dépendant de \mathbf{f}_θ . On va donc potentiellement vouloir différentier par rapport à θ la fonction $S : \theta \rightarrow \mathbf{x}^{n+1}$ ou $S : \mathbf{x}^{n+1}$ est la solution du schéma d'Euler implicite. Cette fonction S n'est pas explicite par rapport à θ . Pour cette raison, on va voir comment différentier des fonctions implicites.

Il peut exister plusieurs exemples:

- Les problèmes d'optimisations où on a:

$$h(\theta) = g(\mathbf{x}^*(\theta), \theta), \quad \text{avec } \mathbf{x}^*(\theta) = \operatorname{argmax}_{\mathbf{x} \in \mathcal{E}} f(\mathbf{x}, \theta) \quad (17.18)$$

donc on souhaite différentier une fonction qui dépend du résultat d'un processus d'optimisation qui dépend de paramètre. Il y a donc une dépendance implicite. Prenons le cas d'une régression linéaire Ridge de la forme

$$\theta^* = \operatorname{argmax} (\|Y - \theta X\|^2 + \lambda \|\theta\|^2)$$

Les paramètres optimaux θ^* sont une fonction implicite du paramètre de régularisation λ on peut donc l'écrire $\theta^*(\lambda)$. Si ensuite on a un critère pour juger du meilleur paramètre λ on peut vouloir le dériver par rapport à λ pour ensuite le déterminer par une méthode de gradient. On se retrouve typiquement dans notre cas.

- La recherche de zéro d'équations nonlinéaires de la forme:

$$f(\mathbf{x}, \theta) = 0$$

qui englobent notamment les schémas implicites dont on a parlé précédemment.

La solution d'un problème d'optimisation ou la recherche de zéro d'une fonction est effectué en général en utilisant un algorithme itératif. Cet algorithme peut s'interpréter comme une composition de fonctions simples qui sont ou qu'on peut rendre dérivable. On peut donc a priori dériver ses algorithmes avec de la différentiation automatique (typiquement rétrograde). On parle de **Dérouler l'algorithme**. On peut voir cette approche comme une variante de **discrétiser**

puis optimiser. En pratique cela peut être très coûteux car on va devoir sauvegarder les états intermédiaires comme montré précédemment. Cela peut devenir très coûteux en mémoire. Une autre approche est de voir si on peut calculer directement le VJP. On pourrait voir cela comme l'approche **optimiser puis discrétiser**. On va introduire ses méthodes pour les problèmes implicites.

17.4.4.1 Différentier à travers des problèmes d'optimisation

- Repartons du problème (17.18) avec $\mathbf{f}, \mathbf{g} : \mathcal{E} \times \Theta \rightarrow \mathbb{R}$. Si on dérive notre fonction h on obtient

$$\partial_\theta h(\theta) = \partial_1 g(\mathbf{x}^*(\theta), \theta) \partial_\theta \mathbf{x}^*(\theta) + \partial_2 g(\mathbf{x}^*(\theta), \theta)$$

- Le seul point bloquant de ce calcul est de déterminer $\partial_\theta \mathbf{x}^*(\theta)$. On va maintenant introduire un cas simplifié où le problème peut être résolu plus facilement.

- On va se restreindre ici au cas $g(\mathbf{x}, \theta) = f(\mathbf{x}, \theta)$. Dans ce cas le problème devient

$$h(\theta) = f(\operatorname{argmax}_{\mathbf{x} \in \mathcal{E}} f(\mathbf{x}, \theta), \theta).$$

- Hors évaluer une fonction en son argmax est équivalent au maximum donc notre problème ici devient un problème de recherche de maximum:

$$h(\theta) = \max_{\mathbf{x} \in \mathcal{E}} f(\mathbf{x}, \theta)$$

- ce qui a priori sera plus simple à différentier. On cherche donc le gradient d'un problème d'optimisation paramétrique par rapport à ses paramètres.

Définition 17.70. Théorème de Danskin. Soit $\mathbf{f} : \mathcal{E} \times \Theta \rightarrow \mathbb{R}$ avec Θ un ensemble convexe. Posons

$$h(\theta) = \max_{\mathbf{x} \in \mathcal{E}} f(\mathbf{x}, \theta)$$

et

$$\mathbf{x}^*(\theta) := \operatorname{argmax}_{\mathbf{x} \in \mathcal{E}} f(\mathbf{x}, \theta)$$

- Si \mathbf{f} est concave en \mathbf{x} , convexe en θ , et atteint son unique maximum en $\mathbf{x}^*(\theta)$ alors la fonction h est différentiable de gradient:

$$\nabla_\theta h(\theta) = \nabla_2 f(\mathbf{x}^*(\theta), \theta)$$

Le Théorème de Danskin dit finalement qu'on peut dériver simplement f comme si \mathbf{x}^* ne dépendait pas de θ ce qui rend le calcul très facile. Le théorème de Rockafellar nous donne le même résultat avec des hypothèses la fonction f différente.

Maintenant on va revenir a un cadre plus géénral d'optimisation sous contrainte et rappeler un résultat qui nous permettra de construire la différentiation de ce problème d'optimisation par rapport aux paramètres

Proposition 17.71. On considère le problème suivant avec $f, g, h : \mathcal{E} \times \Theta \rightarrow \mathbb{R}$:

$$\mathbf{x}^*(\theta) = \underset{\mathbf{x}}{\operatorname{argmin}} f(\mathbf{x}, \theta)$$

sous les conditions que $g(\mathbf{x}, \theta) \leq 0$
 $h(\mathbf{x}, \theta) = 0$

La solution de ce problème est donnée celle de $G(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}; \theta) = 0$ avec

$$G(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}; \theta) = \begin{bmatrix} \nabla_{\mathbf{x}} f(\mathbf{x}; \theta) + \partial_{\mathbf{x}} g(\mathbf{x}; \theta)^T \boldsymbol{\lambda} + \partial_{\mathbf{x}} h(\mathbf{x}; \theta)^T \boldsymbol{\nu} \\ \boldsymbol{\lambda} \circ g(\mathbf{x}; \theta) \\ h(\mathbf{x}; \theta) \end{bmatrix}$$

► Preuve.

On commence par appliquer les conditions de Karush-Kuhn-Tucker. \A compléter

Ce résultat nous montre donc que dériver la solution d'un problème d'optimisation va tout naturellement se transformer en un problème de dérivation d'une fonction implicite liée au gradient (cas sans contrainte) ou au Lagrangien (cas avec contrainte). On va maintenant traiter le cas des fonctions implicites qui était un des enjeux de cette section.

17.4.4.2 Différentier à travers des fonctions implicites

On considère une équation nonlinéaire $F(\mathbf{x}; \theta) = \mathbf{0}$ avec $F : \mathcal{E} \times \Theta \rightarrow \mathcal{E}$. Si on considère $\mathbf{x}^*(\theta)$ la solution de l'équation. On cherche donc à calculer $\partial_{\theta} \mathbf{x}^*(\theta)$. Cela nous permettra de traiter un certains nombres de problème comme les produits d'optimisation introduit précédemment. On va introduire formellement la solution avant de la justifier.

On part de $F(\mathbf{x}^*(\theta); \theta) = \mathbf{0}$ et on dérive par rapport à θ ce qui donne:

$$\partial_1 F(\mathbf{x}^*(\theta), \theta) \partial_{\theta} \mathbf{x}^*(\theta) + \partial_2 F(\mathbf{x}^*(\theta), \theta) = \mathbf{0}$$

Donc la dérivée que l'on cherche doit satisfaire:

$$-\partial_1 F(\mathbf{x}^*(\theta), \theta) \partial_{\theta} \mathbf{x}^*(\theta) = \partial_2 F(\mathbf{x}^*(\theta), \theta)$$

Pour calculer la dérivée il nous faut donc inverser $\partial_1 F(\mathbf{x}^*(\theta), \theta)$. Ceci est assuré sous certaines conditions par le théorème des fonctions implicites que l'on rappelle.

Définition 17.72. Théorème des fonctions implicites. Soit $F : \mathcal{E} \times \Theta \rightarrow \mathcal{E}$. On suppose $F(\mathbf{x}, \theta)$ une fonction continument différentiable dans le voisinage de (\mathbf{x}_0, θ_0) tel que $F(\mathbf{x}_0, \theta_0) = \mathbf{0}$ et $\partial_1 F(\mathbf{x}_0, \theta_0)$ soit inversible. Alors il existe un voisinage U_0 de θ_0 dans lequel il existe une fonction $\mathbf{x}^*(\theta)$ tel que:

- $\mathbf{x}^*(\theta_0) = \mathbf{x}_0$,
- $F(\mathbf{x}^*(\theta), \theta) = \mathbf{0}, \forall \theta \in U_0$,
- On a:

$$\partial_{\theta} \mathbf{x}^*(\theta) = -\partial_1 F(\mathbf{x}^*(\theta), \theta)^{-1} \partial_2 F(\mathbf{x}^*(\theta), \theta).$$

Cela nous permet d'obtenir facilement le Jvp et VJP pour les fonctions implicites.

Proposition 17.73. JVP et VJP pour les fonctions implicites. Soit $\mathbf{x}^*(\theta)$ la solution d'une équation implicite de type $F(\mathbf{x}, \theta) = \mathbf{0}$ avec $F : \mathcal{E} \times \Theta \rightarrow \mathcal{E}$. On pose

$$A = -\partial_1 F(\mathbf{x}^*(\theta), \theta), \quad B = \partial_2 F(\mathbf{x}^*(\theta), \theta)$$

Le JVP $\mathbf{t} = \partial \mathbf{x}^*(\theta) \mathbf{v}$ est solution de

$$A \mathbf{t} = B \mathbf{v}$$

et le VJP $\partial \mathbf{x}^*(\theta)^* \mathbf{u}$ est donné par

$$\partial \mathbf{x}^*(\theta)^* \mathbf{u} = B^* \mathbf{r}$$

avec \mathbf{r} solution de:

$$A^* \mathbf{r} = \mathbf{u}$$

► Preuve.

Le JVP s'obtient immédiatement en multipliant le résultat du théorème des fonctions implicites par \mathbf{v} . \\\ TO DO: Ajouter VJP

L'approche dit "déroulée" peut se voir comme l'approche DpO appliquée aux fonctions implicites la ou l'utilisation du mode VJP introduite juste au dessus peut être vu comme l'approche OpD. Dans ce cadre l'approche OpD est plus efficace car elle évite un coût mémoire important lié au "dérouler" des algorithmes de résolution comme les méthodes de point fixes ou de Newton.

17.4.4.3 Application en algèbre linéaire

Dans cette section on va proposer des exemples en algèbre linéaire. On va supposer qu'on va avoir ici des matrices ou des systèmes linéaires qui dependent de paramètres θ et une fonction de coût qui va dépendre de calcul utilisant ces matrices et membres de droites. On souhaite donc dériver notre fonction de cout par rapport à θ en utilisant la différentiation de fonctions implicites.

Dans un premier exemple, on va supposer qu'on souhaite résoudre

$$\min_{\theta} J(\mathbf{x}^s), \quad \text{sous contrainte } A(\theta)\mathbf{x}^* = \mathbf{b}(\theta)$$

On souhaite donc calculer le VJP

$$\nabla_{\theta} J(\mathbf{x}^*) = (\partial_{\theta}^* \mathbf{x}^s) \nabla_{\mathbf{x}} J$$

ce qui est égal à la chaine de VJP

$$\nabla_{\theta} J(\mathbf{x}^*) = (\partial_{\theta}^*(A, \mathbf{b}))(\partial_{A, \mathbf{b}}^* \mathbf{x}^s) \nabla_{\mathbf{x}} J$$

Le VJP associé à $(\partial_{\theta}^*(A, \mathbf{b}))$ sera naturellement calculer par l'autodifférentiation appliqué aux fonctions qui construisent A, \mathbf{b} en fonction de θ . Par contre pour le calcul VJP associé à la solution \mathbf{x}^s la situation est différente. Si on code un algorithme d'inversion comme une méthode de gradient Conjugué ou une décomposition LU vous allez devoir différentier à travers des centaines voir milliers d'appels de fonctions. Il est donc plus efficace de dériver le VJP par la méthode des fonctions implicites directement.

Proposition 17.74. VJP d'une résolution de système linéaire. Soit une système linéaire $A\mathbf{x} = \mathbf{b}$. On nomme la solution \mathbf{x}^s du solution du système. le vecteur de sortie est nommé \mathbf{u} . Le VJP associé à la solution est donné par

$$\frac{\partial^* \mathbf{x}^s}{\partial A} \mathbf{u} = -\mathbf{r} \otimes \mathbf{x}^s$$

et

$$\frac{\partial^* \mathbf{x}^s}{\partial \mathbf{b}} \mathbf{u} = \mathbf{r}$$

avec $A^t \mathbf{r} = \mathbf{u}$ le système adjoint.

► Preuve.

On commence par écrire notre problème sous la forme d'une équation implicite:

$$F(\mathbf{x}, A, \mathbf{b}) = A\mathbf{x} - \mathbf{b} = 0$$

Il nous fait maintenant calculer les différentielle. Pour la première c'est facile. Il est en effet immédiat que

$$\partial_1 F(\mathbf{x}, A, \mathbf{b}) = A$$

Maintenant on va chercher les différentielles par rapport à \mathbf{b} et A .

- Différentielle par rapport à \mathbf{b} :

$$F(\mathbf{x}, A, \mathbf{b} + \delta \mathbf{b}) = A\mathbf{x} - \mathbf{b} - \delta \mathbf{b}$$

donc

$$(\partial_3 F) \delta \mathbf{b} = F(\mathbf{x}, A, \mathbf{b} + \delta \mathbf{b}) - F(\mathbf{x}, A, \mathbf{b})$$

ce qui donne

$$(\partial_3 F) \delta \mathbf{b} = -\delta \mathbf{b}$$

donc $(\partial_3 F) = -I_d$ et donc $(\partial_3 F)^* = -I_d$

- Différentielle par rapport à A . On écrit dans un premier temps:

$$F(\mathbf{x}, A + \delta A, \mathbf{b}) = (A + \delta A)\mathbf{x} - \mathbf{b}$$

On cherche donc à déterminer l'application linéaire:

$$(\partial_2 F) \delta A : \mathcal{M}_{n,n}(\mathbb{R}) \rightarrow \mathbb{R}^n$$

On a donc

$$(\partial_2 F) \delta A = F(\mathbf{x}, A + \delta A, \mathbf{b}) - F(\mathbf{x}, A, \mathbf{b}) = \delta A \mathbf{x}$$

On va maintenant calculer le dual dont on a besoin pour appliquer les résultats sur les fonctions implicites. On cherche donc à déterminer:

$$(\partial_2 F)^* \delta A : \mathbb{R}^n \rightarrow \mathcal{M}_{n,n}(\mathbb{R}).$$

On va donc calculer

$$\langle (\partial_2 F) \delta A, \mathbf{y} \rangle_{\mathbb{R}^n} = \langle \delta A \mathbf{x}, \mathbf{y} \rangle_{\mathbb{R}^n}$$

On veut maintenant se ramener à un produit scalaire dans $\mathcal{M}_{n,n}(\mathbb{R})$ pour obtenir l'application linéaire duale [Définition 17.41](#):

$$\langle \delta A \mathbf{x}, \mathbf{y} \rangle_{\mathbb{R}^n} = \text{Tr}(\delta A \mathbf{x} \mathbf{y}^t) = \langle \delta A, \mathbf{y} \otimes \mathbf{x} \rangle_{\mathcal{M}_{n,n}(\mathbb{R})}$$

on a donc

$$(\partial_2 F)^* : \mathbf{y} \rightarrow \mathbf{y} \otimes \mathbf{x}$$

Maintenant qu'on a obtenu nos différentielles on peut conclure avec la proposition [Proposition 17.73](#). L'état adjoint \mathbf{r} est solution de

$$-\partial_1 F(\mathbf{x}, A, \mathbf{b})^* \mathbf{r} = -A^t \mathbf{r} = \mathbf{u}$$

Ensuite on va obtenir le résultat en appliquant la suite de la formule d'une VJP [Proposition 17.73](#):

$$\frac{\partial^* \mathbf{x}^s}{\partial \mathbf{b}} \mathbf{u} = (\partial_3 F)^* \mathbf{r} = -\mathbf{r}$$

et

$$\frac{\partial^* \mathbf{x}^s}{\partial A} \mathbf{u} = (\partial_2 F)^* \mathbf{r} = \mathbf{r} \otimes \mathbf{x}.$$

En posant l'adjoint solution de $= A^t \mathbf{r} = \mathbf{u}$ on obtient le résultat final.

d'une EDO. Utiliser la programmation différentiable pour calculer ce même gradient peut être vu comme l'approche DpO. On va maintenant discuter l'intérêt des deux.

Un premier avantage de l'approche DpO est une fois l'ensemble du code programmé de façon différentiable on peut calculer le gradient de l'EDO, modifier le type de fonction coût, différentier les solveurs sans travail supplémentaire. Cependant si on se restreint au cadre des EDO l'approche OpD semble a première vue meilleure puisqu'on a pas a stocker la solution de l'EDO primal grace au dual augmenté de la méthode OdeNet on se retrouve avec un cout mémoire en $O(m)$ au lieu $O(mT)$ pour la méthode DpO. On va cepencent montrer avec des exemples que l'approche OpD ne marche pas toujours et donc qu'une approche DpO avec des "points de contrôles" peut être totalement pertinente.

Le fait qu'on est pas besoin de stocker la solution de l'EDO dans l'approche vient du fait qu'on peut résoudre l'EDO dans le sens inverse car le flot d'une EDO est l'inversible ce qu'on ne peut a priori pas faire avec le mode DpO car les schémas ne pas a priori réversible. Cependant la réversibilité de l'EDO n'est pas automatique. En pratique il y a des conditions de régularité sur f_θ pour avoir de la réversibilité dans un ouvert. Un première exemple est l'EDO:

$$\frac{d}{dt} z(t) = z^3(t)$$

qui admet une solution $\varphi^0(t, z_0) = \frac{z_0}{\sqrt{1-2z_0^2 t}}$ on a donc une solution qui est bien posé jusqu'au temps $t < \frac{1}{2z_0^2}$. On voit donc que l'existence de la solution et donc la réversibilité va dépendre du temps final et de l'ensemble des conditions initiales. Un second exemple est l'EDO

$$\frac{d}{dt} z(t) = -\lambda z(t)$$

dont la solution existe sur \mathbb{R}^+ . Si on résout l'EDO dans le sens inverse on se retrouve a résoudre

$$\frac{d}{dt} z(t) = \lambda z(t)$$

Si λ est grande il s'agit d'une exponentielle fortement croissante. Pour ce problème les schémas seront instable sauf à utiliser des pas de temps infiniment petit. Ces deux exemples montrent qu'en pratique résoudre l'EDO dans le sens inverse peut s'avérer compliquer. Dans [\[1.46\]](#) les auteurs montrent

17.4.5 Régularisation

En construction

17.4.6 Différentiation et ODE

Maintenant que nous avons introduits les outils de différentiation automatique on peut considérer le cas des EDO. Comme introduit au début du chapitre on a vu que l'approche OdeNet revient à la méthode OpD pour le calcul du gradient

aussi la présence d'instabilité si f_θ est un réseau convolutif.

Résumé.

Pour ces raisons les auteurs de [\[1.46\]](#) concluent que l'approche DpO couplé avec des "points de contrôles" est une solution efficace.

Une approche alternative qui pas forcément toujours efficace est d'utiliser des schémas réversibles. En effet ces schémas nous disent que

$$\Phi_{\Delta t}^{-1} = \Phi_{-\Delta t}$$

Cela permet de pouvoir calculer rapidement l'état z_{k-1} à partir de l'état z_k et donc de ne pas stocker les solutions de l'équation primale. Cette approche peut être notamment intéressante pour les EDO conservatives ou les schémas réversibles sont plus fréquents.

