

IFT 3395/6390 Hiver 2008
Fondements de l'apprentissage machine
Professeur: Pascal Vincent

Note: N'attendez pas la date de remise du devoir 2 pour commencer à travailler sur le devoir 3 (qui en est la suite)!

Pour remettre tous vos fichiers (dont un `rapport.pdf`) connectez vous sur la machine remise (`ssh remise`) puis utilisez la commande `remise`.

- Si vous êtes en IFT3395 faites par ex.:
`remise ift3395 tp2 rapport.pdf code.r`
- Si vous êtes en IFT6390 faites par ex.:
`remise ift6390 tp2 rapport.pdf code.r`

Devoir 2

à remettre au plus tard le lundi 31 mars 2008 à minuit.

1 Séparabilité linéaire

On suppose qu'on dispose d'un ensemble de données $D_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ avec $x_i \in \mathbb{R}^d$ et $y_i \in \{-1, +1\}$ indiquant la classe.

- a) Écrivez la forme d'une *fonction discriminante linéaire* g pour ce problème en précisant ses paramètres, ainsi que leur nature et leur dimension. A combien de paramètres scalaires cela correspond-t il?
- b) Que veut-on dire lorsque l'on dit que les points de D_n sont linéairement séparables? Expliquez-le d'abord en Français dans vos propres mots, puis sous forme d'une expression mathématique.
- c) Écrivez l'expression du risque empirique (associé à la fonction discriminante g) qui comptabiliserait le taux d'erreur de classification.
- d) On définit la distance entre un point x et un ensemble de points comme la distance entre x et le point le plus proche de l'ensemble. On cherche la distance entre un point x et la frontière de décision induite par la fonction discriminante g . Exprimez le calcul de cette distance comme un problème d'optimisation (minimisation avec contrainte) et résolvez-le.
- e) On suppose que D_n comporte au moins un point de chaque classe, que tous les x_i sont distincts, et qu'ils sont répartis sans former aucun "alignement" particulier (aucun $n' \leq n$ points ne tombent sur un sous-espace linéaire de dimension $n' - 2$). Montrez que si $n \leq d + 1$, les points de D_n sont toujours linéairement séparables, avec n'importe quelle affectation des étiquettes binaires (les y_i).

2 Apprentissage des paramètres d'un réseau de neurones

On suppose qu'on dispose d'un ensemble de données $D_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ avec $x_i \in \mathbb{R}^d$ et $y_i \in \{1, \dots, m\}$ indiquant la classe parmi m classes.

Soit un réseau de neurones de type *Perceptron multicouche* avec une seule couche cachée (donc 3 couches en tout si on compte la couche d'entrée et la couche de sortie). La couche cachée est constituée de d_h neurones complètement connectés à la couche d'entrée. Ils ont une non-linéarité de type tangente hyperbolique (\tanh). La couche de sortie est constituée de m neurones, complètement connectés à la couche cachée. Ils ont une non-linéarité de type softmax. La sortie du $j^{\text{ème}}$ neurone de la couche de sortie donnera un score pour la classe j interprété comme la probabilité que l'entrée x soit de cette classe j .

Il vous est fortement conseillé de dessiner le réseau de neurones au fur et à mesure afin que vous puissiez mieux suivre les étapes (mais pas besoin de nous fournir un dessin!)

- a) Soit $\mathbf{W}^{(1)}$ la matrice $d_h \times d$ de poids et soit $\mathbf{b}^{(1)}$ le vecteur de biais caractérisant des connexions synaptiques allant de la couche d'entrée à la couche cachée. Indiquez la dimension de $\mathbf{b}^{(1)}$. Donnez la formule de calcul du vecteur d'activations (i.e. avant non-linéarité) des neurones de la couche cachée \mathbf{h}^a à partir d'une observation d'entrée \mathbf{x} , d'abord sous la forme d'une expression de calcul matriciel, puis détaillez le calcul d'un élément h_j^a . Exprimez le vecteur des sorties des neurones de la couche cachée \mathbf{h}^s en fonction de \mathbf{h}^a .
- b) Soit $\mathbf{W}^{(2)}$ la matrice de poids et soit $\mathbf{b}^{(2)}$ le vecteur de biais caractérisant les connexions synaptiques allant de la couche cachée à la couche de sortie. Indiquez les dimensions de $\mathbf{W}^{(2)}$ et $\mathbf{b}^{(2)}$. Donnez la formule de calcul du vecteur d'activations des neurones de la couche de sortie \mathbf{o}^a à partir de leurs entrées \mathbf{h}^s sous la forme d'une expression de calcul matriciel, puis détaillez le calcul de \mathbf{o}_k^a .
- c) La sortie des neurones de sortie est donnée par

$$\begin{aligned} \mathbf{o}^s &= \text{softmax}(\mathbf{o}^a) \\ \mathbf{o}_k^s &= \frac{\exp(\mathbf{o}_k^a)}{\sum_{k'=1}^m \exp(\mathbf{o}_{k'}^a)} \end{aligned}$$

Montrez que les \mathbf{o}_k^s sont positifs et somment à 1. Pourquoi est-ce important?

- d) Le réseau de neurones calcule donc, pour un vecteur d'entrée \mathbf{x} , un vecteur de scores (probabilités) $\mathbf{o}^s(\mathbf{x})$. La probabilité, calculée par le réseau de neurones, qu'une observation \mathbf{x} soit de la classe y est donc donnée par la $y^{\text{ième}}$ sortie $\mathbf{o}_y^s(\mathbf{x})$. Ceci suggère d'utiliser la fonction de perte:

$$L(\mathbf{x}, y) = -\log \mathbf{o}_y^s(\mathbf{x})$$

L'entraînement du réseau de neurones va consister à trouver les paramètres du réseau qui minimisent le risque empirique \hat{R} correspondant à cette fonction de perte. Formulez \hat{R} . Indiquez précisément de quoi est constitué l'ensemble θ des paramètres du réseau. Indiquez à combien de paramètres scalaires n_θ cela correspond. Formulez le problème d'optimisation qui correspond à l'entraînement du réseau permettant de trouver une valeur optimale des paramètres.

- e) Pour trouver la solution à ce problème d'optimisation, on va utiliser une technique de descente de gradient. Exprimez sous forme d'un bref pseudo-code la technique de descente de gradient (batch) pour ce problème.
- f) Notez que le calcul du vecteur de gradient du risque empirique \hat{R} par rapport à l'ensemble des paramètres θ peut s'exprimer comme

$$\begin{pmatrix} \frac{\partial \hat{R}}{\partial \theta_1} \\ \vdots \\ \frac{\partial \hat{R}}{\partial \theta_{n_\theta}} \end{pmatrix} = \frac{1}{n} \sum_{i=1}^n \begin{pmatrix} \frac{\partial L(\mathbf{x}_i, y_i)}{\partial \theta_1} \\ \vdots \\ \frac{\partial L(\mathbf{x}_i, y_i)}{\partial \theta_{n_\theta}} \end{pmatrix}$$

Il suffit donc de savoir calculer le gradient du coût $C = L(\mathbf{x}, y)$ encouru pour un exemple (\mathbf{x}, y) , que l'on définit comme:

$$\frac{\partial C}{\partial \theta} = \begin{pmatrix} \frac{\partial C}{\partial \theta_1} \\ \vdots \\ \frac{\partial C}{\partial \theta_{n_\theta}} \end{pmatrix} = \begin{pmatrix} \frac{\partial L(\mathbf{x}, y)}{\partial \theta_1} \\ \vdots \\ \frac{\partial L(\mathbf{x}, y)}{\partial \theta_{n_\theta}} \end{pmatrix}$$

Pour cela on va appliquer la technique de rétropropagation du gradient, en partant du coût C , et en remontant de proche en proche vers la sortie \mathbf{o} puis vers la couche cachée \mathbf{h} et enfin vers l'entrée \mathbf{x} .

Commencez donc par calculer le vecteur de dérivées partielles du coût C par rapport aux sorties des neurones de la couche de sortie:

$$\frac{\partial C}{\partial \mathbf{o}^s} = \begin{pmatrix} \frac{\partial C}{\partial \mathbf{o}_1^s} \\ \vdots \\ \frac{\partial C}{\partial \mathbf{o}_m^s} \end{pmatrix}$$

Dans ce qui suit, quand on demande de "calculer" des gradients ou dérivées partielles, il s'agit simplement d'exprimer leur calcul en fonction d'éléments déjà calculés aux questions précédentes (ne substituez pas les expressions de dérivées partielles déjà calculées lors des questions d'avant!)

- g) Calculez les dérivées partielles par rapport aux activations des neurones de sortie (aussi appelées *sensibilité* des neurones de sortie):

$$\frac{\partial C}{\partial \mathbf{o}^a} = \begin{pmatrix} \frac{\partial C}{\partial \mathbf{o}_1^a} \\ \vdots \\ \frac{\partial C}{\partial \mathbf{o}_m^a} \end{pmatrix}$$

Notez que comme C ne dépend des activations \mathbf{o}^a qu'indirectement par l'intermédiaire des sorties \mathbf{o}^s il nous faut utiliser la règle de dérivation en chaîne:

$$\frac{\partial C}{\partial \mathbf{o}_k^a} = \sum_{k'=1}^m \frac{\partial C}{\partial \mathbf{o}_{k'}^s} \frac{\partial \mathbf{o}_{k'}^s}{\partial \mathbf{o}_k^a}$$

- h) Calculez les gradients par rapport aux paramètres $\mathbf{W}^{(2)}$ et $\mathbf{b}^{(2)}$ de la couche de sortie. Comme C ne dépend des $\mathbf{W}_{kj}^{(2)}$ et $\mathbf{b}_k^{(2)}$ qu'au travers de \mathbf{o}_k^a la règle de dérivation en chaîne nous donne:

$$\frac{\partial C}{\partial \mathbf{W}_{kj}^{(2)}} = \frac{\partial C}{\partial \mathbf{o}_k^a} \frac{\partial \mathbf{o}_k^a}{\partial \mathbf{W}_{kj}^{(2)}}$$

et

$$\frac{\partial C}{\partial \mathbf{b}_k^{(2)}} = \frac{\partial C}{\partial \mathbf{o}_k^a} \frac{\partial \mathbf{o}_k^a}{\partial \mathbf{b}_k^{(2)}}$$

- i) Calculez les dérivées partielles du coût C par rapport aux sorties des neurones de la couche cachée. Comme C dépend d'un neurone caché \mathbf{h}_j^s au travers des activations de tous les neurones de sortie \mathbf{o}^a reliés à ce neurone caché, la règle de dérivation en chaîne nous donne:

$$\frac{\partial C}{\partial \mathbf{h}_j^s} = \sum_{k=1}^m \frac{\partial C}{\partial \mathbf{o}_k^a} \frac{\partial \mathbf{o}_k^a}{\partial \mathbf{h}_j^s}$$

- j) Calculez les dérivées partielles par rapport aux activations des neurones de la couche cachée. Comme C ne dépend de l'activation \mathbf{h}_j^a d'un neurone de la couche cachée qu'au travers de la sortie \mathbf{h}_j^s de ce neurone, la règle de dérivation en chaîne donne:

$$\frac{\partial C}{\partial \mathbf{h}_j^a} = \frac{\partial C}{\partial \mathbf{h}_j^s} \frac{\partial \mathbf{h}_j^s}{\partial \mathbf{h}_j^a}$$

Notez que $\mathbf{h}_j^s = \tanh(\mathbf{h}_j^a)$, où la tangente hyperbolique s'applique élément par élément. La formule de la tangente hyperbolique est: $\tanh(z) = \frac{\sinh z}{\cosh z} = \frac{e^z - e^{-z}}{e^z + e^{-z}} = \frac{e^{2z} - 1}{e^{2z} + 1}$. Comme étape intermédiaire, commencez par démontrer que $\frac{\partial \tanh(z)}{\partial z} = \tanh'(z) = 1 - \tanh^2(z)$.

- k) Calculez les gradients par rapport aux paramètres $\mathbf{W}^{(1)}$ et $\mathbf{b}^{(1)}$ de la couche cachée.
- l) Calculez les dérivées partielles du coût C par rapport au vecteur d'entrée \mathbf{x} .
- m) Nous allons maintenant considérer un risque empirique **régularisé**: $\tilde{R} = \hat{R} + \lambda \mathcal{L}(\theta)$, où θ est le vecteur de tous les paramètres du réseau et $\mathcal{L}(\theta)$ calcule une pénalité scalaire en fonction des paramètres θ , plus ou moins importante selon une préférence à priori qu'on a sur les valeurs de θ . λ est un hyper-paramètre (scalaire, positif ou nul) qui contrôle le compromis entre trouver des valeurs des paramètres qui minimisent le risque empirique ou qui minimisent cette pénalité. On va considérer ici une régularization de type "weight decay" quadratique qui pénalise la norme carrée (norme L_2) des poids (mais pas des biais):

$$\begin{aligned} \mathcal{L}(\theta) &= \|\mathbf{W}^{(1)}\|^2 + \|\mathbf{W}^{(2)}\|^2 \\ &= \sum_{i,j} \left(\mathbf{w}_{ij}^{(1)} \right)^2 + \sum_{i',j'} \left(\mathbf{w}_{i'j'}^{(2)} \right)^2 \end{aligned}$$

On veut en fait minimiser le risque régularisé \tilde{R} plutôt que \hat{R} . Que devient le gradient par rapport aux différents paramètres?