

IFT-3395 Devoir 3

Olivier St-Laurent, Maxime Daigle

2018-11-09

Question 1 Relations et dérivées de quelques fonction de base

1. Montrez que $\text{sigmoid}(x) = \frac{1}{2}(\tanh(\frac{1}{2}x) + 1)$

$$\text{sigmoid}(x) = \frac{1}{2}(\tanh(\frac{1}{2}x) + 1) \iff 2\text{sigmoid}(x) - 1 = \tanh(\frac{x}{2})$$

$$\begin{aligned} 2\text{sigmoid}(x) - 1 &= \frac{2 - 1 - \exp(-x)}{1 + \exp(-x)} = \frac{1 - \exp(-x)}{1 + \exp(x)} = \frac{\exp(\frac{x}{2})}{\exp(\frac{x}{2})} \left(\frac{1 - \exp(-x)}{1 + \exp(-x)} \right) \\ &= \frac{\exp(\frac{x}{2}) - \exp(\frac{x}{2})\exp(-x)}{\exp(\frac{x}{2}) + \exp(\frac{x}{2})\exp(-x)} = \frac{\exp(\frac{x}{2}) - \exp(\frac{-x}{2})}{\exp(\frac{x}{2}) + \exp(\frac{-x}{2})} = \tanh(\frac{x}{2}) \end{aligned}$$

2. Montez que $\ln \text{sigmoid}(x) = -\text{softplus}(-x)$

$$\ln \text{sigmoid}(x) = \ln \frac{1}{1 + \exp(-x)} = \ln(1) - \ln(1 + \exp(-x)) = 0 + \ln(1 + \exp(-x)) = -\text{softplus}(-x)$$

3. Montrez que $\frac{d \text{sigmoid}}{dx}(x) = \text{sigmoid}(x)(1 - \text{sigmoid}(x))$

$$\begin{aligned} \frac{d \text{sigmoid}}{dx}(x) &= \frac{((1 + \exp(-x))^{-1})}{dx} = \frac{-1}{(1 + e^{-x})^2} (-e^{-x}) = \left(\frac{1}{1 + e^{-x}} \right) \left(\frac{e^{-x}}{1 + e^{-x}} \right) \\ &= \text{sigmoid}(x) \left(\frac{1 + e^{-x}}{1 + e^{-x}} - \frac{1}{1 + e^{-x}} \right) = \text{sigmoid}(x)(1 - \text{sigmoid}(x)) \end{aligned}$$

4. Montrez que la dérivée de tanh est : $\tanh'(x) = 1 - \tanh^2(x)$

$$\begin{aligned} \frac{d}{dx} \left(\frac{e^x - e^{-x}}{e^x + e^{-x}} \right) &= \frac{(e^x - e^{-x})'(e^x + e^{-x}) - (e^x - e^{-x})(e^x + e^{-x})'}{(e^x + e^{-x})^2} = \frac{(e^x + e^{-x})^2 - (e^x - e^{-x})^2}{(e^x + e^{-x})^2} \\ &= \frac{(e^x + e^{-x})^2}{(e^x + e^{-x})^2} - \frac{(e^x - e^{-x})^2}{(e^x + e^{-x})^2} = 1 - \tanh^2(x) \end{aligned}$$

5. Expliquer la relation entre le gradient du risque empirique et les erreurs du modèle sur l'ensemble d'entraînement

Le gradient du risque empirique permet de trouver le minimum global des erreurs du modèle sur l'ensemble d'entraînement. Autrement dit, le point où le gradient est égal à 0 est le point où les erreurs du modèle sont les plus petites sur l'ensemble d'entraînement.

Question 1.2 Régression linéaire régularisée (“ridge regression”)

1. Exprimez le gradient du risque régularisé. En quoi diffère-t-il du gradient du risque empirique non régularisé ?

$$\begin{aligned}\nabla \tilde{R}(\theta) &= \frac{\delta \tilde{R}(\theta)}{\delta \theta} \\&= \frac{\delta}{\delta \theta} (\hat{R} + \lambda \mathcal{L}(\theta)) \\&= \frac{\delta}{\delta \theta} \left(\sum_{i=1}^n (w^T x^{(i)} + b - t^{(i)})^2 + \lambda \sum_{k=1}^d w_k^2 \right) \\&= \frac{\delta}{\delta \theta} \left(\sum_{i=1}^n (w^T x^{(i)} + b - t^{(i)})^2 \right) + \frac{\delta}{\delta \theta} \left(\lambda \sum_{k=1}^d w_k^2 \right) \\&= \nabla \hat{R}(\theta) + \frac{\delta}{\delta \theta} \lambda \left(\sum_{k=1}^d w_k^2 \right)\end{aligned}$$

Donc, la différence entre $\nabla \tilde{R}(\theta)$ et $\nabla \hat{R}(\theta)$ est l'addition de $\frac{\delta}{\delta \theta} \lambda (\sum_{k=1}^d w_k^2)$

2. Donner le pseudo code détaillé de l'algorithme qui cherche les paramètres optimaux pour minimiser \tilde{R} par descente de gradient batch.

$$\text{loop } \theta \leftarrow \theta - \eta \frac{\delta \tilde{R}}{\delta \theta}$$

$$= \theta - \eta \left[\sum_{i=1}^n \frac{\delta}{\delta \theta} (w^T x^{(i)} + b - t^{(i)})^2 + \lambda \frac{\delta}{\delta \theta} \left(\sum_{k=1}^d w_k^2 \right) \right]$$

L'algorithme est donc :

$\theta = \text{initialize_randomly}()$

stop = False

while not stop :

$$\text{gradient} = \sum_{i=1}^n \frac{\delta}{\delta \theta} (w^T x^{(i)} + b - t^{(i)})^2 + \lambda \frac{\delta}{\delta \theta} \left(\sum_{k=1}^d w_k^2 \right)$$

$$\theta = \theta - \eta \cdot \text{gradient}$$

$$\text{stop} = || \text{gradient} || < \varepsilon$$

3. exprimer le risque empirique et son gradient sous forme matricielle

$$\theta^* = (X^T X)^{-1} X^T t$$

Question 1.3 Régression avec un pré-traitement non-linéaire fixe

1. Écrire de manière détaillée la forme paramétrique qu'on obtient pour $\tilde{f}(x)$ dans le cas une dimension ($x \in \mathbb{R}$) si on utilise $\phi = \phi_{poly^k}$

$$\tilde{f}(x) = f(\phi_{poly^k}(x)) = w^T \phi_{poly^k}(x) + b = w^T \begin{bmatrix} x \\ x^2 \\ \vdots \\ x^k \end{bmatrix} + b$$

2. Préciser quels sont les paramètres et leur dimensionalité

Les paramètres sont \mathbf{w} et \mathbf{b} . Pour $\phi_{poly^k} \in \mathbb{R}^k$, $w \in \mathbb{R}^k$ est un vecteur. $b \in \mathbb{R}$ est un scalaire.