



Introduction au Deep Learning avec PyTorch

Partie 1 : Introduction

Laurent Risser

Ingénieur de Recherche à l'Institut de Mathématiques de Toulouse et au 3IA ANITI

lrissier@math.univ-toulouse.fr

Partie 1 : Apprentissage supervisé et optimisation (*≈ 30 minutes*)

Partie 2 : Architecture et apprentissage des réseaux de neurones (*≈ 40 minutes*)

Partie 3 : Apprentissage avec PyTorch (*≈ 30 minutes*)

Partie 4 : TPs sous PyTorch (*≈ 60 minutes*)

On connait les notes de n élèves au cours de l'année scolaire 2019-2020 ainsi que leur notes à un concours de fin d'année.

L'année suivante, **on aimerait prédire** les notes au concours de la promotion 2020-2021 en fonction de leurs notes au cours de l'année.

	Maths	Info	Français	Concours
Elève 1	12	15	09	14
Elève 2	05	09	12	07
...
Elève i	x_i^1	x_i^2	x_i^3	y_i
...
Elève n	10	12	15	11

	Maths	Info	Français	Concours
Nouvel élève	13	14	11	?

On connait les notes de n élèves au cours de l'année scolaire 2019-2020 ainsi que leur notes à un concours de fin d'année.

L'année suivante, **on aimerait prédire** les notes au concours de la promotion 2020-2021 en fonction de leurs notes au cours de l'année.

	Maths	Info	Français	Concours
Elève 1	12	15	09	14
Elève 2	05	09	12	07
...
Elève i	x_i^1	x_i^2	x_i^3	y_i
...
Elève n	10	12	15	11

	Maths	Info	Français	Concours
Nouvel élève	13	14	11	?

Posons les notations :

- Notes de l'élève $i \in \{1, \dots, n\}$ durant l'année 2019-2020 : $x_i = (x_i^1, x_i^2, \dots, x_i^p) \in \mathbb{R}^p$
- Notes de l'élève i au concours 2019-2020 : $y_i \in \mathbb{R}$
- Notes de l'élève new durant l'année 2020-2021 : $x_{new} = (x_{new}^1, x_{new}^2, \dots, x_{new}^p) \in \mathbb{R}^p$
- Fonction pour prédire y_{new} en fonction des x_{new} : $\widehat{y_{new}} = h_{\Theta}(x_{new})$

Prédiction de y_{new}

Fonction $\mathbb{R}^p \rightarrow \mathbb{R}$ avec les paramètres Θ que l'on va apprendre avec les $(x_i, y_i)_{i=1, \dots, n}$

Notations : n observations d'apprentissage : $x_i = (x_i^1, x_i^2, \dots, x_i^p)$ et y_i avec $i = 1, \dots, n$ et fonction de prédiction h_Θ de paramètres Θ

1ere question clé : Quel modèle choisir pour h_Θ ?

- Nous allons commencer un modèle très simple : La **régression linéaire** !

$$\widehat{y}_{new} = h_\Theta(x_{new}) = w_0 + \sum_{j=1}^p w_j x_{new}^j \quad \text{dont les paramètres sont } \Theta = \{w_0, w_1, \dots, w_p\}$$

- Prédiction pour le nouvel élève ?

	Maths	Info	Français	Concours
Nouvel élève	13	14	11	?

Prenons $w_0 = 0$, $w_1 = 0.33$, $w_2 = 0.33$ et $w_3 = 0.33$ → Alors $\widehat{y}_{new} = 12.54$

Notations : n observations d'apprentissage : $x_i = (x_i^1, x_i^2, \dots, x_i^p)$ et y_i avec $i = 1, \dots, n$ et fonction de prédiction h_Θ de paramètres Θ

2eme question clé : Critère pour apprendre les paramètres Θ ?

Les meilleurs paramètres $\hat{\Theta}$ minimisent un **risque empirique** sur les $(x_i, y_i)_{i=1, \dots, n}$, par exemple :

$$\begin{aligned}
 \hat{\Theta} &= \arg \min_{\Theta = \{w_0, \dots, w_p\}} \frac{1}{n} \sum_{i=1}^n \text{loss}(h_\Theta(x_i), y_i) \\
 &= \arg \min_{\Theta = \{w_0, \dots, w_p\}} \frac{1}{n} \sum_{i=1}^n (h_\Theta(x_i) - y_i)^2 \\
 &= \arg \min_{\Theta = \{w_0, \dots, w_p\}} \frac{1}{n} \sum_{i=1}^n \left(w_0 + \sum_{j=1}^p w_j x_i^j - y_i \right)^2
 \end{aligned}$$

Risque empirique $R_\Theta((x_i, y_i)_{i=1, \dots, n})$

Notations : n observations d'apprentissage : $x_i = (x_i^1, x_i^2, \dots, x_i^p)$ et y_i avec $i = 1, \dots, n$ et fonction de prédiction h_{Θ} de paramètres Θ

2eme question clé : Critère pour apprendre les paramètres Θ ?

Les meilleurs paramètres $\hat{\Theta}$ minimisent un **risque empirique** sur les $(x_i, y_i)_{i=1, \dots, n}$, par exemple :

$$\begin{aligned} \hat{\Theta} &= \arg \min_{\Theta = \{w_0, \dots, w_p\}} \frac{1}{n} \sum_{i=1}^n \text{loss}(h_{\Theta}(x_i), y_i) \\ &= \arg \min_{\Theta = \{w_0, \dots, w_p\}} \frac{1}{n} \sum_{i=1}^n (h_{\Theta}(x_i) - y_i)^2 \\ &= \arg \min_{\Theta = \{w_0, \dots, w_p\}} \frac{1}{n} \sum_{i=1}^n \left(w_0 + \sum_{j=1}^p w_j x_i^j - y_i \right)^2 \end{aligned}$$

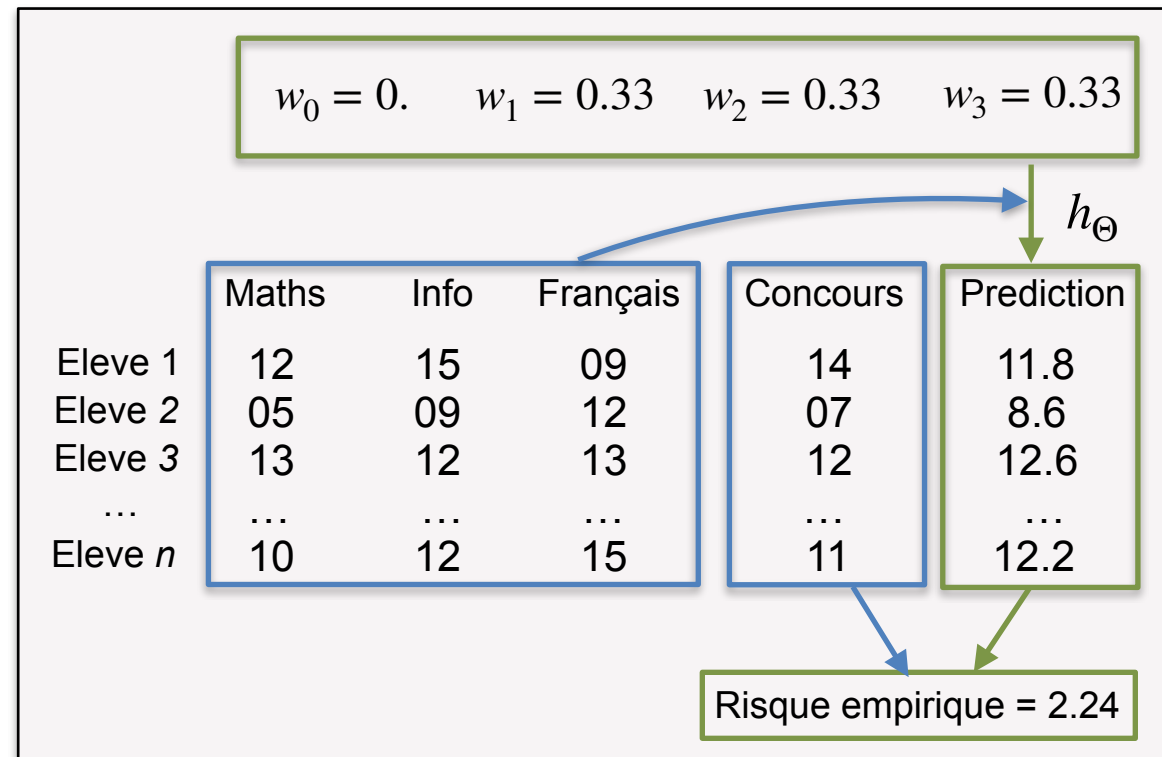
	Maths	Info	Français	Concours
Eleve 1	12	15	09	14
Eleve 2	05	09	12	07
Eleve 3	13	12	13	12
...
Eleve n	10	12	15	11

Notations : n observations d'apprentissage : $x_i = (x_i^1, x_i^2, \dots, x_i^p)$ et y_i avec $i = 1, \dots, n$ et fonction de prédiction h_Θ de paramètres Θ

2eme question clé : Critère pour apprendre les paramètres Θ ?

Les meilleurs paramètres $\hat{\Theta}$ minimisent un **risque empirique** sur les $(x_i, y_i)_{i=1, \dots, n}$, par exemple :

$$\begin{aligned} \hat{\Theta} &= \arg \min_{\Theta = \{w_0, \dots, w_p\}} \frac{1}{n} \sum_{i=1}^n \text{loss}(h_\Theta(x_i), y_i) \\ &= \arg \min_{\Theta = \{w_0, \dots, w_p\}} \frac{1}{n} \sum_{i=1}^n (h_\Theta(x_i) - y_i)^2 \\ &= \arg \min_{\Theta = \{w_0, \dots, w_p\}} \frac{1}{n} \sum_{i=1}^n \left(w_0 + \sum_{j=1}^p w_j x_i^j - y_i \right)^2 \end{aligned}$$

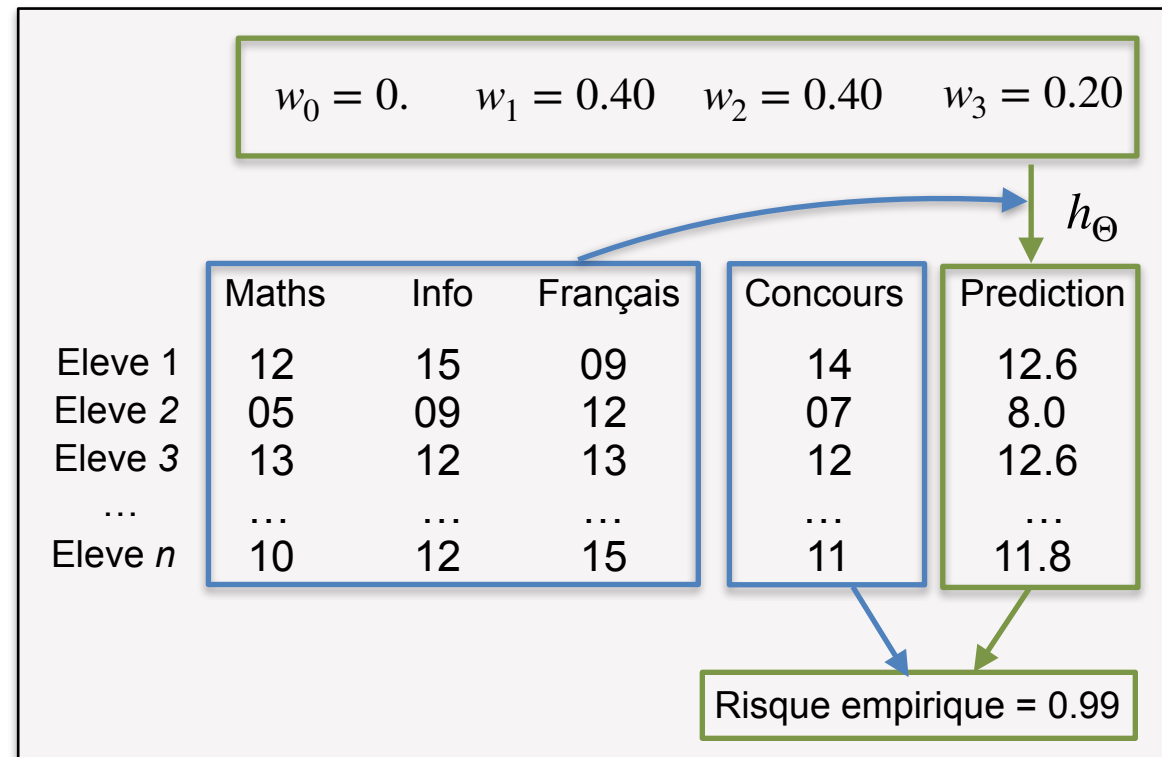


Notations : n observations d'apprentissage : $x_i = (x_i^1, x_i^2, \dots, x_i^p)$ et y_i avec $i = 1, \dots, n$ et fonction de prédiction h_Θ de paramètres Θ

2eme question clé : Critère pour apprendre les paramètres Θ ?

Les meilleurs paramètres $\hat{\Theta}$ minimisent un **risque empirique** sur les $(x_i, y_i)_{i=1, \dots, n}$, par exemple :

$$\begin{aligned} \hat{\Theta} &= \arg \min_{\Theta = \{w_0, \dots, w_p\}} \frac{1}{n} \sum_{i=1}^n \text{loss}(h_\Theta(x_i), y_i) \\ &= \arg \min_{\Theta = \{w_0, \dots, w_p\}} \frac{1}{n} \sum_{i=1}^n (h_\Theta(x_i) - y_i)^2 \\ &= \arg \min_{\Theta = \{w_0, \dots, w_p\}} \frac{1}{n} \sum_{i=1}^n \left(w_0 + \sum_{j=1}^p w_j x_i^j - y_i \right)^2 \end{aligned}$$

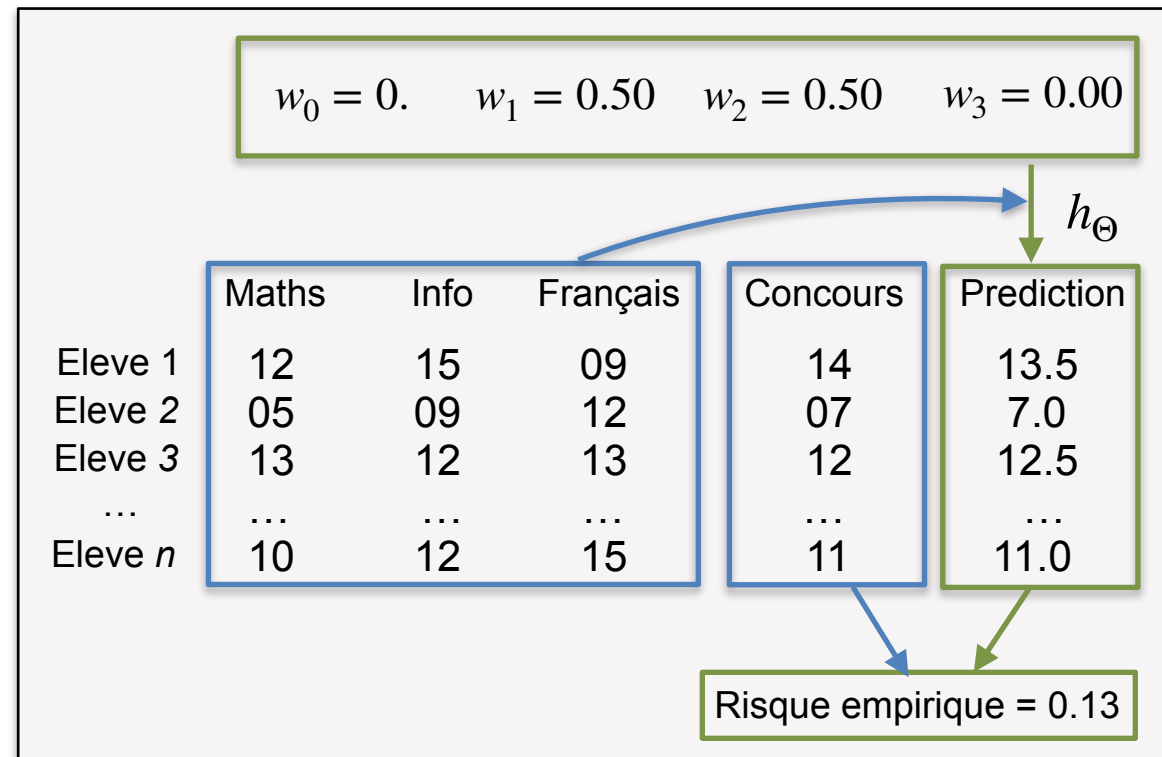


Notations : n observations d'apprentissage : $x_i = (x_i^1, x_i^2, \dots, x_i^p)$ et y_i avec $i = 1, \dots, n$ et fonction de prédiction h_Θ de paramètres Θ

2eme question clé : Critère pour apprendre les paramètres Θ ?

Les meilleurs paramètres $\hat{\Theta}$ minimisent un **risque empirique** sur les $(x_i, y_i)_{i=1, \dots, n}$, par exemple :

$$\begin{aligned} \hat{\Theta} &= \arg \min_{\Theta = \{w_0, \dots, w_p\}} \frac{1}{n} \sum_{i=1}^n \text{loss}(h_\Theta(x_i), y_i) \\ &= \arg \min_{\Theta = \{w_0, \dots, w_p\}} \frac{1}{n} \sum_{i=1}^n (h_\Theta(x_i) - y_i)^2 \\ &= \arg \min_{\Theta = \{w_0, \dots, w_p\}} \frac{1}{n} \sum_{i=1}^n \left(w_0 + \sum_{j=1}^p w_j x_i^j - y_i \right)^2 \end{aligned}$$



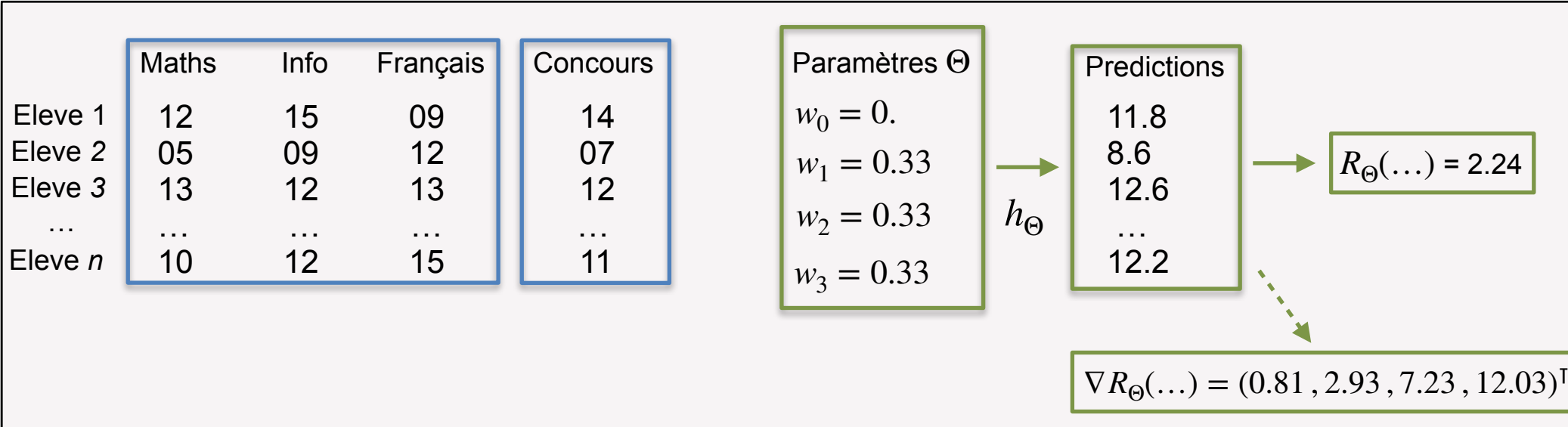
Notations : n observations d'apprentissage : $x_i = (x_i^1, x_i^2, \dots, x_i^p)$ et y_i avec $i = 1, \dots, n$ et fonction de prédiction h_Θ de paramètres Θ

3eme question clé : Minimisation automatique du risque en fonction de Θ

$$R_\Theta((x_i, y_i)_{i=1, \dots, n}) = \frac{1}{n} \sum_{i=1}^n (h_\Theta(x_i) - y_i)^2$$

$$\nabla R_\Theta(\dots) = \left(\frac{\partial R_\Theta(\dots)}{\partial w_0}, \frac{\partial R_\Theta(\dots)}{\partial w_1}, \dots, \frac{\partial R_\Theta(\dots)}{\partial w_p} \right)^T = \frac{1}{n} \sum_{i=1}^n \left(2(h_\Theta(x_i) - y_i), 2x_i^1(h_\Theta(x_i) - y_i), \dots, 2x_i^p(h_\Theta(x_i) - y_i) \right)^T$$

↑
Gradient de R_Θ



Notations : n observations d'apprentissage : $x_i = (x_i^1, x_i^2, \dots, x_i^p)$ et y_i avec $i = 1, \dots, n$ et fonction de prédiction h_{Θ} de paramètres Θ

3eme question clé : Minimisation automatique du risque en fonction de Θ

$$R_{\Theta}((x_i, y_i)_{i=1, \dots, n}) = \frac{1}{n} \sum_{i=1}^n (h_{\Theta}(x_i) - y_i)^2 \quad \nabla R_{\Theta}(\dots) = \frac{2}{n} \sum_{i=1}^n \left((h_{\Theta}(x_i) - y_i), \dots, x_i^p (h_{\Theta}(x_i) - y_i) \right)^{\top}$$

Principe de la descente de gradient :

$$\Theta_{it1} = (0.0, 0.33, 0.33, 0.33)^{\top}$$



$$R_{\Theta_{it1}}(\dots) = 2.24$$

$$\nabla R_{\Theta_{it1}}(\dots) = (0.81, 2.93, 7.23, 12.03)^{\top}$$

$$\Theta_{it2} = \Theta_{it1} - \lambda \nabla R_{\Theta_{it1}}(\dots) = (-0.00, 0.33, 0.32, 0.32)^{\top}$$



$$R_{\Theta_{it2}}(\dots) = 2.09$$

$$\nabla R_{\Theta_{it2}}(\dots) = (0.28, -2.47, 1.23, 5.10)^{\top}$$



Notations : n observations d'apprentissage : $x_i = (x_i^1, x_i^2, \dots, x_i^p)$ et y_i avec $i = 1, \dots, n$ et fonction de prédiction h_{Θ} de paramètres Θ

3eme question clé : Minimisation automatique du risque en fonction de Θ

$$R_{\Theta}((x_i, y_i)_{i=1, \dots, n}) = \frac{1}{n} \sum_{i=1}^n (h_{\Theta}(x_i) - y_i)^2 \quad \nabla R_{\Theta}(\dots) = \frac{2}{n} \sum_{i=1}^n \left((h_{\Theta}(x_i) - y_i), \dots, x_i^p (h_{\Theta}(x_i) - y_i) \right)^{\top}$$

Principe de la descente de gradient :

$\Theta_{it1} = (0.0, 0.33, 0.33, 0.33)^{\top}$	→	$R_{\Theta_{it1}}(\dots) = 2.24$ $\nabla R_{\Theta_{it1}}(\dots) = (0.81, 2.93, 7.23, 12.03)^{\top}$
$\Theta_{it2} = \Theta_{it1} - \lambda \nabla R_{\Theta_{it1}}(\dots) = (-0.00, 0.33, 0.32, 0.32)^{\top}$	→	$R_{\Theta_{it2}}(\dots) = 2.09$ $\nabla R_{\Theta_{it2}}(\dots) = (0.28, -2.47, 1.23, 5.10)^{\top}$
$\Theta_{it3} = \Theta_{it2} - \lambda \nabla R_{\Theta_{it2}}(\dots) = (-0.00, 0.33, 0.32, 0.31)^{\top}$	→	$R_{\Theta_{it3}}(\dots) = 2.04$ $\nabla R_{\Theta_{it3}}(\dots) = (0.17, -3.48, 0.05, 3.71)^{\top}$

Notations : n observations d'apprentissage : $x_i = (x_i^1, x_i^2, \dots, x_i^p)$ et y_i avec $i = 1, \dots, n$ et fonction de prédiction h_{Θ} de paramètres Θ

3eme question clé : Minimisation automatique du risque en fonction de Θ

$$R_{\Theta}((x_i, y_i)_{i=1, \dots, n}) = \frac{1}{n} \sum_{i=1}^n (h_{\Theta}(x_i) - y_i)^2 \qquad \nabla R_{\Theta}(\dots) = \frac{2}{n} \sum_{i=1}^n \left((h_{\Theta}(x_i) - y_i), \dots, x_i^p (h_{\Theta}(x_i) - y_i) \right)^{\top}$$

Principe de la descente de gradient :

$$\Theta_{it1} = (0.0, 0.33, 0.33, 0.33)^{\top}$$



$$R_{\Theta_{it1}}(\dots) = 2.24$$

$$\nabla R_{\Theta_{it1}}(\dots) = (0.81, 2.93, 7.23, 12.03)^{\top}$$

$$\Theta_{it2} = \Theta_{it1} - \lambda \nabla R_{\Theta_{it1}}(\dots) = (-0.00, 0.33, 0.32, 0.32)^{\top}$$



$$R_{\Theta_{it2}}(\dots) = 2.09$$

$$\nabla R_{\Theta_{it2}}(\dots) = (0.28, -2.47, 1.23, 5.10)^{\top}$$

$$\Theta_{it3} = \Theta_{it2} - \lambda \nabla R_{\Theta_{it2}}(\dots) = (-0.00, 0.33, 0.32, 0.31)^{\top}$$



$$R_{\Theta_{it3}}(\dots) = 2.04$$

$$\nabla R_{\Theta_{it3}}(\dots) = (0.17, -3.48, 0.05, 3.71)^{\top}$$

⋮

$$\Theta_{it1000} = \Theta_{it999} - \lambda \nabla R_{\Theta_{it999}}(\dots) = (-0.00, 0.48, 0.53, -0.01)^{\top}$$



$$R_{\Theta_{it3}}(\dots) = 0.08$$

$$\nabla R_{\Theta_{it3}}(\dots) = (-0.01, 0.09, -0.12, 0.04)^{\top}$$

Notations : n observations d'apprentissage : $x_i = (x_i^1, x_i^2, \dots, x_i^p)$ et y_i avec $i = 1, \dots, n$ et fonction de prédiction h_Θ de paramètres Θ

3eme question clé : Minimisation automatique du risque en fonction de Θ

$$R_\Theta((x_i, y_i)_{i=1, \dots, n}) = \frac{1}{n} \sum_{i=1}^n (h_\Theta(x_i) - y_i)^2 \qquad \nabla R_\Theta(\dots) = \frac{2}{n} \sum_{i=1}^n \left((h_\Theta(x_i) - y_i), \dots, x_i^p (h_\Theta(x_i) - y_i) \right)^T$$

Score à convergence :

$$\Theta_{it1000} = \Theta_{it999} - \lambda \nabla \Theta_{it999}(\dots) = (-0.00, 0.48, 0.53, -0.01)^T \quad \rightarrow \quad R_{\Theta_{it3}}(\dots) = 0.08$$

On a convergé vers une paramétrisation qui marche très bien **sur le jeu d'apprentissage** !

Prédiction :

	Maths	Info	Français	Concours
Nouvel élève	13	14	11	$\widehat{y}_{new} = 13.49$

... dans la vrai vie, il faudra aussi valider l'apprentissage avant de l'utiliser sur de nouvelles données

Notations : n observations d'apprentissage : $x_i = (x_i^1, x_i^2, \dots, x_i^p)$ et y_i avec $i = 1, \dots, n$ et fonction de prédiction h_{Θ} de paramètres Θ

Pour aller plus loin : Descente de gradient stochastique

$$\begin{aligned}\nabla R_{\Theta}(\dots) &= \frac{2}{n} \sum_{i=1}^n \left((h_{\Theta}(x_i) - y_i), \dots, x_i^p (h_{\Theta}(x_i) - y_i) \right)^{\top} \\ &\approx \frac{2}{\#B} \sum_{i \in B} \left((h_{\Theta}(x_i) - y_i), \dots, x_i^p (h_{\Theta}(x_i) - y_i) \right)^{\top}\end{aligned}$$

Mini-batch (exemple $B = \{1,2,3\}$)

- Meilleure rapidité !
- Moins de mémoire nécessaire !

Batch-training → descente de gradient avec un mini-batch qui évolue itération après itération

- Meilleure exploration de l'espace des paramètres !

Notations : n observations d'apprentissage : $x_i = (x_i^1, x_i^2, \dots, x_i^p)$ et y_i avec $i = 1, \dots, n$ et fonction de prédiction h_Θ de paramètres Θ

Pour aller plus loin : Descente de gradient stochastique

Initialise $\Theta = (w_0, w_1, \dots, w_p)$

For $epoch = 1$ to **5**

For $b = 1$ to n with step **10**

$B = \{b, b + 1, \dots, b + 9\}$

Approche $\nabla R_\Theta(\dots)$ avec $\nabla R_\Theta(\dots) \approx \frac{2}{\#B} \sum_{i \in B} \left((h_\Theta(x_i) - y_i), \dots, x_i^p (h_\Theta(x_i) - y_i) \right)$

$\Theta = \Theta - \lambda \nabla R_\Theta(\dots)$

End For

End For

Descente de gradient avec **5 epochs** et des **mini-batches** de taille **10**

