

advanced optimization algorithms

① gradient descent | momentum = 0.9
Learning rate: fixed or Learning Rate Schedule

$$w_j = w_j - \underbrace{\alpha}_{\text{learning rate}} \cdot \frac{\partial}{\partial w_j} f(w, b)$$

② ADAM algorithm

adaptive moment estimation

not just one α

$$w_1 = w_1 - \alpha_1 \frac{\partial}{\partial w_1} f(w, b)$$

⋮

$$w_{10} = w_{10} - \alpha_{10} \frac{\partial}{\partial w_j} f(w, b)$$

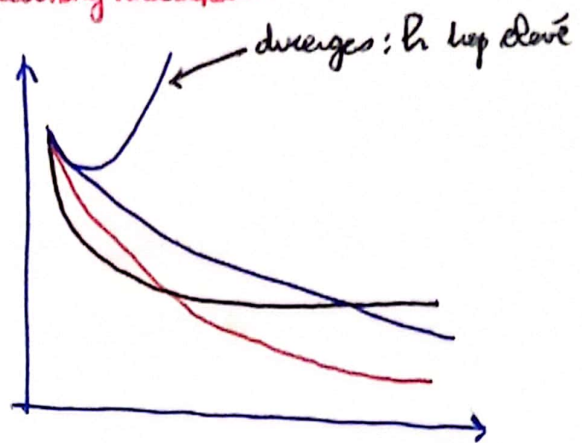
$$b = b - \alpha_{11} \frac{\partial}{\partial w_j} f(w, b)$$

⇒ ① if w_j keeps moving in same direction

⇒ increase α_j ⇒ go faster in that direction

② if w_j keeps oscillating =

⇒ reduce α_j



ordre croissant du lr :

bleue → rouge → verte

Optimization algorithms

Batch VS mini-batch gradient descent:

Vectorization allows you to efficiently compute on m examples

$$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix}$$

$$y = [y^{(1)}, \dots, y^{(m)}]$$

if we have $m = 5,000,000$?

the algo will be slow

↓ solution

split the entire dataset into mini-batches of 1,000 examples each

$$\begin{array}{l} X^{[1]} : \text{first mini batch} \\ X^{[2]} : \text{second mini batch} \\ \vdots \\ X^{[5000]} : 5000 - \text{mini batch} \end{array} \quad \begin{array}{l} y^{[1000]} \\ y^{[2000]} \\ \vdots \\ y^{[5000]} \end{array}$$

↓

for $t = 1$ to 5000:

• forward prop on $X^{[t]}$

$$z^{[t]} = w^{[t]} X^{[t]} + b^{[t]}$$

$$A^{[t]} = g^{[t]}(z^{[t]})$$

$$A^{[t]} = g^{[t]}(z^{[t]})$$

• compute cost: $J = \frac{1}{1000} \sum_{i=1}^{1000} \mathcal{L}(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum \|w^{[t]}\|_F^2$

size of the mini batch

from $X^{[t]}$ and $y^{[t]}$

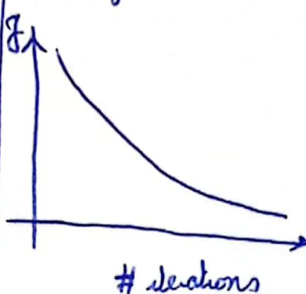
• Backpropagation

$$w^{[t]} := w^{[t]} - \alpha dw^{[t]}$$

}

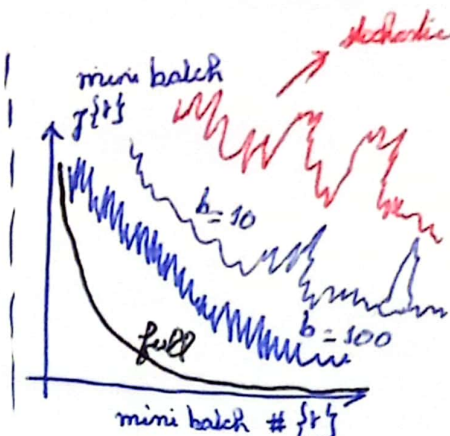
this step is called 1 epoch: one pass through all the training set (one iteration of gradient descent)

Batch gradient descent



the cost J must decrease in every iteration

1 iteration per epoch



mini-batch the cost might not decrease in every iteration

but if we plot $J^{[t]}$ $\left| \begin{array}{l} n \text{ iterations} \\ m \text{ per epoch} \end{array} \right.$

parameter to choose: mini batch size

if mini batch size = m : batch gradient descent

$$\{X^{[t]}, y^{[t]}\} = \{X, y\}$$

if mini batch size = 1: stochastic gradient descent: n iterations per epoch

every example is its own mini batch (doesn't converge to the min: oscillates around min)

see opti course
lose speed from vectorization

In practice: mini batch should be between 1 and m

un mini batch trop petit entraîne un bruit trop grand sur l'estimation du gradient
⇒ empêche la convergence

⇒ best choice: 50-100 mini batch

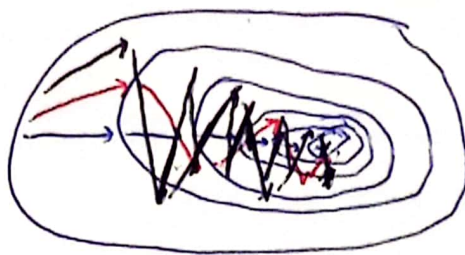
if we have small ~~batch~~ set : use batch gradient descent
($m \leq 2000$)

else :

typical mini batch size : 64, 128, 256, 512, 1024
(usually a power of 2)

⚠ make sure that mini batch fit in CPU/CPU memory

descente de gradient batch



descente de gradient mini batch

stochastique : $\nearrow \theta_j^{(k+1)} = \theta_j^{(k)} - \eta_k \frac{\partial \ell_i(\theta^{(k)})}{\partial \theta_j}$

convergence du stochastique gradient descent

1/ $\sum_k \eta_k < +\infty$ et $\sum_{k=1} \eta_k^2 = +\infty$

2/ des hypothèses de régularité sur J

Alors la suite $\theta^{(k)}$ converge vers un min local de J

Remarque

exponentially weighted averages

example 3 Temperature in London

$$\theta_1: 40^\circ\text{F} \quad 4^\circ\text{C}$$

$$\theta_2: 45^\circ\text{F} \quad 5^\circ\text{C}$$

$$\theta_{180} = 60^\circ\text{F} \quad 15^\circ\text{C}$$

$$\theta_{181} = 56^\circ\text{F}$$



$$V_0 = 0$$

$$V_1 = 0.9 V_0 + 0.1 \theta_1$$

$$V_2 = 0.9 V_1 + 0.1 \theta_2$$

$$V_t = 0.9 V_{t-1} + 0.1 \theta_t$$

if we plot V_t : we obtain exponentially weighted averages

$$V_t = \beta \cdot V_{t-1} + (1-\beta) \theta_t$$

$$\approx \frac{1}{1-\beta} \text{ days temperature}$$

$$\beta = 0.9 : \approx \text{averaging over } \frac{1}{1-0.9} = 10 \text{ days}$$

$$\beta = 0.98 : \approx \text{averaging over } \frac{1}{1-0.98} = 50 \text{ days}$$

momentum: pour réduire la variance et limiter les oscillations + éviter de redescendre après dans un min local

$$V_{100} = 0.5 \theta_{100} + 0.5 \times 0.98 \times 0.99 + 0.5 \times (0.9)^2 \theta_{98}$$

$$+ \dots + 0.5 \times (0.99)$$

$$(0.9)^{10} \approx 0.35 = \frac{1}{e}$$

$$(1-\epsilon)^{1/\epsilon} = \frac{1}{e}$$

Rian correction in

we start very slow due to initialization $V_0 = 0$
 \Rightarrow at the beginning

$$\text{instead of taking } V_t \text{ we take } \frac{V_t}{1 - \beta_0^t}$$

Gradient descent with momentum

we compute an exponentially weighted average on the gradients

\Rightarrow compute faster: tends to smooth the gradient
 descent: less oscillations

momentum:

on iteration t :

compute dw, db on current mini batch

$$V \begin{cases} Vdw = \beta \cdot Vdw + (1-\beta) dw \\ Vdb = \beta \cdot Vdb + (1-\beta) db \end{cases}$$

$$W := W - \alpha \cdot Vdw$$

$$b := b - \alpha \cdot Vdb$$

velocity: direction dans laquelle les paramètres vont être modifiés

hyperparameters: α, β

if $\beta = 0.9 \approx$ average over 10 gradients

RMS prop algorithm

on iteration t :

compute dw , db on current mini batch

$$Sdw = \beta Sdw + (1-\beta) \cdot \underbrace{dw^2}_{\text{element wise}}$$

$$Sdb = \beta Sdb + (1-\beta) \cdot db^2$$

$$w = w - \alpha \cdot \frac{dw}{\sqrt{Sdw}}$$

$$b = b - \alpha \cdot \frac{db}{\sqrt{Sdb}}$$

we hope that Sdw ~~is~~ to be small if we want learning to be fast in the direction of w

we hope that Sdb to be large if we want

learning to slow in the direction of b

Adam optimization algorithm

$$V_{dw} = 0 \quad S_{dw} = 0 \quad V_{db} = 0 \quad S_{db} = 0$$

on iteration t

compute dw, db using current minibatch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw$$

$$V_{db} = \beta_2 V_{db} + (1 - \beta_2) db$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) \underbrace{dw^2}_{\text{element wise}}$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) \cdot db^2 \text{ squaring}$$

$$V_{dw}^{\text{corrected}} = \frac{V_{dw}}{(1 - \beta_1^t)}$$

$$V_{db}^{\text{corrected}} = \frac{V_{db}}{(1 - \beta_2^t)}$$

$$S_{dw}^{\text{corrected}} = \frac{S_{dw}}{(1 - \beta_2^t)}$$

$$S_{db}^{\text{corrected}} = \frac{S_{db}}{(1 - \beta_2^t)}$$

$$w = w - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}} + \epsilon}}$$

$$b = b - \alpha \cdot \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}} + \epsilon}}$$

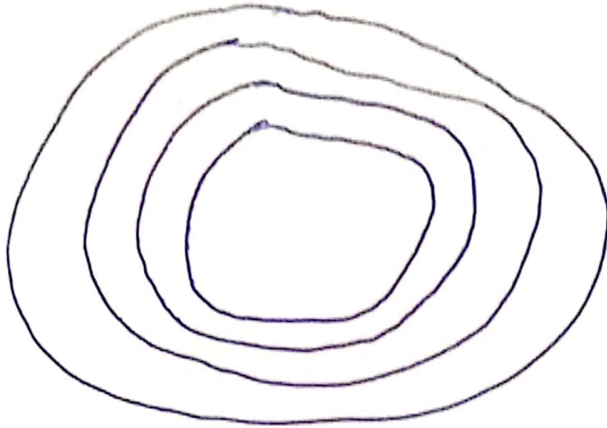
Hyperparameters: α : needs to be tuned

$$\left. \begin{array}{l} \beta_1 \approx 0.9 \\ \beta_2 \approx 0.99 \\ \epsilon \approx 10^{-8} \end{array} \right\} \begin{array}{l} \text{weak default} \\ \text{values usually} \end{array}$$

Adam: adaptive moment estimation

Learning rate decay

we slowly decrease alpha while learning



1 epoch = 1 pass through the data

$$\alpha = \frac{1}{1 + \text{decay-rate} \times \text{epoch num}} \cdot \alpha_0$$

other methods

exponential

$$\alpha = 0.95^{\text{epoch num}} \cdot \alpha_0$$

$$\alpha = \frac{K}{\sqrt{\text{epoch-num}}} \cdot \alpha_0$$

$$\alpha = \frac{K}{\sqrt{x}} \cdot \alpha_0$$