



deeplearning.ai

# Optimization Algorithms

---

## Mini-batch gradient descent

Applying machine learning is a highly empirical process, is a highly iterative process. In which you just had to train a lot of models to find one that works really well.

So it's really helps to train models quickly.

Training on large dataset is just slow. So what you find is that having fast optimization algorithms, having good optimization algorithm can really speed up your efficiency and your team.

Viewing that as processing your entire batch of training samples all at the same time.

In contrast, process is single mini batch  $X^{t+1}$ ,  $y^{t+1}$  at the same time

# Batch vs. mini-batch gradient descent

$x, y$

$x^{t+1}, y^{t+1}$

You have to process your entire training set before you take one little step of gradient descent. And the you have to process the entire training sets of five million training samples again before you take another step of gradient descent.

Vectorization allows you to efficiently compute on  $m$  examples.

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(500)} & | & x^{(1001)} & \dots & x^{(2000)} & | & \dots & | & \dots & x^{(m)} \end{bmatrix}$$

$(n_x, m)$   $X^{\{1\}}$   $(n_x, 1000)$   $X^{\{2\}}$   $(n_x, 1000)$   $X^{\{5,000\}}$   $(n_x, 1000)$

nx features #

X superscript with curly braces

$$y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(500)} & | & y^{(1001)} & \dots & y^{(2000)} & | & \dots & | & \dots & y^{(m)} \end{bmatrix}$$

$(1, m)$   $y^{\{1\}}$   $(1, 1000)$   $y^{\{2\}}$   $(1, 1000)$   $y^{\{5,000\}}$   $(1, 1000)$

What if  $m = 5,000,000$ ?

5,000 mini-batches of 1,000 each

Mini-batch  $t$ :  $x^{t+1}, y^{t+1}$

$x^{(i)}$  ith training sample

$z^{[l]}$  index the different layers of the neural network

$x^{\{t\}}, y^{\{t\}}$  index into different mini batches

# Mini-batch gradient descent

repeat {  
for  $t = 1, \dots, 5000$  { 5000 batches

Forward prop on  $X^{\{t\}}$ .

$$Z^{\{t\}} = W^{\{t\}} X^{\{t\}} + b^{\{t\}}$$

$$A^{\{t\}} = g^{\{t\}}(Z^{\{t\}})$$

Z Vectorizing connotation

$$\vdots$$

$$A^{\{t\}} = g^{\{t\}}(Z^{\{t\}})$$

Vectorized implementation  
(1000 examples)

Compute cost  $J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^L \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_{\mathbf{a}} \|W^{\{t\}}\|_F^2$

Backprop to compute gradients w.r.t  $J^{\{t\}}$  (using  $X^{\{t\}}, Y^{\{t\}}$ )

$$W^{\{t+1\}} := W^{\{t\}} - \alpha dW^{\{t\}}, \quad b^{\{t+1\}} := b^{\{t\}} - \alpha db^{\{t\}}$$

update

"1 epoch"

single

pass through training set.

1 step of grad desc  
using  $X^{\{t+1\}}, Y^{\{t+1\}}$ .  
(as if  $m=1000$ )

rather than having an explicit for loop over all 1,000 examples, you would use vectorization to process 1,000 examples sort of all at the same time.

X, Y

Frobenius norm which is square