



deeplearning.ai

Optimization Algorithms

Mini-batch gradient descent

Batch vs. mini-batch gradient descent

Vectorization allows you to efficiently compute on m examples.

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(1000)} & | & x^{(1001)} & \dots & x^{(2000)} & | & \dots & | & \dots & x^{(m)} \end{bmatrix}$$

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

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(1000)} & | & y^{(1001)} & \dots & y^{(2000)} & | & \dots & | & \dots & y^{(m)} \end{bmatrix}$$

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

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

5,000 mini-batches of 1,000 each

Mini-batch t : $X^{\{t\}}, Y^{\{t\}}$

$x^{(i)}$ i example

$z^{[l]}$ l layer

$X^{\{t\}}, Y^{\{t\}}$ t mini-batch

Mini-batch gradient descent

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

Forward prop on $X^{(t)}$.

$$Z^{(t)} = W^{(t)} X^{(t)} + b^{(t)}$$

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

\vdots

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

Vectorized implementation
(1000 examples)

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

$\swarrow \searrow$ from $X^{(t)}, Y^{(t)}$

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

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

1 step of gradient descent
using $X^{(t+1)}, Y^{(t+1)}$.
(as if $t=1000$)

X, Y

}
}

"1 epoch"

pass through training set.