



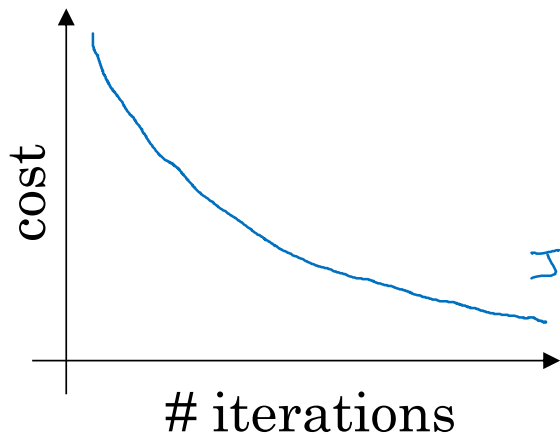
deeplearning.ai

Optimization Algorithms

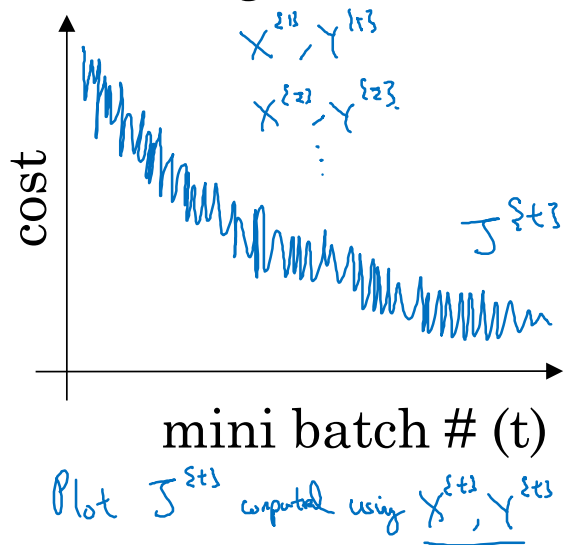
Understanding
mini-batch
gradient descent

Training with mini batch gradient descent

Batch gradient descent



Mini-batch gradient descent



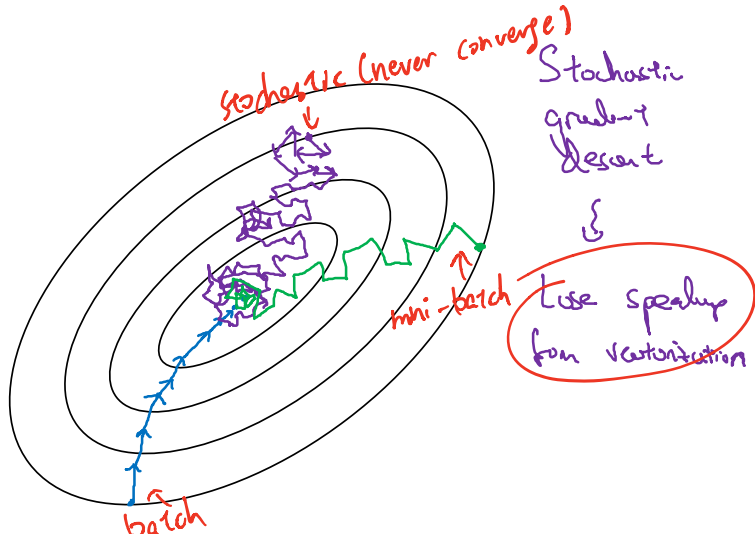
Choosing your mini-batch size

→ If mini-batch size = m : Batch gradient descent.

$$(X^{(1)}, Y^{(1)}) = (X, Y)$$

→ If mini-batch size = 1 : ^{Big m} Stochastic gradient descent. Every example is its own mini-batch.
 $(X^{(1)}, Y^{(1)}) = (x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$ mini-batch.

→ In practice: Somewhere in-between 1 and m



In-between
(mini-batch size
not too big/small)

Batch
gradient descent
(mini-batch size = m)

Fastest learning.

- Vectorization.
(~1000)
- Make passes without
processing entire training set.

Too long
per iteration

Choosing your mini-batch size

If small toy set : Use batch gradient descent.
($m \leq 2000$)

Typical mini-batch sizes:

→ 64, 128, 256, 512

2^6

2^7

2^8

2^9

common

2^n

$\frac{1024}{2^{10}}$

Make sure mini-batch fits in CPU/GPU memory.
 $X^{(t)}, Y^{(t)}$