deeplearning.ai
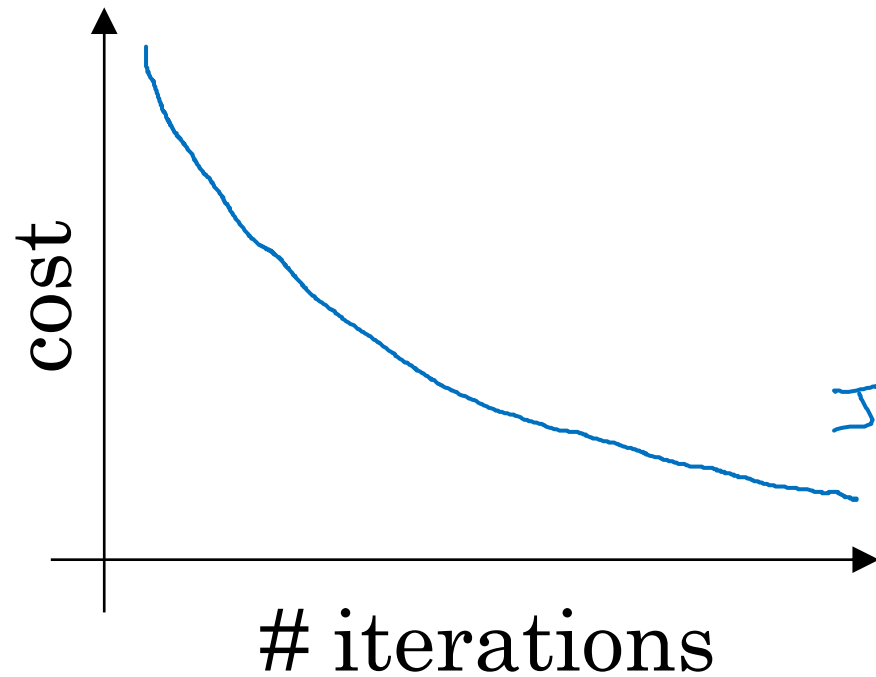
Optimization
Algorithms

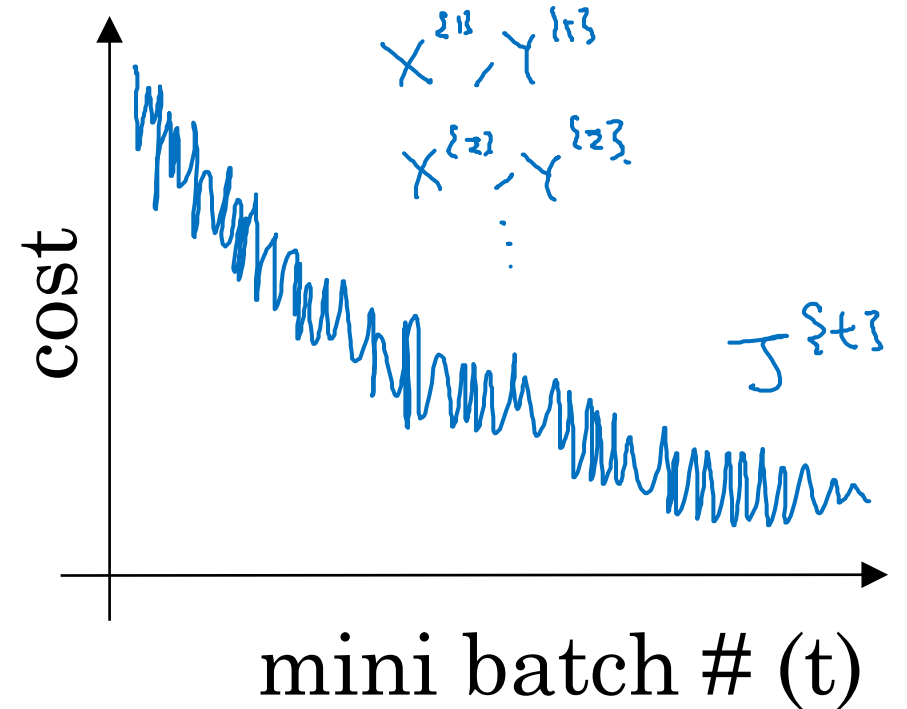Understanding
mini-batch
gradient descent

# Training with mini batch gradient descent

### Batch gradient descent



cost vs. # iterations — $J$

### Mini-batch gradient descent



cost vs. mini batch # (t)

$X^{\{1\}}, Y^{\{1\}}$

$X^{\{2\}}, Y^{\{2\}}$

$J^{\{t\}}$

Plot $J^{\{t\}}$ computed using $X^{\{t\}}, Y^{\{t\}}$

# Choosing your mini-batch size

$\rightarrow$ If mini-batch size $= m$ : Batch gradient descent. $(X^{\{1\}}, Y^{\{1\}}) = (X, Y)$

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

In practice: Somewhn in-between $\underline{1}$ and $\underline{m}$



Stochastic
gradient
descent
$\}$

Lose speedup
from vectorization

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

Fastest learning.
• Vectorization.
  ($\sim 1000$)
• Make progress without
  processing entire tray set.

Batch
gradient descent
(mini-batch size $= m$)
$\}$

Too long
per iteration

Andrew Ng

# Choosing your mini-batch size

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

Typical mini-batch sizes:

$$\rightarrow \quad \underbrace{64 \quad , \quad 128, \quad 256, \quad 512}_{2^6 \qquad 2^7 \qquad 2^8 \qquad 2^9} \qquad \frac{1024}{2^{10}}$$

Make sure mini-batch fit in CPU/GPU memory.
$$X^{\{t\}}, Y^{\{t\}}$$