# Optimization Algorithms

---

# Understanding mini-batch gradient descent

deeplearning.ai

# Training with mini batch gradient descent

## Batch gradient descent



cost

# iterations

$J$

## Mini-batch gradient descent



cost

mini batch # (t)

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

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

osci

$J^{\{t\}}$

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

It's ok if it doesn't go down on every iteration. But it should trend downwards

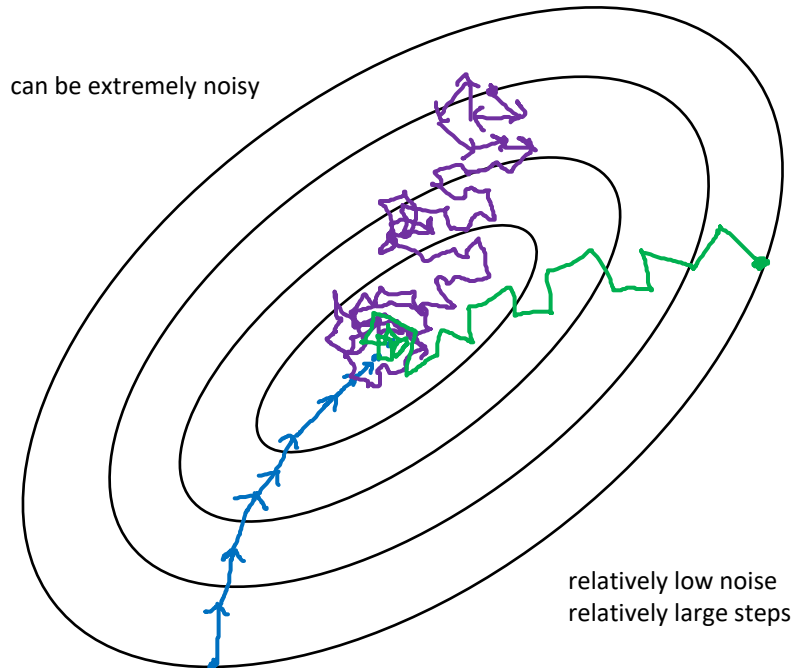Andrew Ng

# Choosing your mini-batch size

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

$\rightarrow$ If min-batch size = 1 : Stochastic gradient descent. Every example is it own

stochastic gradient descent

$(X^{\{1\}}, Y^{\{1\}}) = (x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$ min-batch.

In practice: Somewhn in-between 1 and m

As stochastic gradient descent won't ever converge, it'll always just kind of oscillate and wander around the region of the minimum. But it won't ever just head to the minimum and stay there

can be extremely noisy

relatively low noise
relatively large steps

Stochastic
gradient
descent

noise can ameliorated
or can be reduced by
just using a smaller
learning rate.

Lose speedup
from vectorization

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

Fastest learning.
- Vectorization.
  (~1 000)
- Make progress without
  processing entire training set.

And then it doesn't always exactly converge or oscillate in a very small region. if that's the issue you can always reduce the learning rate slowly.

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:

Because of the way computer memory is laid out and accessed, sometimes your code runs faster if you mini-batch size is a power of 2.

$\longrightarrow$ $64$, $128$, $256$, $512$ $\qquad \dfrac{1024}{2^{10}}$

$2^6 \qquad 2^7 \qquad 2^8 \qquad 2^9$

a little bit common

Make sure mini-batch fit in CPU/GPU memory.

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

Andrew Ng