# Training with mini batch gradient descent



Batch gradient descent

cost

$J$

# iterations

Mini-batch gradient descent

cost

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

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

$J^{\{t\}}$

mini batch # (t)

Plot $J^{\{t\}}$ computed using $X^{\{t\}}, Y^{\{t\}}$
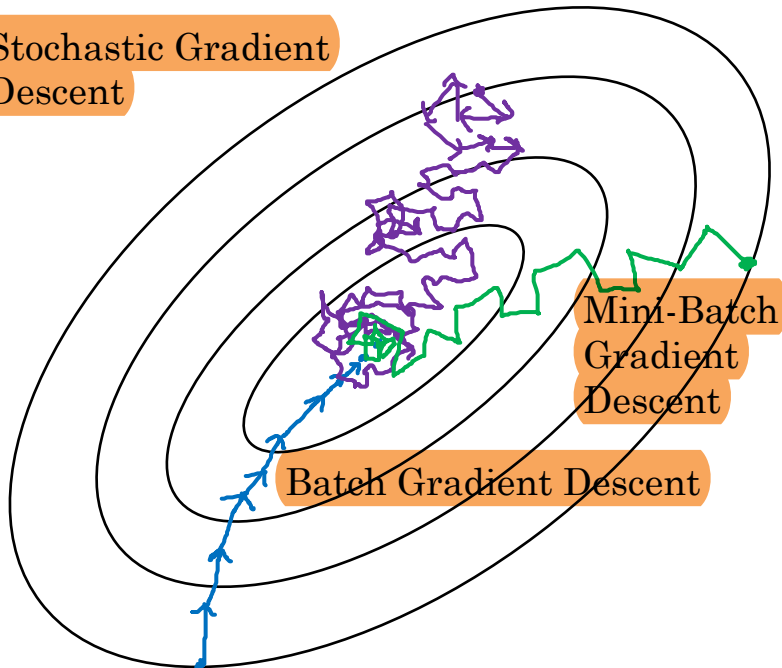
# Choosing your mini-batch size

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

→ 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: Somewh in-between $\underline{1}$ and $\underline{m}$



Stochastic Gradient Descent

Mini-Batch Gradient Descent

Batch Gradient Descent

Stochastic
gradient
Descent
{
Lose speedup
from vectorization

In-between
(mini-batch size
not too big/small)
}
Fastest learning.
• Vectorization.
  (~1000)
• Make progress without
  processing entire training 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.
$(m \le 2000)$

Typical mini-batch sizes :

$\longrightarrow$ $\underbrace{64 , 128, 256, 512}$   $\frac{1024}{2^{10}}$
      $2^6$    $2^7$   $2^8$   $2^9$

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