

Variation of Gradient Descent Optimization

Batch Gradient Descent vs stochastic gradient descent vs Mini Batch GD

Batch GD

Stochastic GD

Mini Batch GD

→ Gradient is computed over entire dataset before making an update to model parameters

→ gradient is computed over 1 training example before making update

→ subset of training data is created before making update

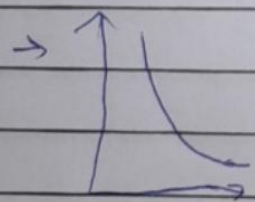
100

100

10 10 10

$\theta = \theta - \eta \text{grad}$
 \downarrow
 1 epoch grad
 1 update $\theta - \eta \text{grad}^2$
 100 update $\theta - \eta \text{grad}^2$

1 epoch = 10 mini batch
 $\theta = \theta - \eta g^1$
 $\theta = \theta - \eta g^2$



loss dec smoothly

→ faster convergence but noisy updates
 → unable to use vectorization

more frequent updates
 faster convergence

Batch Gradient Descent

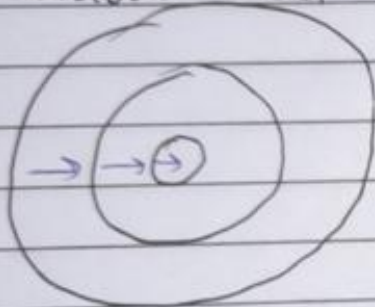
Pseudocode

1 epoch $|x, y|$

1 epoch iterate over entire data

for e in range (MAX-ITRS):
 $grad = gradient(x, y)$
 $\theta = \theta - \eta \cdot grad$

Contour Plot



Loss Function



Mini Batch Gradient Descent

Pseudocode

1 epoch $|b_1| |b_2| |b_3| \dots |$
 x', y'

batch no = m

batch size

for e in range (max-its):
 shuffle(data)

for batch in all-batches:

$x', y' = load_batch()$

$grad = calc_grad(x, y)$

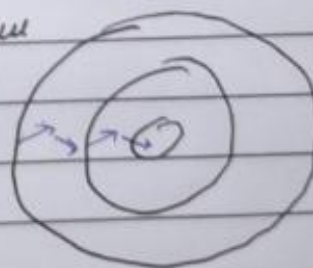
$\theta = \theta - \eta \cdot grad$

freq
of
updates
↓

faster more computation
convergence

→ updates can be noisy
 → achieve local minima
 much earlier

Contour plot



loss function



noisy loss
curve