



The most common optimization algorithms

1.Gradient Descent

□

compute the gradient accuracy how much data we use to
time

1.1BGD

□

taset the entire training d
a local minimum for non-convex surfaces
very slow do not fit in memory update our
model online

1.2SGD

□

$(x(i),y(i))$
new and potentially better local minima
faster can be used to learn online.

1.3 Mini-batch gradient descent

□

Challenges

2. Momentum(动量)

helps accelerate SGD

□

□

3. NAG(牛顿动量)

□

4. Adagrad

main weakness

5. Adadelta

w of accumulated past gradients to some fixed size w.

restricts the window

RMSprop

Adam

□ □

How to Choose

input data is sparse
learning-rate methods

the adaptive

RMSprop

Adagrad

its radically diminishing l

arning rates

Sprop

Adam

adds bias-correction and momentum to R

Adam

ll choice.

Insofar, Adam might be the best over

SGD usually achieves to find a minimum

much more reliant on

robust initialization and annealing schedule

ork

fast convergence and train a deep or complex neural net
adaptive learning rate methods