

Gradient descent

$$L(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N L^{(i)}(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N L(y^{(i)}, \vec{\theta}, t^{(i)})$$

↑
parameters of the model (include biases and weights)

loss computed over training example $(x^{(i)}, t^{(i)})$

model's prediction on input $x^{(i)}$ given that its parameters are set to θ .

$$\frac{\partial L(\theta)}{\partial \theta} = \frac{1}{N} \sum_{i=1}^N \frac{\partial L^{(i)}(\theta)}{\partial \theta}$$

by linearity

requires summing over all of the training examples.

⇒ batch training

Batch training is impractical when there are many training examples (when N is large).

Stochastic Gradient Descent

$$\vec{\theta} \leftarrow \vec{\theta} - \alpha \frac{\partial L^{(i)}}{\partial \theta}$$

loss for a single training example

→ SGD can make progress even before it has looked at all of the training data.

→ If training example is picked at random, Stochastic GD is an unbiased estimate of batch GD:

$$E\left[\frac{\partial L^{(i)}}{\partial \theta}\right] = \frac{1}{N} \sum_{i=1}^N \frac{\partial L^{(i)}}{\partial \theta} = \frac{\partial L(\theta)}{\partial \theta}$$

⚠ SGD analyses a single training example for each update, so it does not take advantage of vectorization.

Minibatch SGD

→ tradeoff speed (by including more training examples in each step of SGD) with memory usage (by not including all training examples as done in batch GD).

→ At each step, we sample a minibatch of training examples.
A minibatch contains a "small" number of examples.

→ Each pass through the entire training set is called an epoch

→ Minibatch SGD reduces variance in the updates:

$$\text{Var}\left[\frac{1}{S} \sum_{i=1}^S \frac{\partial L^{(i)}}{\partial \theta}\right] = \frac{1}{S^2} \text{Var}\left[\sum_{i=1}^S \frac{\partial L^{(i)}}{\partial \theta}\right]$$

update on a minibatch of S training examples

$$= \frac{1}{S^2} S \text{Var}\left[\frac{\partial L^{(i)}}{\partial \theta}\right]$$

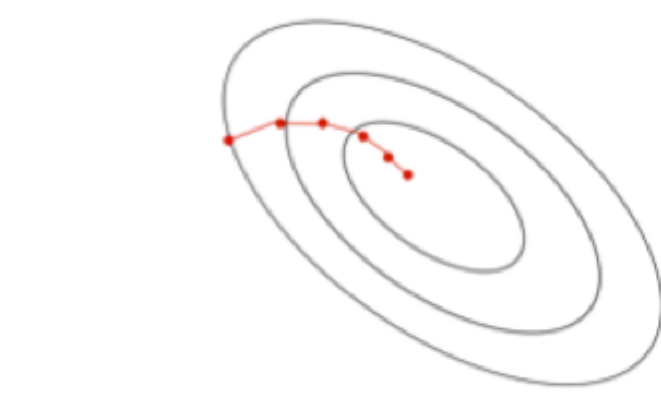
$$= \frac{1}{S} \text{Var}\left[\frac{\partial L^{(i)}}{\partial \theta}\right]$$

→ Size of the minibatch (number of training examples) is a hyperparameter.

↳ too large: minibatch SGD update will involve too much memory

↳ too small: won't benefit from vectorization.

Batch GD



moves directly downhill

Minibatch SGD



Takes steps in a noisy direction, but also moves downhill on average

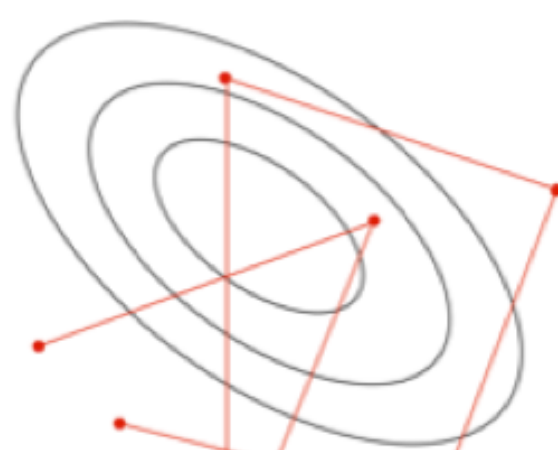
Influence of the learning rate:

small learning rate



little fluctuations

large learning rate



more fluctuations

Typical strategy:

- start with a larger learning rate
- after a few epochs, begin decreasing the learning rate.

⚠ Reducing the learning rate can reduce the loss suddenly (because there are less fluctuations in the SGD)

