

Basic Optimization Algorithms

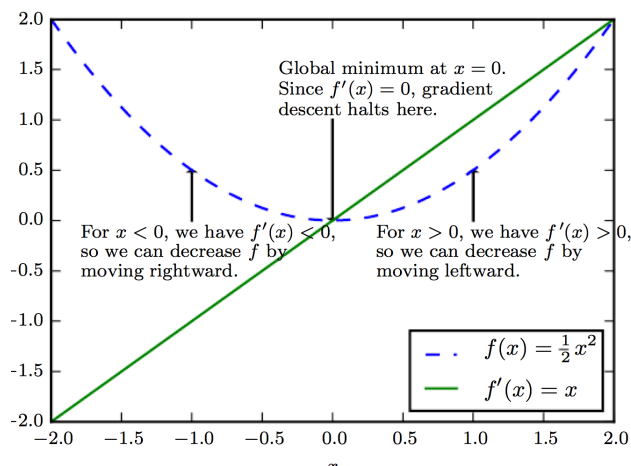
Sargur N. Srihari
srihari@cedar.buffalo.edu

Topics

- Importance of Optimization in machine learning
- How learning differs from optimization
- Challenges in neural network optimization
- **Basic Optimization Algorithms**
 - SGD, Momentum, Nesterov Momentum
- Parameter initialization strategies
- Algorithms with adaptive learning rates
 - AdaGrad, RMSProp, Adam
- Approximate second-order methods
- Optimization strategies and meta-algorithms

1. Stochastic Gradient Descent

- Gradient descent follows gradient of entire training set downhill



Criterion $f(\mathbf{x})$ minimized by moving from current solution in direction of the negative of gradient

- SGD: Accelerated by minibatches downhill
 - Wide use for ML in general and for deep learning
 - Average gradient on a minibatch is an estimate of the gradient

SGD follows gradient estimate downhill

Algorithm: SGD update at training iteration k

Require: Learning rate ϵ_k .

Require: Initial parameter θ

while stopping criterion not met **do**

Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

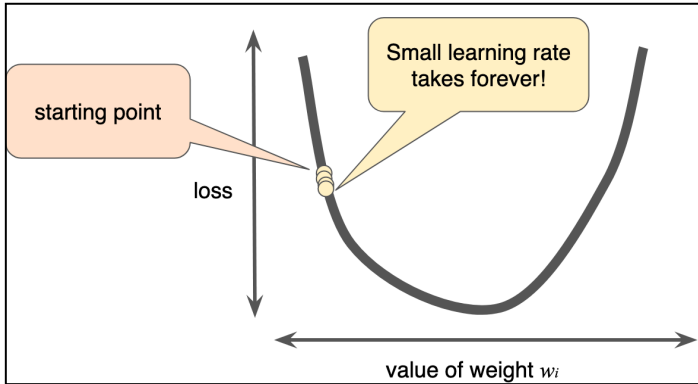
Apply update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$

end while

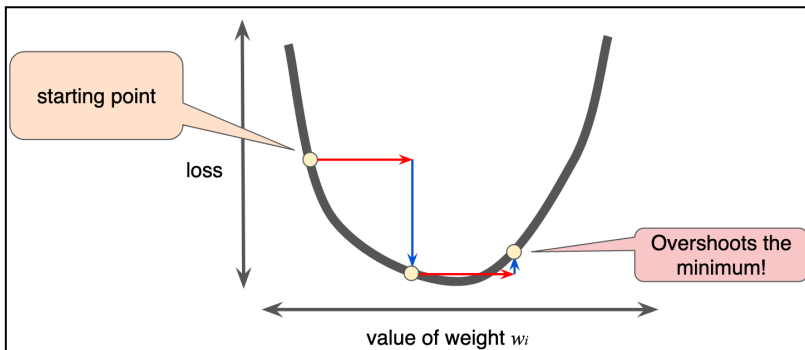
A crucial parameter is the learning rate ϵ

At iteration k it is ϵ_k

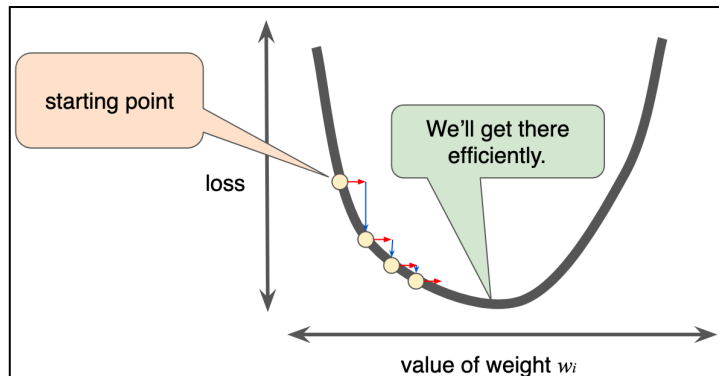
Choice of learning rate



Too small learning rate will take too long



Too large, the next point will perpetually bounce haphazardly across the bottom of the well



If gradient is small then you can safely try a larger learning rate, which compensates for the small gradient and results in a larger step size

Learning rate in Keras

- Keras provides SGD class to implement SGD optimizer with learning rate and momentum
 - The default learning rate is 0.01 and no momentum used

```
1 from keras.optimizers import SGD
2 ...
3 opt = SGD()
4 model.compile(..., optimizer=opt)
```

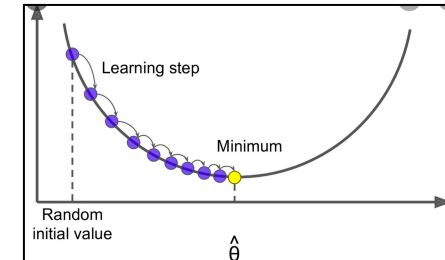
Need for decreasing learning rate

- True gradient of total cost function
 - Becomes small and then 0
 - We can use a fixed learning rate
- But SGD has a source of noise
 - Random sampling of m training samples
 - Gradient does not vanish even when arrive at a minimum
 - Sufficient condition for SGD convergence
$$\sum_{k=1}^{\infty} \epsilon_k = \infty \quad \sum_{k=1}^{\infty} \epsilon_k^2 < \infty$$
 - Common to decay learning rate linearly until iteration τ : $\epsilon_k = (1 - \alpha)\epsilon_0 + \alpha\epsilon_\tau$ with $\alpha = k/\tau$
 - After iteration τ , it is common to leave ϵ constant
 - Often a small positive value in the range 0.0 to 1.0

Learning Rate Decay

- Decay learning rate

$$\tau: \varepsilon_k = (1-\alpha)\varepsilon_0 + \alpha\varepsilon_\tau \text{ with } \alpha = k/\tau$$



- Learning rate is calculated at each update
 - (e.g. end of each mini-batch) as follows:

```
1 lr = initial_lr * (1 / (1 + decay * iteration))
```

- Where *lr* is learning rate for current epoch
- initial_lr* is specified as an argument to SGD
- decay* is the decay rate which is greater than zero and
- iteration* is the current update number

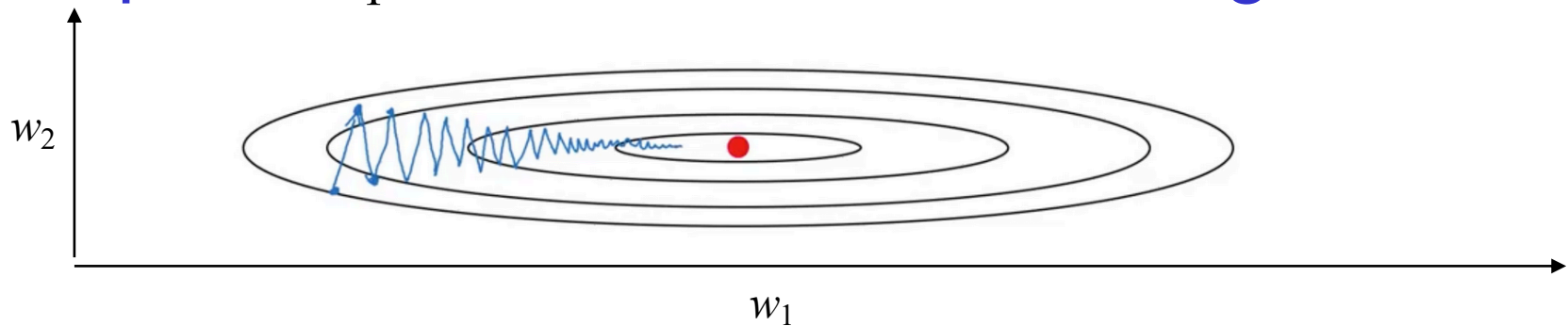
```
1 from keras.optimizers import SGD
2 ...
3 opt = SGD(lr=0.01, momentum=0.9, decay=0.01)
4 model.compile(..., optimizer=opt)
```


2. The Momentum method

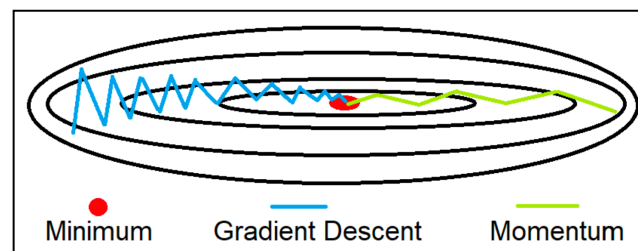
- SGD is a popular optimization strategy
- But it can be slow
- Momentum method accelerates learning, when:
 - Facing high curvature
 - Small but consistent gradients
 - Noisy gradients
- Algorithm accumulates moving average of past gradients and move in that direction, while exponentially decaying

Gradient Descent with momentum

- Gradient descent with momentum converges faster than standard gradient descent
- Taking large steps in w_2 direction and small steps in w_1 direction slows down algorithm



- Momentum reduces oscillation in w_2 direction



- We can set a higher learning rate

Momentum definition

- Introduce variable \mathbf{v} , or velocity
- It is the direction and speed at which parameters move through parameter space
- Momentum in physics is mass times velocity
- The momentum algorithm assumes unit mass
- A hyperparameter $\alpha \in [0,1)$ determines exponential decay

Momentum update rule

- The update rule is given by

$$\begin{aligned} \mathbf{v} &\leftarrow \alpha \mathbf{v} - \varepsilon \nabla_{\theta} \left(\frac{1}{m} \sum_{i=1}^m L(f(x^{(i)}; \theta), y^{(i)}) \right) \\ \theta &\leftarrow \theta + \mathbf{v} \end{aligned}$$

- The velocity \mathbf{v} accumulates the gradient elements $\nabla_{\theta} \left(\frac{1}{m} \sum_{i=1}^m L(f(x^{(i)}; \theta), y^{(i)}) \right)$
- The larger α is relative to ε , the more previous gradients affect the current direction
- The SGD algorithm with momentum is next

SGD algorithm with momentum

Algorithm: SGD with momentum

Require: Learning rate ϵ , momentum parameter α .

Require: Initial parameter θ , initial velocity v .

while stopping criterion not met **do**

Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

Compute gradient estimate: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

Compute velocity update: $\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \mathbf{g}$

Apply update: $\theta \leftarrow \theta + \mathbf{v}$

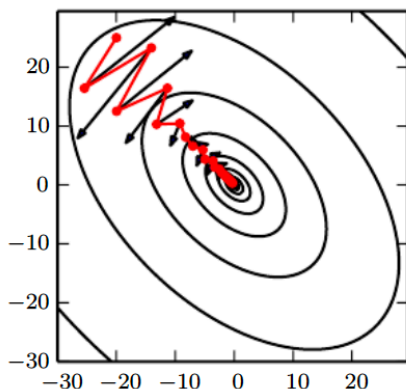
end while

Keras: The learning rate can be specified via the *lr* argument and the momentum can be specified via the *momentum* argument.

```
1 from keras.optimizers import SGD
2 ...
3 opt = SGD(lr=0.01, momentum=0.9)
4 model.compile(..., optimizer=opt)
```

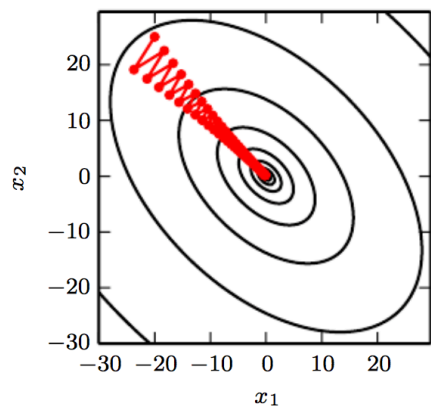
Momentum

- SGD with momentum



Contour lines depict a quadratic loss function
With a poorly conditioned Hessian matrix
Red path cutting across the contours depicts path followed by momentum learning rule as it minimizes this function

- Comparison to SGD without momentum



At each step we show path that would be taken by SGD at that step
Poorly conditioned quadratic objective
Looks like a long narrow valley with steep sides
Wastes time

3. Nesterov Momentum

- A variant to accelerate gradient, with update

$$\begin{aligned} \mathbf{v} &\leftarrow \alpha \mathbf{v} - \epsilon \nabla_{\boldsymbol{\theta}} \left[\frac{1}{m} \sum_{i=1}^m L(\mathbf{f}(\mathbf{x}^{(i)}; \boldsymbol{\theta} + \alpha \mathbf{v}), \mathbf{y}^{(i)}) \right], \\ \boldsymbol{\theta} &\leftarrow \boldsymbol{\theta} + \mathbf{v}, \end{aligned}$$

- where parameters α and ϵ play a similar role as in the standard momentum method
- Difference between Nesterov and standard momentum is where gradient is evaluated.
 - Nesterov gradient is evaluated after the current velocity is applied.
 - Thus one can interpret Nesterov as attempting to add a correction factor to the standard method of momentum

SGD with Nesterov Momentum

- A variant of the momentum algorithm
 - Nesterov's accelerated gradient method
- Applies a correction factor to standard method

Algorithm: SGD with Nesterov momentum

Require: Learning rate ϵ , momentum parameter α .

Require: Initial parameter θ , initial velocity v .

while stopping criterion not met **do**

Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding labels $\mathbf{y}^{(i)}$.

Apply interim update: $\tilde{\theta} \leftarrow \theta + \alpha v$

This line is added from plain momentum

Compute gradient (at interim point): $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(\mathbf{x}^{(i)}; \tilde{\theta}), \mathbf{y}^{(i)})$

Compute velocity update: $v \leftarrow \alpha v - \epsilon \mathbf{g}$

Apply update: $\theta \leftarrow \theta + v$

end while

Keras SGD with Nesterov Momentum

- An optimizer is one of the two arguments required for compiling a Keras model:

```
from keras import optimizers

model = Sequential()
model.add(Dense(64, kernel_initializer='uniform', input_shape=(10,)))
model.add(Activation('softmax'))

sgd = optimizers.SGD(lr=0.01, decay=1e-6, momentum=0.9, nesterov=True)
model.compile(loss='mean_squared_error', optimizer=sgd)
```

Arguments

learning_rate: float ≥ 0 . Learning rate.

momentum: float ≥ 0 . Parameter that accelerates SGD in the relevant direction and dampens oscillations.

nesterov: boolean. Whether to apply Nesterov momentum.