

# Optimization in Deep Learning

R.C. Luo

Oct 2020

“Optimization is more of an art than a science.” This note is a summary of the 8th chapter of the textbook *Deep Learning* [1].

## 1 Neural Network Optimization

Neural network optimizations are in general more indirect than pure optimizations. In practice, bunches of methods are introduced regarding how to shuffle, pack and feed the training data in order to make them function more like test data.

Deep learning wants to minimize the **generalization error**

$$J^*(\theta) = E_{(x,y) \sim p_{\text{data}}} [L(f(x; \theta), y)],$$

but can only observe and minimize the **empirical error**

$$J(\theta) = E_{(x,y) \sim \hat{p}_{\text{data}}} [L(f(x; \theta), y)] = \frac{1}{m} \sum_{i=1}^m L(f(x^{(i)}; \theta), y^{(i)}).$$

## 2 Challenges

Some critical challenges in optimizations of neural networks:

- **Ill Conditioning.** First-order optimizations chop off the second-order term in Taylor expansion, but for ill conditioned Hessian matrix, this kind of chopping may be problematic ( $\frac{1}{2}\epsilon^2 g^T H g > \epsilon g^T g$ ).
- **Local Minima.** It remains open questions for neural networks of practical interest whether local minima can be the decent approximation of global minima, whether trainings are easily attracted to local minima of high cost, and whether training dynamics can escape poor local minima.
- **Plateaus, Saddle Points.** For large neural networks and higher dimensional parameter space, it's even more common to encounter saddle points in comparison with local minima. Unmodified Newton method can be easily trapped by saddle points, while gradient descent dynamics seems

to be able to rapidly escape from saddle points. Other flat regions of high cost (plateaus) pose major problems for almost all numerical optimization algorithms.

- **Cliffs and Exploding Gradients.** Established techniques: gradient clipping.
- **Long-Term Dependencies.** Established techniques: LSTM, GRU, Attention.
- **Inexact Gradients.** Most modern neural networks are designed to keep robustness under gradient noises. If one really wants to avoid this problem, choosing a surrogate loss function is suggested.
- **Poor Correspondence between Local and Global Structure.** Gradient descent and essentially all optimizers take small, local moves. However, this sometimes may not define a well path to a solution (due to ill conditionings, flat regions, as well as greediness). The most ideal thing may be to initialize the parameters perfectly onto a path to a solution.
- **Theoretical Limits.** Last but not least, the theoretical proofs are very limited in this field and the bounds are too loose and of few practical uses.

## 3 Algorithms

### 3.1 Gradient Descent

**Gradient descent** and its variants are probably the most used optimization algorithms. Algorithm 1 shows how to follow the gradient downhill.

---

**Algorithm 1** Minibatch SGD update at training iteration  $k$

---

**Require:** Learning rate  $\epsilon_k$ .

**Require:** Initial parameter  $\theta$ .

**while** stopping criterion not met **do**

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

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

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

**end while**

---

The **batch size**  $b$  is a hyper-parameter. When  $b = \# \text{training samples}$ , it is **full-batch gradient descent**. When  $b = 1$ , it is **stochastic gradient descent**. Conventionally,  $b$  is set to some power of 2 between 32 and 256.

Another crucial parameter is the **learning rate**  $\epsilon$ . The simplest way to set  $\epsilon$  to some small constant is workable. It's more common to gradually decay the

learning rate, for example

$$\epsilon_k = (1 - \frac{k}{r})\epsilon_0 + \frac{k}{r}\epsilon_r.$$

### 3.2 Momentum

The method of **momentum** is designed to accelerate learning, especially in the face of high curvature, small but consistent gradients, or noisy gradients. The momentum algorithm accumulates an exponentially decaying moving average of past gradients and continues to move in their direction.

The SGD algorithm with momentum accumulation is shown below.

---

#### Algorithm 2 SGD with momentum

---

**Require:** Learning rate  $\epsilon$ , momentum parameter  $\alpha (= 0.5, 0.9, 0.99)$ .

**Require:** Initial parameter  $\theta$ , initial velocity  $v$ .

**while** stopping criterion not met **do**

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

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

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

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

**end while**

---

**Nesterov momentum** is attempting to add a *correction factor* to the standard method of momentum.

---

#### Algorithm 3 SGD with Nesterov momentum

---

**Require:** Learning rate  $\epsilon$ , momentum parameter  $\alpha (= 0.5, 0.9, 0.99)$ .

**Require:** Initial parameter  $\theta$ , initial velocity  $v$ .

**while** stopping criterion not met **do**

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

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

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

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

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

**end while**

---

### 3.3 Adaptive Gradient

Learning rate is a crucial hyperparameter in neural networks. Some methods automatically adapt the learning rates through the course of learning. This section will give a brief review over related algorithms.

The **AdaGrad** algorithm is designed to decay the learning rate. It has desirable theoretical properties in the context of convex optimizations, but only fairly works for deep learning optimizations.

---

**Algorithm 4** The AdaGrad algorithm

---

**Require:** Global Learning rate  $\epsilon$ .

**Require:** Initial parameter  $\theta$ .

**Require:** Small constant  $\delta (= 10^{-7})$ , for numerical stability.

Initialize gradient accumulation variable  $\mathbf{r} = \mathbf{0}$

**while** stopping criterion not met **do**

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

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

    Accumulate squared gradient:  $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$

    Compute update:  $\Delta \theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$ . (Division and square root applied element-wisely)

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

**end while**

---

The **RMSProp** algorithm modifies Adagrad to perform better in the non-convex setting by changing the gradient accumulation into an exponentially weighted moving average.

---

**Algorithm 5** The RMSProp algorithm

---

**Require:** Global Learning rate  $\epsilon$ , decay rate  $\rho$ .

**Require:** Initial parameter  $\theta$ .

**Require:** Small constant  $\delta (= 10^{-6})$ , for numerical stability.

Initialize gradient accumulation variable  $\mathbf{r} = \mathbf{0}$

**while** stopping criterion not met **do**

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

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

    Accumulate squared gradient:  $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$

    Compute update:  $\Delta \theta \leftarrow -\frac{\epsilon}{\sqrt{\delta + \mathbf{r}}} \odot \mathbf{g}$ . ( $\frac{1}{\sqrt{\delta + \mathbf{r}}}$  applied element-wisely)

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

**end while**

---

The **Adam** algorithm is another adaptive learning rate optimization algorithm with its name derived from “adaptive moments”.

### 3.4 Second-Order Method

The most widely used second-order method is Newton’s method. However, it seems not very appropriate in the context of deep neural network trainings regarding computational cost, saddle points and numerical stability.

---

**Algorithm 6** The Adam algorithm

---

**Require:** Step size  $\epsilon (= 0.001)$ .

**Require:** Exponential decay rates  $\rho_1 (= 0.99)$  and  $\rho_2 (= 0.999)$  in  $[0, 1]$ .

**Require:** Small constant  $\delta (= 10^{-8})$ , for numerical stability.

**Require:** Initial parameter  $\theta$ .

Initialize 1st and 2nd moment variables  $\mathbf{s} = \mathbf{0}$ ,  $\mathbf{r} = \mathbf{0}$

Initialize time step  $t = 0$

**while** stopping criterion not met **do**

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

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

$t \leftarrow t + 1$

    Update biased first moment estimate:  $\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$

    Update biased second moment estimate:  $\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$

    Correct bias in first moment:  $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$

    Correct bias in second moment:  $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$

    Compute update:  $\Delta \theta \leftarrow -\epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \delta}}$ . (operations applied element-wisely)

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

**end while**

---

**Conjugate gradient** is a method to efficiently avoid the calculation of the inverse Hessian by iteratively descending conjugate directions. The basic idea is to ensure the orthogonality of two sequential search directions  $\mathbf{d}_t$  and  $\mathbf{d}_{t-1}$ :

$$\mathbf{d}_t^T \mathbf{H} \mathbf{d}_{t-1} = 0.$$

At the training iteration  $t$ , the next direction  $\mathbf{d}_t$  takes the form

$$\mathbf{d}_t = \nabla_{\theta} J(\theta) + \beta_t \mathbf{d}_{t-1}.$$

Two popular ways to compute  $\beta_t$  are:

1. **Fletcher-Reeves.**

$$\beta_t = \frac{\nabla_{\theta} J(\theta_t)^T \nabla_{\theta} J(\theta_t)}{\nabla_{\theta} J(\theta_{t-1})^T \nabla_{\theta} J(\theta_{t-1})}$$

2. **Polak-Ribière.**

$$\beta_t = \frac{(\nabla_{\theta} J(\theta_t) - \nabla_{\theta} J(\theta_{t-1}))^T \nabla_{\theta} J(\theta_t)}{\nabla_{\theta} J(\theta_{t-1})^T \nabla_{\theta} J(\theta_{t-1})}$$

The conjugate gradient algorithm is given in algorithm 8.

---

**Algorithm 7** Newton's method

---

**Require:** Initial parameter  $\theta_0$ **Require:** Training set of  $b$  examples

```
while stopping criterion not met do
  Compute gradient  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$ 
  Compute Hessian:  $\mathbf{H} \leftarrow \frac{1}{m} \nabla_{\theta}^2 \sum_i L(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$ 
  Compute Hessian (pseudo)inverse:  $\mathbf{H}^{-1}$ 
  Compute update:  $\Delta\theta \leftarrow -\mathbf{H}^{-1}\mathbf{g}$ 
  Apply update:  $\theta \leftarrow \theta + \Delta\theta$ 
end while
```

---

---

**Algorithm 8** Conjugate gradient method

---

**Require:** Initial parameters  $\theta_0$ .**Require:** Training set of  $m$  examples.Initialize  $\rho_0 = \mathbf{0}$ Initialize  $g_0 = 0$ Initialize  $t = 1$ 

```
while stopping criterion not met do
  Initialize the gradient  $\mathbf{g}_t = \mathbf{0}$ 
  Compute gradient  $\mathbf{g}_t = \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$ 
  Compute  $\beta_t = \frac{(\mathbf{g}_t - \mathbf{g}_{t-1})^T \mathbf{g}_t}{\mathbf{g}_{t-1}^T \mathbf{g}_{t-1}}$  (Polak-Ribière)
  Compute search direction:  $\rho_t = -\mathbf{g}_t + \beta_t \rho_{t-1}$ 
  Perform line search to find:  $\epsilon^* = \operatorname{argmin}_{\epsilon} \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}; \theta_t + \epsilon \rho_t), y^{(i)})$ 
  Apply update:  $\theta_{t+1} = \theta_t + \epsilon^* \rho_t$ 
   $t \leftarrow t + 1$ 
end while
```

---

## 4 Parameter Initialization

Parameter initialization is of critical importance but not yet well understood. The textbook lists several heuristic methods, most of which are based on experiments and are not theoretically proven.

Some heuristics for initialization of weights are listed as follows:

1. **Break symmetry.** Two children units should have different initial parameters.
2. **Larger/smaller weights.** It's a contradiction. Large weights are favored by optimizations, while smaller weights are favored by regularizations.
3. **Norm-preserving.** Initialize orthogonal weight matrices and set gain factors to scale the activation functions. This may preserve the norms and allow trainings of deeper networks.

4. **Sparse initialization.** Each unit can have exactly  $k$  non-zero weights.
5. **Layerwise initialization.** It may be easier to train a deep network layerwisely with some prior features on each layer.

However, these ideas often fail to lead optimal performance because:

1. We may be using the wrong criteria. These kinds of specific designs may not be beneficial to whole network.
2. During training process, the initial properties may get lost.
3. The criteria might succeed at improving the speed of optimization but inadvertently increase generalization error.

Some heuristics for initialization of bias are listed as follows:

1. For output layer, it's usual to set the bias vector  $\mathbf{b}$  to approximate the marginal class distribution  $\mathbf{c}$ :

$$\text{softmax}(\mathbf{b}) = \mathbf{c}.$$

2. For ReLU units, it's usual to set the bias to some small positive numbers in case of saturation.
3. Sometimes a unit controls whether other units are able to participate in a function. For instance, for gates ( $h \in [0, 1]$ ), it's common to open the gate at first (set  $h = 1$ ).

Besides these simple constant or random methods of initializing model parameters, it is also possible to initialize model parameters by other related, *pretrained* model parameters.

## 5 Meta-Algorithms

## References

- [1] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*, volume 1. MIT Press, 2016.