

# Lecture 12 Training Neural Networks

Shiwei Lan<sup>1</sup>

<sup>1</sup>School of Mathematical and Statistical Sciences  
Arizona State University

STP598 Machine Learning and Deep Learning  
Fall 2021

Training NN

S.Lan

Gradient  
Descent  
Optimization

Training Deep  
Models

Stochastic Gradient  
Descents

Algorithms with  
Adaptive Learning Rates

## 1 Gradient Descent Optimization

## 2 Training Deep Models

Stochastic Gradient Descents

Algorithms with Adaptive Learning Rates

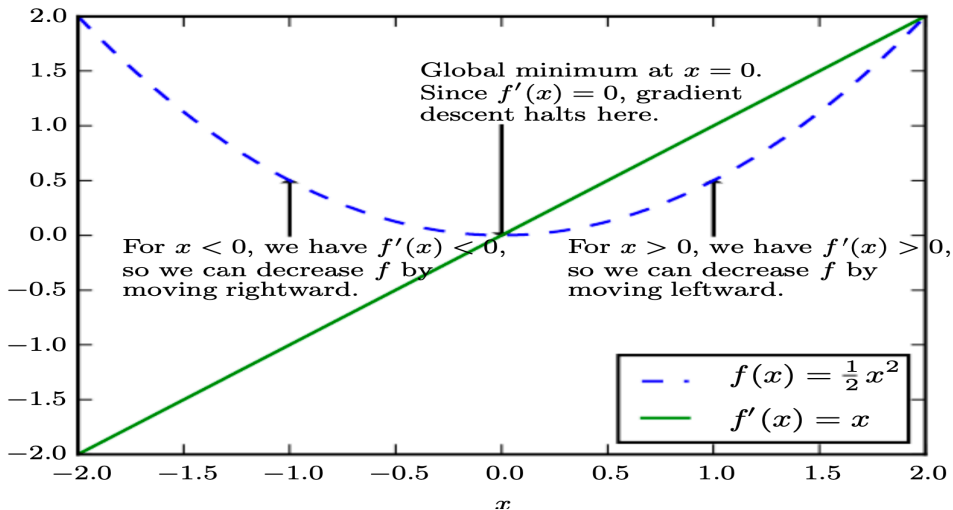
- Most deep learning algorithms involve optimization of some sort.
- Training neural network often relies on minimizing some objective (cost/loss) function  $f(x)$ .
- To minimize  $f$ , we would like to find the direction  $u$  in which  $f$  decreases the fastest by using the directional derivative:

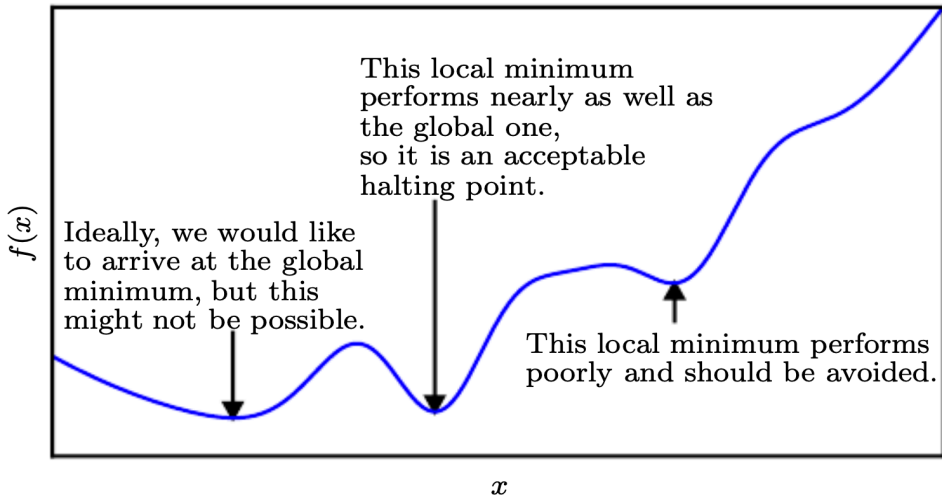
$$\min_{u, \|u\|=1} \left. \frac{\partial}{\partial \alpha} \right|_{\alpha=0} f(x + \alpha u) = \min_{u, \|u\|=1} u^T \nabla_x f(x) = \|\nabla_x f(x)\|_2 \min_{u, \|u\|=1} \cos \theta$$

- The minimal is achieved when  $\theta = \pi$ , i.e. the direction  $u = -\nabla_x f(x)$  is the **steepest descent** or **gradient descent**.
- Then we update the state by

$$x' = x - \varepsilon \nabla_x f(x)$$

where  $\varepsilon$  is called *learning rate*.





There are other challenges like:

- overflow/underflow, e.g. softmax function.
- ill-conditioning:  $f(x) = A^{-1}x$  where  $A \in \mathbb{R}^{n \times n}$  with eigenvalues  $\{\lambda_i\}$ , then condition number is  $\max_{i,j} |\lambda_i/\lambda_j|$ .
- complex landscape, e.g. plateaus, saddle points, cliffs...
- expensive gradients: large data volume.

Training NN

S.Lan

Gradient  
Descent  
Optimization

Training Deep  
Models

Stochastic Gradient  
Descents

Algorithms with  
Adaptive Learning Rates

## 1 Gradient Descent Optimization

## 2 Training Deep Models

- Stochastic Gradient Descents
- Algorithms with Adaptive Learning Rates

- Learning  $\neq$  pure optimization.
- In most machine learning scenarios, we care about some performance measure  $P$ , defined with respect to test set.
- We reduce a different cost function  $J(\theta)$  in the hope that doing so will improve  $P$ . This is in contrast to pure optimization with  $J$  as the goal.
- Typically, the cost function is defined as an expectation of some loss function  $L(\cdot, \cdot)$ , namely, **risk**,

$$J(\theta) = \mathbb{E}_{(x,y) \sim p_{data}} L(f(x; \theta), y)$$

- In reality, we often minimize the an approximate version, **empirical risk**,

$$\tilde{J}(\theta) = \mathbb{E}_{(x,y) \sim \hat{p}_{data}} L(f(x; \theta), y) = \frac{1}{N} \sum_{i=1}^N L(f(x^{(i)}; \theta), y^{(i)})$$

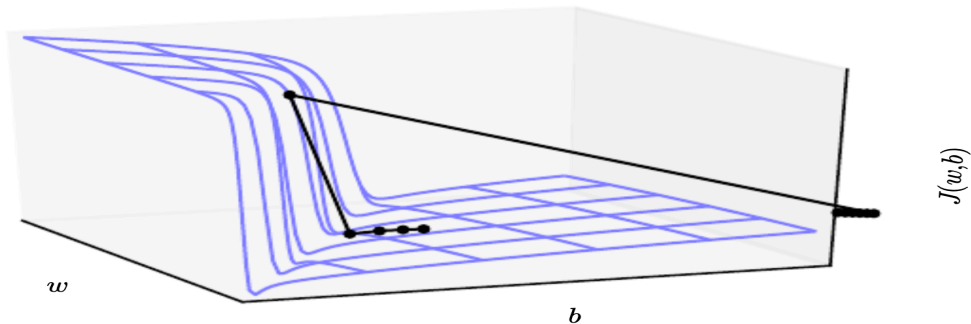


- Empirical risk minimization is prone to overfitting. In stead, we often consider a surrogate loss function, e.g. negative log-likelihood, i.e.

$$\theta_{ML} = \arg \max_{\theta} \sum_{i=1}^N \log p_{model}(x^{(i)}, y^{(i)}; \theta)$$

- To combat the issue of expensive gradients when  $N$  is large, a small batch of data size  $m$  is (randomly) chosen to approximate the gradient in gradient descent algorithms:

$$\theta' = \theta - \frac{N\varepsilon}{m} \nabla_{\theta} \log p_{model}(x^{(i)}, y^{(i)}; \theta)$$



- Like general gradient descent optimization, training neural network also faces the same challenges including cliffs, or exploding gradients.
- To alleviate such issue, **gradient clipping** is adopted when the norm of gradient  $\|g\| > \text{max\_norm}$  for some threshold  $\text{max\_norm}$ :

$$g \leftarrow g \frac{\text{max\_norm}}{\|g\|}$$

- Very deep models involve the composition of several functions or layers.
- In practice, when we update all of the layers simultaneously, unexpected results can happen because many functions composed together are changed simultaneously, e.g.  $\hat{y} = xw_1w_2 \cdots w_l$  where  $h_i = h_{i-1}w_i$ , then the gradient in back-propagation could be either too small or too large.
- To solve this issue, **batch normalization** is adopted.
- Given a minibatch of activations  $\mathbf{H}$ , we normalize  $\mathbf{H}$ , we replace it with

$$\mathbf{H}' = \frac{\mathbf{H} - \mu}{\sigma}$$

where we have

$$\mu = \frac{1}{m} \sum_i \mathbf{H}_i, \quad \sigma = \sqrt{\delta + \frac{1}{m} \sum_i (\mathbf{H} - \mu)_i^2}, \quad \delta \approx 10^{-8}$$

- At test time,  $\mu$  and  $\sigma$  may be replaced by running averages that were collected during training time.

- Stochastic gradient descent (SGD) and its variants are probably the most used optimization algorithms for machine learning in general and for deep learning in particular.
- In practice, it is common to decay the learning rate  $\varepsilon$  linearly in the minibatch gradient descent until iteration  $\tau$ :

$$\varepsilon_k = (1 - \alpha)\varepsilon_0 + \alpha\varepsilon_\tau, \quad \alpha = k/\tau.$$

such that the convergence condition,  $\sum_{k=1}^{\infty} \varepsilon_k = \infty$ ,  $\sum_{k=1}^{\infty} \varepsilon_k^2 < \infty$ , is met.

---

**Algorithm 8.1** Stochastic gradient descent (SGD) update at training iteration  $k$

---

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

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

**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:  $\theta \leftarrow \theta - \varepsilon \hat{\mathbf{g}}$

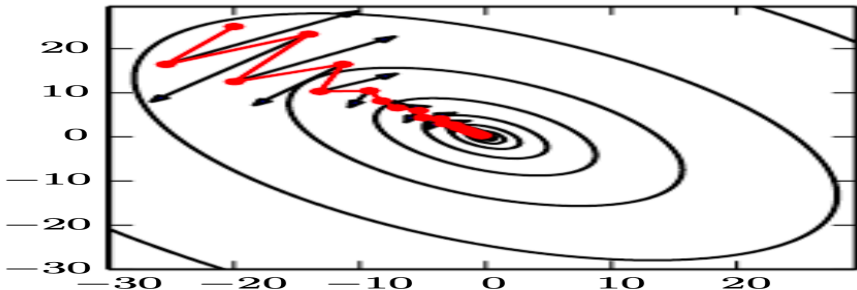
**end while**

---

- While stochastic gradient descent remains a very popular optimization strategy, learning with it can sometimes be slow.
- The method of momentum (Polyak, 1964) is designed to accelerate learning, especially in the face of high curvature, small but consistent gradients, or noisy gradients.

$$v \leftarrow \alpha v - \varepsilon g(\theta), \quad \alpha \in [0, 1)$$

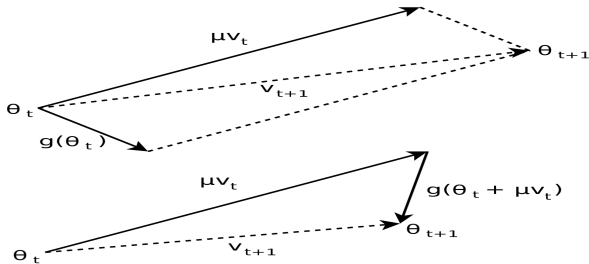
$$\theta \leftarrow \theta + v$$



- Sutskever et al. (2013) introduced a variant of the momentum algorithm that was inspired by Nesterov's accelerated gradient method (Nesterov, 1983, 2004).

$$v \leftarrow \alpha v - \varepsilon g(\theta + \alpha v), \quad \alpha \in [0, 1)$$

$$\theta \leftarrow \theta + v$$



**Figure 1. (Top) Classical Momentum (Bottom) Nesterov Accelerated Gradient**

---

**Algorithm 8.4** The AdaGrad algorithm

---

**Require:** Global learning rate  $\epsilon$

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

**Require:** Small constant  $\delta$ , perhaps  $10^{-7}$ , for numerical stability

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

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

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-wise)

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

**end while**

---

---

**Algorithm 8.5** The RMSProp algorithm

---

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

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

**Require:** Small constant  $\delta$ , usually  $10^{-6}$ , used to stabilize division by small numbers.

Initialize accumulation variables  $\mathbf{r} = 0$

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

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

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

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

**end while**

---



---

## Algorithm 8.7 The Adam algorithm

---

**Require:** Step size  $\epsilon$  (Suggested default: 0.001)

**Require:** Exponential decay rates for moment estimates,  $\rho_1$  and  $\rho_2$  in  $[0, 1)$ .  
(Suggested defaults: 0.9 and 0.999 respectively)

**Require:** Small constant  $\delta$  used for numerical stabilization. (Suggested default:  $10^{-8}$ )

**Require:** Initial parameters  $\theta$

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

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

$t \leftarrow t + 1$

    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:  $\Delta \theta = -\epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \delta}}$  (operations applied element-wise)

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

**end while**

---

Training NN

S.Lan

Gradient  
Descent  
Optimization

Training Deep  
Models

Stochastic Gradient  
Descents

Algorithms with  
Adaptive Learning Rates

Find more at <https://github.com/jettify/pytorch-optimizer>