

# ECE 449/590 – OOP and Machine Learning

## Lecture 18 Stochastic Gradient Descent

Professor Jia Wang  
Department of Electrical and Computer Engineering  
Illinois Institute of Technology

October 26, 2022

Gradient-Based Optimization

Stochastic Gradient Descent

# Reading Assignment

- ▶ This lecture: Deep Learning 4.3, 5.9, 8.1.3, 8.3, 8.4
- ▶ Next lecture: Deep Learning 6

Gradient-Based Optimization

Stochastic Gradient Descent

# Unconstrained Optimization

- ▶ Minimize some function  $f(\mathbf{x})$  by altering  $\mathbf{x}$ .
  - ▶ Use  $-f(\mathbf{x})$  if you would like to maximize  $f(\mathbf{x})$ .
  - ▶  $f(\mathbf{x})$  is the objective function, a.k.a. criterion, cost function, loss function, and error function.
- ▶ The optimal solution  $\mathbf{x}^* = \arg \min f(\mathbf{x})$ .
  - ▶ Does  $\mathbf{x}^*$  exist?
  - ▶ How to find it? What do we know/can we compute about  $f$ ?
  - ▶ What if we cannot find it?

# Critical Points

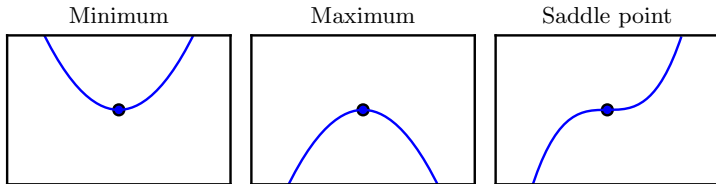


Figure 4.2

(Goodfellow 2017)

- ▶ Consider the 1-D case to minimize  $y = f(x)$ .
- ▶ Critical points  $\frac{dy}{dx} = f'(x) = 0$ .
  - ▶ Not necessary a minimum.

# Gradient Descent (1-D)

- ▶ First order Taylor expansion:  $f(x + \epsilon) \approx f(x) + \epsilon f'(x)$
- ▶ If  $f'(x) < 0$ , increase  $x$ ; if  $f'(x) > 0$ , decrease  $x$ .
  - ▶ For small enough  $\epsilon$ ,  $f(x)$  will decrease.
- ▶ Iterative optimization: apply multiple steps
  - ▶ In each step, update  $x$  to  $x + \epsilon$  by a small enough  $\epsilon$ .
  - ▶ Eventually  $f(x)$  is minimized.
- ▶ Challenges
  - ▶ What if  $f'(x) = 0$  (or too small)?
  - ▶ How large should  $\epsilon$  take?
  - ▶ How to compute  $f'(x)$ ?
  - ▶ What if  $f(x)$  has no minimum?

# Gradient Descent

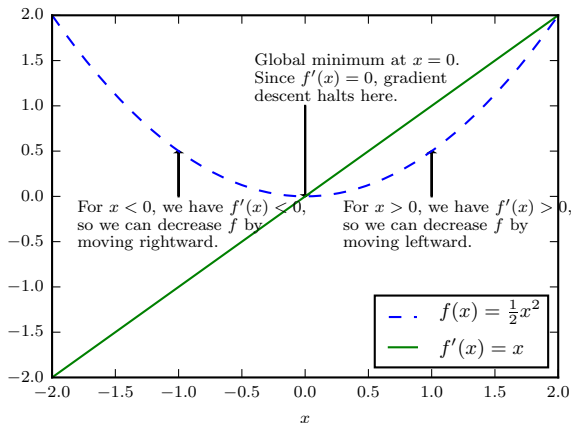


Figure 4.1



# Approximate Optimization

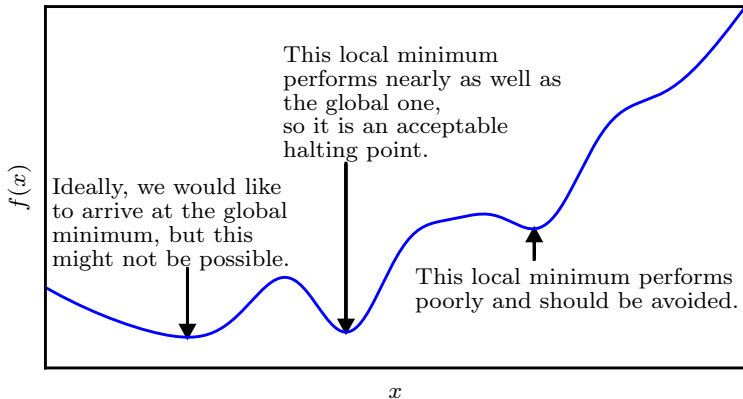


Figure 4.3

# Gradient Descent

- ▶ First order Taylor expansion:  
$$f(\mathbf{x} + \delta\mathbf{x}) \approx f(\mathbf{x}) + \delta\mathbf{x}^\top \nabla_{\mathbf{x}} f(\mathbf{x})$$
  - ▶  $f(\mathbf{x})$  will decrease for small enough  $\delta\mathbf{x}$  with  $\delta\mathbf{x}^\top \nabla_{\mathbf{x}} f(\mathbf{x}) < 0$ .
- ▶ Let  $\delta\mathbf{x} = \alpha\mathbf{u}$  for scalar  $\alpha$  and unit vector  $\mathbf{u}$ .
- ▶ Can we minimize  $\mathbf{u}^\top \nabla_{\mathbf{x}} f(\mathbf{x})$  as a function of  $\mathbf{u}$ ?
  - ▶ For a small enough fixed  $\alpha$ , make  $\delta\mathbf{x}^\top \nabla_{\mathbf{x}} f(\mathbf{x})$  as small as possible in order to make  $f(\mathbf{x} + \delta\mathbf{x})$  as small as possible.
- ▶ Gradient descent:  $\mathbf{u}^* = -\frac{\nabla_{\mathbf{x}} f(\mathbf{x})}{\|\nabla_{\mathbf{x}} f(\mathbf{x})\|}$ 
  - ▶ A.k.a. method of steepest descent.
- ▶ For machine learning, usually we skip to compute the norm for  $\mathbf{u}^*$  and choose  $\delta\mathbf{x} = -\epsilon \nabla_{\mathbf{x}} f(\mathbf{x})$ .
  - ▶  $\epsilon$ : learning rate
  - ▶ Do we need to find the best  $\epsilon$  for each step?

# Gradient Descent and Poor Conditioning

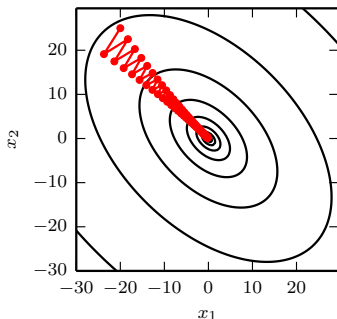


Figure 4.6

(Goodfellow 2017)

- If for each step we could find  $\epsilon^*$  to minimize  $f(\mathbf{x} - \epsilon \nabla_{\mathbf{x}} f(\mathbf{x}))$ , it may still take a lot of steps to reach the minimal  $f$ .

# Beyond Gradient Descent

- ▶ Second order Taylor expansion:

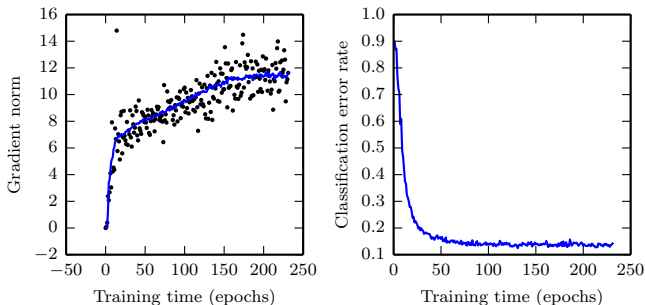
$$f(\mathbf{x} + \delta\mathbf{x}) \approx f(\mathbf{x}) + \delta\mathbf{x}^\top \nabla_{\mathbf{x}} f(\mathbf{x}) + \frac{1}{2} \delta\mathbf{x}^\top \mathbf{H}_{\mathbf{x}} f(\mathbf{x}) \delta\mathbf{x}$$

- ▶ The Hessian matrix:  $\mathbf{H}_{\mathbf{x}} f(\mathbf{x})_{i,j} = \frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x})$
- ▶ Best step size for gradient descent  $\delta\mathbf{x} = -\epsilon^* \nabla_{\mathbf{x}} f(\mathbf{x})$ :

$$\epsilon^* = \frac{\nabla_{\mathbf{x}} f(\mathbf{x})^\top \nabla_{\mathbf{x}} f(\mathbf{x})}{\nabla_{\mathbf{x}} f(\mathbf{x})^\top \mathbf{H}_{\mathbf{x}} f(\mathbf{x}) \nabla_{\mathbf{x}} f(\mathbf{x})}$$

- ▶ Newton's method:  $\delta\mathbf{x}^* = -(\mathbf{H}_{\mathbf{x}} f(\mathbf{x}))^{-1} \nabla_{\mathbf{x}} f(\mathbf{x})$ 
  - ▶ Allow to move in a different direction than gradient descent.
  - ▶ Still, there are issues related with non-minimum critical points and  $\mathbf{H}_{\mathbf{x}} f(\mathbf{x})$  being too close to 0.
- ▶ Practical challenges for machine learning
  - ▶ How to calculate  $\mathbf{H}_{\mathbf{x}} f(\mathbf{x})$ ?
  - ▶ How to store  $\mathbf{H}_{\mathbf{x}} f(\mathbf{x})$ ?

# We usually don't even reach a local minimum



(Goodfellow 2017)

- For many machine learning problems, we don't even need  $f'(x)$  to approach 0.
- For simplicity, just think about to minimize  $f(x) = \log(|x - c|)$

Gradient-Based Optimization

Stochastic Gradient Descent

# Loss Function Revisited

- ▶ For training, loss function is the average of individual ones.

$$J(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, y \sim p_{data}} L(\mathbf{x}, y; \boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m L(\mathbf{x}^{(i)}, y^{(i)}; \boldsymbol{\theta})$$

- ▶ The gradient can be computed as the average too.

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} L(\mathbf{x}^{(i)}, y^{(i)}; \boldsymbol{\theta})$$

- ▶ But this will be too time consuming for gradient descent.
  - ▶ Need to visit every training example for each step.

# Stochastic Gradient Descent (SGD)

- ▶ Instead of computing the gradient  $\nabla_{\theta} J(\theta)$  from all training examples, we may approximate it by sampling.
  - ▶ Indeed,  $\frac{1}{m} \sum_{i=1}^m L(\mathbf{x}^{(i)}, y^{(i)}; \theta)$  is an approximation of the actual loss function  $\mathbb{E}_{\mathbf{x}, y \sim p_{data}} L(\mathbf{x}, y; \theta)$ .
  - ▶ For machine learning, finding the actual minimum is not as important as reducing the loss function.
- ▶ Stochastic Gradient Descent (SGD) and Minibatch
  - ▶ Sample  $m'$  examples  $\mathbb{B} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m')}\}$  with a fixed  $m' \ll m$  no matter how large  $m$  is.
  - ▶  $\nabla_{\theta} J(\theta) \approx \frac{1}{m'} \sum_{i=1}^{m'} \nabla_{\theta} L(\mathbf{x}^{(i)}, y^{(i)}; \theta)$
- ▶ In practice, training of neural network models are organized into epochs.
  - ▶ The epoch begins by randomly shuffling training examples.
  - ▶ Then each step consumes  $m'$  examples.
  - ▶ The epoch ends after  $\frac{m}{m'}$  steps when all examples are consumed once.



# Discussions on Minibatch Algorithms

- ▶ Larger batches provide a more accurate estimate of the gradient, but with less than linear returns.
- ▶ Multicore architectures are usually underutilized by extremely small batches.
- ▶ Typically, all examples need to be available from the memory so that they could be processed in parallel. This may limit batch size in certain hardware.
- ▶ Some kinds of hardware, especially GPUs, achieve better runtime with specific sizes of arrays, e.g. powers of 2.
- ▶ Small batches introducing noises to the training process may benefit learning as a whole because they may work as regularization to reduce generalization errors.

# The Learning Rate

- ▶ It is common for SGD to use different learning rates  $\epsilon_k$  for different steps  $k$ .
- ▶ In theory, for SGD to converge, it is sufficient that

$$\sum_{k=1}^{\infty} \epsilon_k = \infty, \text{ and } \sum_{k=1}^{\infty} \epsilon_k^2 < \infty$$

- ▶ In practice, reduce learning rates as training progresses.
  - ▶ Larger rates help to update weights faster in the beginning.
  - ▶ Smaller rates help to keep what have already been learned.
  - ▶ E.g. linear decay until step  $\tau$ :  $\epsilon_k = (1 - \frac{k}{\tau})\epsilon_0 + \frac{k}{\tau}\epsilon_{\tau}$

# Momentum

- ▶ Noisy gradients due to minibatch or poorly conditioned Hessian may cause gradient descent to follow a zig-zaging path.
  - ▶ Lead to slow convergence as steps are cancelling each other.
- ▶ Momentum: use an exponentially decaying moving average of past gradients to update weights.
  - ▶ Velocity updates:  $\mathbf{v} \leftarrow \alpha \mathbf{v} - \frac{\epsilon}{m'} \sum_{i=1}^{m'} \nabla_{\boldsymbol{\theta}} L(\mathbf{x}^{(i)}, y^{(i)}; \boldsymbol{\theta})$
  - ▶ Weights updates:  $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{v}$
- ▶ The actual step sizes are large to accelerate learning if past gradients are aligned.

# Parameter Initialization Strategies

- ▶ While SGD takes multiple steps to reduce  $J(\theta)$ , if a good  $\theta$  is chosen for the first step, the iterative process may take less steps to converge.
  - ▶ In extreme cases, bad initial  $\theta$ 's may prevent SGD to converge to a reasonable minimum.
- ▶ Symmetry in neural network models
  - ▶ Many nodes in neural network models have identical inputs and drive the same output.
  - ▶ The training process is expected to assign different weights to such nodes so they would learn different features.
  - ▶ However, if their weights are the same in the beginning, gradient descent will compute the same gradient for them and they will remain the same through the learning process.
- ▶ “Break symmetry”: initialize parameters randomly

# Summary

- ▶ Machine learning introduces many unique challenges to optimization.
- ▶ SGD with minibatch works well when training neural network models.
  - ▶ Choices of learning rate, momentum, and initialization may still affect the efficiency of the learning process.