

ECE 595: Machine Learning I

Lecture 5.1: Gradient Descent

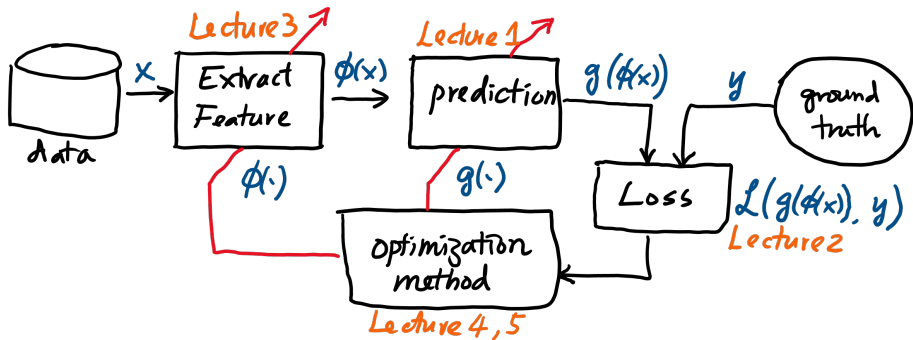
Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



Outline



Outline

Mathematical Background

- Lecture 4: Intro to Optimization
- **Lecture 5: Gradient Descent**

Lecture 5: Gradient Descent

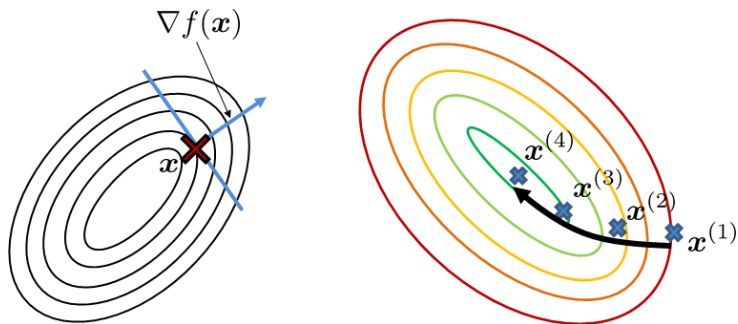
- **Gradient Descent**
 - Descent Direction
 - Step Size
 - Convergence
- Stochastic Gradient Descent
 - Difference between GD and SGD
 - Why does SGD work?

Gradient Descent

The algorithm:

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha^{(t)} \nabla f(\mathbf{x}^{(t)}), \quad t = 0, 1, 2, \dots,$$

where $\alpha^{(t)}$ is called the **step size**.



Why is the direction $-\nabla f(\mathbf{x})$?

- Recall (Lecture 4): If \mathbf{x}^* is optimal, then

$$\underbrace{\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} [f(\mathbf{x}^* + \epsilon \mathbf{d}) - f(\mathbf{x}^*)]}_{\geq 0, \forall \mathbf{d}} = \nabla f(\mathbf{x}^*)^T \mathbf{d}$$
$$\implies \nabla f(\mathbf{x}^*)^T \mathbf{d} \geq 0, \quad \forall \mathbf{d}$$

- But if $\mathbf{x}^{(t)}$ is not optimal, then we want

$$f(\mathbf{x}^{(t)} + \epsilon \mathbf{d}) \leq f(\mathbf{x}^{(t)})$$

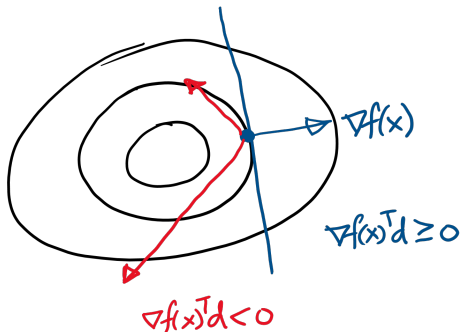
- So,

$$\underbrace{\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} [f(\mathbf{x}^{(t)} + \epsilon \mathbf{d}) - f(\mathbf{x}^{(t)})]}_{\leq 0, \text{ for some } \mathbf{d}} = \nabla f(\mathbf{x}^{(t)})^T \mathbf{d}$$
$$\implies \nabla f(\mathbf{x}^{(t)})^T \mathbf{d} \leq 0$$

Descent Direction

Pictorial illustration:

- $\nabla f(\mathbf{x})$ is **perpendicular** to the contour.
- A search direction \mathbf{d} can either be on the positive side $\nabla f(\mathbf{x})^T \mathbf{d} \geq 0$ or negative side $\nabla f(\mathbf{x})^T \mathbf{d} < 0$.
- Only those on the negative side can reduce the cost.
- All such \mathbf{d} 's are called the **descent directions**.



The Steepest \mathbf{d}

Previous slide: If $\mathbf{x}^{(t)}$ is not optimal yet, then some \mathbf{d} will give

$$\nabla f(\mathbf{x}^{(t)})^T \mathbf{d} \leq 0.$$

- So, let us make $\nabla f(\mathbf{x}^{(t)})^T$ as negative as possible.

$$\mathbf{d}^{(t)} = \underset{\|\mathbf{d}\|_2 = \delta}{\operatorname{argmin}} \nabla f(\mathbf{x}^{(t)})^T \mathbf{d},$$

- We need δ to control the magnitude; Otherwise \mathbf{d} is unbounded.
- The solution is

$$\mathbf{d}^{(t)} = -\nabla f(\mathbf{x}^{(t)})$$

- Why? By Cauchy Schwarz,

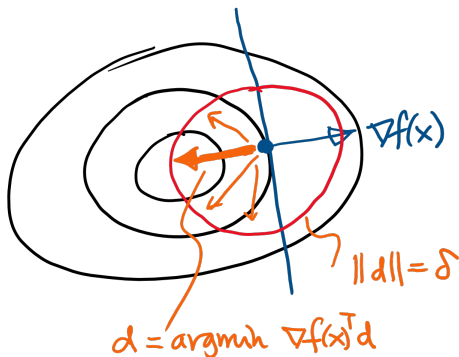
$$\nabla f(\mathbf{x}^{(t)})^T \mathbf{d} \geq -\|\nabla f(\mathbf{x}^{(t)})\|_2 \|\mathbf{d}\|_2.$$

- Minimum attained when $\mathbf{d} = -\nabla f(\mathbf{x}^{(t)})$.
- Set $\delta = \|\nabla f(\mathbf{x}^{(t)})\|_2$.

Steepest Descent Direction

Pictorial illustration:

- Put a ball surrounding the current point.
- All \mathbf{d} 's inside the ball are feasible.
- Pick the one that minimizes $\nabla f(\mathbf{x})^T \mathbf{d}$.
- This direction must be parallel (but opposite sign) to $\nabla f(\mathbf{x})$.



Step Size

The algorithm:

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha^{(t)} \nabla f(\mathbf{x}^{(t)}), \quad t = 0, 1, 2, \dots,$$

where $\alpha^{(t)}$ is called **the step size**.

- **1. Fixed step size**

$$\alpha^{(t)} = \alpha.$$

- **2. Exact line search**

$$\alpha^{(t)} = \underset{\alpha}{\operatorname{argmin}} f\left(\mathbf{x}^{(t)} + \alpha \mathbf{d}^{(t)}\right),$$

- E.g., if $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} + \mathbf{c}^T \mathbf{x}$, then

$$\alpha^{(t)} = -\frac{\nabla f(\mathbf{x}^{(t)})^T \mathbf{d}^{(t)}}{\mathbf{d}^{(t)T} \mathbf{H} \mathbf{d}^{(t)}}.$$

- **3. Inexact line search:**

Amijo / Wolfe conditions. See Nocedal-Wright Chapter 3.1.

Convergence

Let \mathbf{x}^* be the global minimizer. Assume the followings:

- Assume f is **twice differentiable** so that $\nabla^2 f$ exist.
- Assume $0 \preceq \lambda_{\min} \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq \lambda_{\max} \mathbf{I}$ for all $\mathbf{x} \in \mathbb{R}^n$
- Run gradient descent with **exact line search**.

Then, (Nocedal-Wright Chapter 3, Theorem 3.3)

$$\begin{aligned} f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^*) &\leq \left(1 - \frac{\lambda_{\min}}{\lambda_{\max}}\right)^2 \left(f(\mathbf{x}^{(t)}) - f(\mathbf{x}^*)\right) \\ &\leq \left(1 - \frac{\lambda_{\min}}{\lambda_{\max}}\right)^4 \left(f(\mathbf{x}^{(t-1)}) - f(\mathbf{x}^*)\right) \\ &\leq \vdots \\ &\leq \left(1 - \frac{\lambda_{\min}}{\lambda_{\max}}\right)^{2t} \left(f(\mathbf{x}^{(1)}) - f(\mathbf{x}^*)\right). \end{aligned}$$

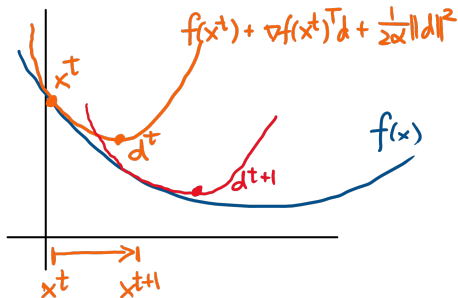
Thus, $f(\mathbf{x}^{(t)}) \rightarrow f(\mathbf{x}^*)$ as $t \rightarrow \infty$.

Understanding Convergence

- Gradient descent can be viewed as **successive approximation**.
- Approximate the function as

$$f(\mathbf{x}^t + \mathbf{d}) \approx f(\mathbf{x}^t) + \nabla f(\mathbf{x}^t)^T \mathbf{d} + \frac{1}{2\alpha} \|\mathbf{d}\|^2.$$

- We can show that the \mathbf{d} that minimizes $f(\mathbf{x}^t + \mathbf{d})$ is $\mathbf{d} = -\alpha \nabla f(\mathbf{x}^t)$.
- This suggests: Use a **quadratic function** to locally approximate f .
- Converge when curvature α of the approximation is not too big.



Advice on Gradient Descent

- Gradient descent is useful because
 - Simple to implement (compared to ADMM, FISTA, etc)
 - Low computational cost per iteration (no matrix inversion)
 - Requires only first order derivative (no Hessian)
 - Gradient is available in deep networks (via back propagation)
- Most machine learning has built-in (stochastic) gradient descents
- Welcome to implement your own, but you need to be careful
 - Convex non-differentiable problems, e.g., ℓ_1 -norm
 - Non-convex problem, e.g., ReLU in deep network
 - Trap by local minima
 - Inappropriate step size, a.k.a. learning rate
- Consider more “transparent” algorithms such as CVX when
 - Formulating problems. No need to worry about algorithm.
 - Trying to obtain insights.