Gradient Descent
00000

Convergence of Gradient Descent
00000000000

Variations
0000000

# CPSC 8420 Advanced Machine Learning
# Week 6: Gradient Descent

Dr. Kai Liu

September 24, 2020

Gradient Descent
○○○○○

Convergence of Gradient Descent
○○○○○○○○○○○

Variations
○○○○○○○

## Motivation

Gradient Descent is one of the most widely used methods in machine learning to optimize solution.

- Most problems in machine learning they don't have a closed solution like Least Squares or PCA, thus an iterative updating algorithm may be more applicable.

- Obtaining a closed solution maybe computationally demanding with operations such as matrix inversion or SVD, *etc*.

- Gradient Descent works well not only for convex problem, but also non-convex one. Though global minimum can not be guaranteed, most of times, local minimum is acceptable.

Gradient Descent
00000

Convergence of Gradient Descent
00000000000

Variations
0000000

## Learning Outcomes

Our goal for today's lecture is to understand:

- Why Gradient Descent decreases the objective

- Convergence rate of GD in various conditions

- Variations of GD

# Gradient Descent

## Differentiable unconstrained minimization

$$\min_x f(x), \quad s.t. \ x \in \mathbb{R}^n \tag{1}$$

where f (objective or cost function) is differentiable

## Iterative descent algorithms

Start with a point $x^0$, and construct a sequence $\{x^t\}$ such that:

$$f(x^{t+1}) < f(x^t), \quad t = 0, 1, 2 \ldots \tag{2}$$
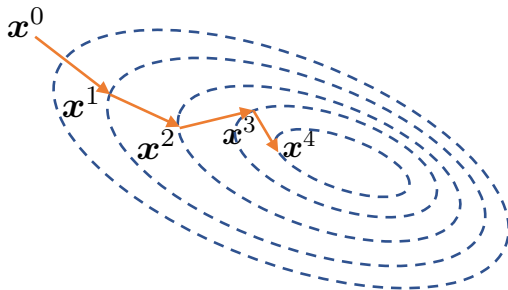
In each iteration, search in descent direction

$$x^{t+1} = x^t + \eta_t d^t, \tag{3}$$

where $d^t$ is the descent direction at $x^t$; $\eta_t > 0$ is the stepsize.

## Gradient descent (GD)

One of the most important examples is **Gradient descent (GD)**:

$$x^{t+1} = x^t - \eta_t \nabla f(x^t). \tag{4}$$



It can be traced to Augustin Louis Cauchy.

## Gradient descent (GD)

One of the most important examples is **Gradient descent (GD)**:

$$x^{t+1} = x^t - \eta_t \nabla f(x^t). \tag{5}$$

descent direction: $d^t = -\nabla f(x^t)$, which can be justified from Taylor Expansion:

$$f(x^{t+1}) \approx f(x^t) + \langle x^{t+1} - x^t, \nabla f(x^t) \rangle + \frac{1}{2} \langle x^{t+1} - x^t, \nabla^2 f(x^t)(x^{t+1} - x^t) \rangle. \tag{6}$$

Gradient Descent
○○○○○

Convergence of Gradient Descent
●○○○○○○○○○○

Variations
○○○○○○○

# Convergence of Gradient Descent

Gradient Descent
00000

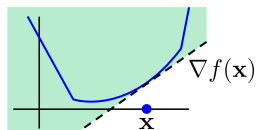Convergence of Gradient Descent
0●00000000

Variations
0000000

## Strongly Convex Strongly Smooth Function

We say a continuously differentiable function $f : \mathbb{R}^p \to \mathbb{R}$ is $\alpha$-strongly convex (SC) and $\beta$-strongly smooth (SS) if for every $x, x^+ \in \mathbb{R}^p$, we have:
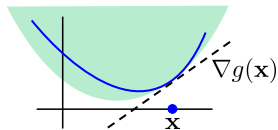
$$\frac{\alpha}{2}\|x^+ - x\|_2^2 \le f(x^+) - f(x) - \langle \nabla f(x), x^+ - x \rangle \le \frac{\beta}{2}\|x^+ - x\|_2^2, \ (7)$$

$\alpha, \beta$ can be taken as the smallest and largest singular value of $\nabla^2 f(x)$.
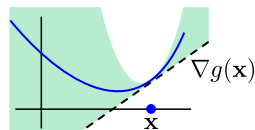
Gradient Descent
00000

Convergence of Gradient Descent
00●00000000

Variations
0000000

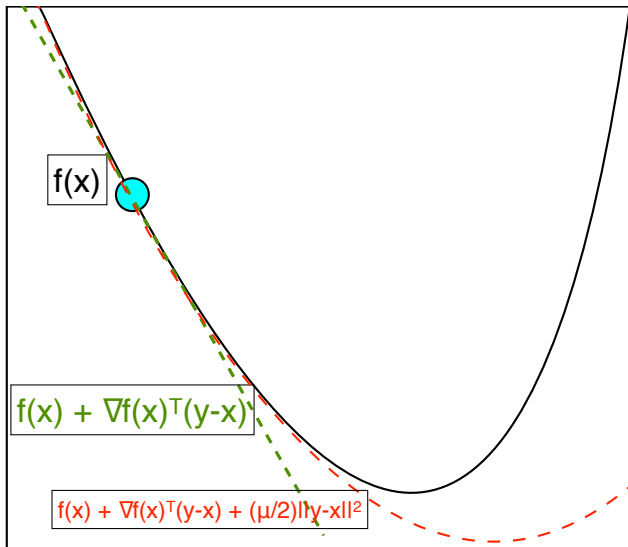# Convex, SC and SS



$-f : \mathbb{R}^d \to \mathbb{R}$
**CONVEX FUNCTION**

$-g : \mathbb{R}^d \to \mathbb{R}$
**STRONGLY CONVEX FUNCTION**

$-g : \mathbb{R}^d \to \mathbb{R}$
**STRONGLY SMOOTH CONVEX FUNCTION**

Gradient Descent
○○○○○

Convergence of Gradient Descent
○○○●○○○○○○○

Variations
○○○○○○○

# Convex, SC and SS



$f(x)$

$f(x) + \nabla f(x)^\top(y-x)$

$f(x) + \nabla f(x)^\top(y-x) + (\mu/2)\|y-x\|^2$

Gradient Descent
ooooo

Convergence of Gradient Descent
oooooooooooo

Variations
ooooooo

## Strongly Convex Strongly Smooth Function

$$\frac{\alpha}{2}\|x^+ - x\|_2^2 \le f(x^+) - f(x) - \langle \nabla f(x), x^+ - x \rangle \le \frac{\beta}{2}\|x^+ - x\|_2^2, \ (8)$$
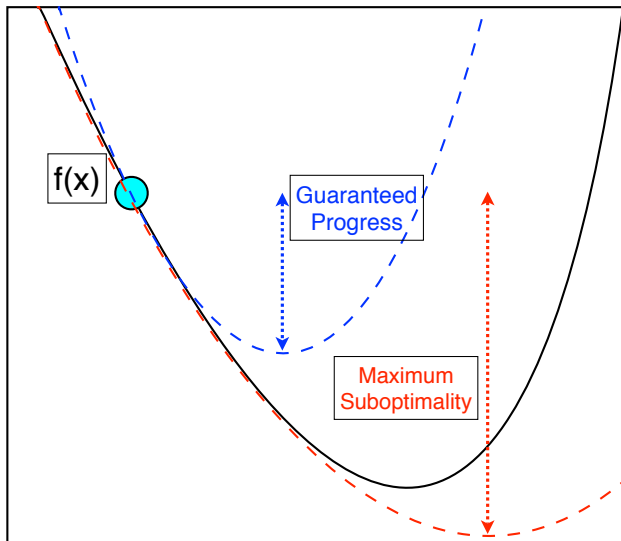
based on which we will have:

$$f(x^+) \le f(x) - \frac{1}{2\beta}\|\nabla f(x)\|^2$$
$$f(x^+) \ge f(x) - \frac{1}{2\alpha}\|\nabla f(x)\|^2. \tag{9}$$

Replacing $x^+$ with $x^*$, then we will have:

$$\frac{1}{2\beta}\|\nabla f(x)\|^2 \le f(x) - f(x^*) \le \frac{1}{2\alpha}\|\nabla f(x)\|^2 \tag{10}$$

# Progress within each Update

Gradient Descent
○○○○○

Convergence of Gradient Descent
○○○○○○●○○○○

Variations
○○○○○○○

## Linear Convergence Rate

$$\begin{aligned}
f(x^+) - f(x^*) &\leq f(x) - f(x^*) - \frac{1}{2\beta}\|\nabla f(x)\|^2 \\
&\leq f(x) - f(x^*) - \frac{\alpha}{\beta}(f(x) - f(x^*)) \\
&= (1 - \frac{\alpha}{\beta})(f(x) - f(x^*)) \\
&\leq exp^{-\frac{\alpha}{\beta}}(f(x) - f(x^*))
\end{aligned} \tag{11}$$

which implies $\frac{f(x^+)-f(x^*)}{f(x)-f(x^*)} = 1 - \frac{\alpha}{\beta}$, is the definition of linear convergence. $\frac{f(x^t)-f(x^*)}{f(x^0)-f(x^*)} \leq exp^{-\frac{\alpha}{\beta}t}$ Then to obtain $\epsilon$-suboptimal result, we need $\mathcal{O}(\kappa log\frac{1}{\epsilon})$ iterations, where we define $\kappa = \frac{\beta}{\alpha}$ to be the **condition number**.

# Revisiting Ridge Regression

Vanilla Least Squares:
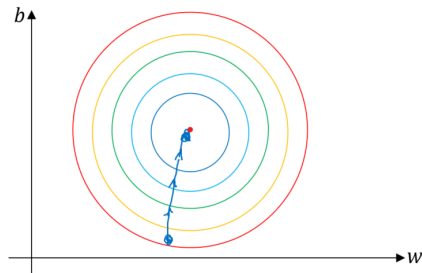
$$\min_\beta \frac{1}{2}\|y - A\beta\|^2 \tag{12}$$

which is a convex problem with respect to $\beta$, and we can verify that
$\nabla^2 f(\beta) = A^T A$, thus $\kappa = \frac{\sigma_{max}(A^T A)}{\sigma_{min}(A^T A)}$.
Vanilla Least Squares:

$$\min_\beta \frac{1}{2}\|y - A\beta\|^2 + \frac{\lambda}{2}\|\beta\|^2 \tag{13}$$

which is also convex with respect to $\beta$, and we can verify that
$\nabla^2 f(\beta) = A^T A + \lambda I$, thus the new $\tilde{\kappa} = \frac{\sigma_{max}(A^T A)+\lambda}{\sigma_{min}(A^T A)+\lambda} < \kappa$.

Gradient Descent
○○○○○

Convergence of Gradient Descent
○○○○○○○○○●○○

Variations
○○○○○○○

# Geometry of $\kappa$ to influence convergence rate

Gradient Descent
00000

Convergence of Gradient Descent
0000000000●0

Variations
0000000

## Non Strongly Convex

$$
\begin{aligned}
f(x^+) - f(x^*) &\leq f(x) - f(x^*) - \frac{1}{2\beta}\|\nabla f(x)\|^2 \\
&\leq \langle \nabla f(x), x - x^* \rangle - \frac{1}{2\beta}\|\nabla f(x)\|^2
\end{aligned}
\tag{14}
$$

on the other hand we have:

$$
\begin{aligned}
\|x^+ - x^*\|^2 &= \|x - \eta\nabla f(x) - x^*\|^2 \\
&= \|x - x^*\|^2 - 2\eta\langle \nabla f(x), x - x^* \rangle + \eta^2\|\nabla f(x)\|^2 \\
&= \|x - x^*\|^2 - 2\eta(\langle \nabla f(x), x - x^* \rangle - \frac{\eta}{2}\|\nabla f(x)\|^2),
\end{aligned}
\tag{15}
$$

then we have
$\langle \nabla f(x), x - x^* \rangle - \frac{\eta}{2}\|\nabla f(x)\|^2 = \frac{1}{2\eta}(\|x - x^*\|^2 - \|x^+ - x^*\|^2)$, and
therefore $f(x^+) - f(x^*) \leq \frac{1}{2\eta}(\|x - x^*\|^2 - \|x^+ - x^*\|^2)$

Gradient Descent
00000

Convergence of Gradient Descent
000000000●

Variations
0000000

## Non Strongly Convex

Summation the equation above from $k = 0$ to $k = T - 1$, we have:
$\sum_{k=0}^{T-1} f(x_{k+1}) - f(x^*) \leq \frac{\|x_0 - x^*\|^2 - \|x_T - x^*\|^2}{2\eta} \leq \frac{\|x_0 - x^*\|^2}{2\eta}$, then
$f(x_T) - f(x^*) \leq \frac{1}{T} \sum_{k=0}^{T-1} f(x_{k+1}) - f(x^*) \leq \frac{\|x_0 - x^*\|^2}{2T\eta}$, which is
sub-linear convergence rate. Then to obtain $\epsilon$-suboptimal result, we
need $T = \mathcal{O}(\frac{1}{\epsilon})$ iterations.

Gradient Descent
00000

Convergence of Gradient Descent
00000000000

Variations
●000000

# Variations

Gradient Descent
○○○○○

Convergence of Gradient Descent
○○○○○○○○○○○

Variations
○●○○○○○

# Heavy-ball



gradient descent



heavy-ball method

Gradient Descent
○○○○○

Convergence of Gradient Descent
○○○○○○○○○○○

Variations
○○●○○○○

## Heavy-ball



B. Polyak

$$\text{minimize}_{\boldsymbol{x} \in \mathbb{R}^n} \quad f(\boldsymbol{x})$$

$$\boldsymbol{x}^{t+1} = \boldsymbol{x}^t - \eta_t \nabla f(\boldsymbol{x}^t) + \underbrace{\theta_t(\boldsymbol{x}^t - \boldsymbol{x}^{t-1})}_{\text{momentum term}}$$

- add inertia to the "ball" (i.e. include a momentum term) to mitigate zigzagging

Gradient Descent
○○○○○

Convergence of Gradient Descent
○○○○○○○○○○○

Variations
○○○●○○○

# Heavy-ball

iteration complexity: $O(\sqrt{\kappa} \log \frac{1}{\varepsilon})$

significant improvement over GD: $O(\sqrt{\kappa} \log \frac{1}{\varepsilon})$ vs. $O(\kappa \log \frac{1}{\varepsilon})$

Gradient Descent
○○○○○

Convergence of Gradient Descent
○○○○○○○○○○○

Variations
○○○○●○○

# Nesterov Accelerate Gradient



— *Nesterov '83*

$$\boldsymbol{x}^{t+1} = \boldsymbol{y}^t - \eta_t \nabla f(\boldsymbol{y}^t)$$
$$\boldsymbol{y}^{t+1} = \boldsymbol{x}^{t+1} + \frac{t}{t+3}(\boldsymbol{x}^{t+1} - \boldsymbol{x}^t)$$

Y. Nesterov

# Nesterov Accelerate Gradient

Iteration complexity: $\mathcal{O}(\frac{1}{\sqrt{\epsilon}})$, much faster than gradient descent: $\mathcal{O}(\frac{1}{\epsilon})$, alternates between gradient updates and proper extrapolation. It is not a descent method (i.e. we may not have $f(x^{t+1}) \leq f(x^t)$). It is one of the **most beautiful and mysterious** results in optimization.

Gradient Descent
00000

Convergence of Gradient Descent
00000000000

Variations
0000000

## Conclusion

|                                | stepsize rule       | convergence rate                                      | iteration complexity                              |
| ------------------------------ | ------------------- | ----------------------------------------------------- | ------------------------------------------------- |
| convex & smooth problems       | $\eta_t = \frac{1}{L}$ | $O\left(\frac{1}{t^2}\right)$                         | $O\left(\frac{1}{\sqrt{\varepsilon}}\right)$      |
| strongly convex & smooth problems | $\eta_t = \frac{1}{L}$ | $O\left(\left(1 - \frac{1}{\sqrt{\kappa}}\right)^t\right)$ | $O\left(\sqrt{\kappa}\log\frac{1}{\varepsilon}\right)$ |