

# Nonlinear Optimization

Maestría en Cómputo Estadístico

Centro de Investigación en Matemáticas A.C.



# What is nonlinear optimization? I

- A nonlinear optimization problem is also referred to as a “nonlinear programming”
- It consists of the optimization of a function subject to constraints, such that any function can be nonlinear
- Linear optimization is a special case of nonlinear optimization



# What is nonlinear optimization? II

- Nonlinear optimization models arise often in science and engineering
- Some remarkable examples:
  - Support vector machines
  - Neural network training
  - Shape optimization
  - Image reconstruction



# Optimality Conditions I

- The nonlinear optimization methods that we will discuss are not always guaranteed to find a global minimum
- These algorithms terminate at a stationary point which usually is at least a local minimum
- For unconstrained nonlinear optimization, there are two conditions that, when satisfied, imply that an optimal solution has been found



# Optimality Conditions II

- A well-known result is that an objective function  $f$  defined and differentiable, a point  $\mathbf{x}^*$  is a local maximum or minimum if its gradient is zero
- The above is known as the *first order optimality condition*, and can be formally stated as:

## First order optimality condition

Let  $f : \mathcal{X} \rightarrow \mathbb{R}$  be a function defined on a set  $\mathcal{X} \subseteq \mathbb{R}^n$ . Suppose that  $\mathbf{x}^* \in \mathcal{X}$  is a local optimum point and that the partial derivatives of  $f$  exist at  $\mathbf{x}^*$ , then  $\nabla f(\mathbf{x}^*) = 0$



# Optimality Conditions III

- A second condition uses the second derivative and is stated as

## Second order optimality condition

Let  $f : \mathcal{X} \rightarrow \mathbb{R}$  be a function defined on a set  $\mathcal{X} \subseteq \mathbb{R}^n$ . Suppose that  $f$  is twice continuously differentiable over  $\mathcal{X}$  and that  $\mathbf{x}^*$  is a stationary point, then

- if  $\mathbf{x}^*$  is a local minimum of  $f$  over  $\mathcal{X}$ , then  $\nabla^2 f(\mathbf{x}^*)$  is positive semi-definite
- if  $\mathbf{x}^*$  is a local maximum of  $f$  over  $\mathcal{X}$ , then  $\nabla^2 f(\mathbf{x}^*)$  is negative semi-definite



# Steepest descent method Method for Minimization I

- The Steepest descent method method is a classical method for nonlinear optimization
- It is one of the simplest and the most fundamental minimization methods for unconstrained optimization
- It makes use of the first-order necessary condition for a local minimizer

$$\nabla f(\mathbf{x}) = 0 \quad (1)$$



# Steepest descent method Method for Minimization II

- The main idea is to consider the search as an interactive procedure
- An initial search point is required and is updated in each iteration based on a step size,  $\alpha$
- The direction is given by  $\nabla f(\mathbf{x})$ , and the step size determines how far we go in that particular direction
- A point is updated as follows:

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





# Steepest descent method Method for Minimization III

---

**Algorithm 1** Steepest descent

---

- 1: Set an initial search point  $\mathbf{x}^{(0)}$
  - 2: Set a step size,  $\alpha \geq 0$
  - 3: Set  $t \leftarrow 0$
  - 4: **while**  $\| \nabla f(\mathbf{x}^{(t)}) \| > \epsilon$  **do**
  - 5:    $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} - \alpha \nabla f(\mathbf{x}^{(t)})$
  - 6:    $t \leftarrow t + 1$
  - 7: **end while**
  - 8: **return**  $\mathbf{x}^{(t)}$
- 



# Example I

- Minimize  $f(\mathbf{x}) = \frac{3}{2}x_1^2 + \frac{1}{2}x_2^2 - x_1x_2 - 2x_1$ . Assume that the initial point is  $\mathbf{x}^{(0)} = (-2, 4)$  and the step size is 0.5

## Solution:

- $\nabla f(\mathbf{x}) = (3x_1 - x_2 - 2, x_2 - x_1)$
- $\nabla f(\mathbf{x}^{(0)}) = (-12, 6)$   
 $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - 0.5\nabla f(\mathbf{x}^{(0)})$   
 $\mathbf{x}^{(1)} = (-2 + 6, 4 - 3) = (4, 1)$



## Example II

$$\textcircled{3} \quad \nabla f(\mathbf{x}^{(1)}) = (9.000, -3.000)$$

$$\mathbf{x}^{(2)} = (4.000 - 4.500, 1.000 + 1.500) = (-0.500, 2.500)$$

$$\textcircled{4} \quad \nabla f(\mathbf{x}^{(2)}) = (-6.000, 3.000)$$

$$\mathbf{x}^{(2)} = (-0.500 + 3.000, 2.500 - 1.500) = (2.500, 1.000)$$

$$\textcircled{5} \quad \nabla f(\mathbf{x}^{(3)}) = (4.500, -1.500)$$

$$\mathbf{x}^{(4)} = (2.500 - 2.250, 1.000 + 0.750) = (0.250, 1.750)$$

$$\textcircled{6} \quad \nabla f(\mathbf{x}^{(4)}) = (-3.000, 1.500)$$

$$\mathbf{x}^{(5)} = (0.250 + 1.500, 1.750 - 0.750) = (1.750, 1.000)$$

$$\textcircled{7} \quad \nabla f(\mathbf{x}^{(5)}) = (2.250, -0.750)$$

$$\mathbf{x}^{(6)} = (1.750 - 1.125, 1.000 + 0.375) = (0.625, 1.375)$$



# Example III

- Repeat the process with:
  - Step size equals to 1.0
  - Step size equals to 0.1



# Steepest Descent

- In this method, the convergence strongly depends on the step size
- The optimal step size can be determined by solving

$$\min f \left( \mathbf{x}^{(t)} - \alpha \nabla f \left( \mathbf{x}^{(t)} \right) \right) \quad (3)$$

- Repeat the previous exercise computing the optimal step size in each iteration



# The Newton's Method I

- The main idea is to start with an arbitrary point, and iteratively try to find the point at which the function evaluates to zero
- Ergo, it iteratively uses the quadratic approximation of the objective function
- An initial solution is required and updated in each iteration as follows:

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \frac{\nabla f(\mathbf{x}^{(t)})}{H_f(\mathbf{x}^{(t)})} \quad (4)$$



# The Newton's Method II

---

**Algorithm 2** Newton's Method

---

- 1: Set an initial search point  $\mathbf{x}^{(0)}$
  - 2: Set  $t \leftarrow 0$
  - 3: **while**  $\| \nabla f(\mathbf{x}^{(t)}) \| > \epsilon$  **do**
  - 4:    $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \frac{\nabla f(\mathbf{x}^{(t)})}{H_f(\mathbf{x}^{(t)})}$
  - 5:    $t \leftarrow t + 1$
  - 6: **end while**
  - 7: **return**  $\mathbf{x}^{(t)}$
- 



# Example I

- Minimize  $f(\mathbf{x}) = \frac{3}{2}x_1^2 + \frac{1}{2}x_2^2 - x_1x_2 - 2x_1$ . Assume that the initial point is  $\mathbf{x}^{(0)} = (-2, 4)$

## Solution:

- Gradient:

$$\nabla f(\mathbf{x}) = (3x_1 - x_2 - 2, x_2 - x_1) \quad (5)$$

- Hessian matrix:

$$\begin{bmatrix} 3 & -1 \\ -1 & 1 \end{bmatrix} \quad (6)$$





## Example II

### Solution:

$$\textcircled{1} \nabla f(\mathbf{x}^{(0)}) = (-12, 6)$$

$$\textcircled{2} \mathbf{x}^{(2)} = (-2.000 + 3.00, 4.000 - 3.000) = (1.000, 1.000)$$

$$\textcircled{3} \nabla f(\mathbf{x}^{(1)}) = (0, 0)$$



# Exercise

Solve the following optimization problems using both Steepest Descent and Newton's Method

①  $x_1^4 + 2x_1x_2 + x_2^4$

②  $x_1 - x_2 + 2x_1^2 + 2x_1x_2 + x_2^2$

③  $2x_1^2 + x_2^2 - 2x_1x_2 + 2x_1^3 + x_1^4$



# Conjugate Gradient Descent I

- Conjugate gradient descent can be regarded as being between the method of steepest descent and Newton's method
- It aims at accelerating the convergence of steepest descent while avoids the computational cost of the Newton's method
- Unlike previous methods, in conjugate gradient descent, at each iteration the direction is modified by a combination of the earlier directions



# Conjugate Gradient Descent II

---

## Algorithm 3 Conjugate Gradient Descent Method

---

- 1: Set an initial search point  $\mathbf{x}^{(0)}$
  - 2: Set a step size,  $\alpha \geq 0$
  - 3: Set  $t \leftarrow 0$
  - 4:  $s(t) \leftarrow -\nabla f(\mathbf{x}^{(t)})$
  - 5:  $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} + \alpha s(t)$
  - 6: **while**  $\|\nabla f(\mathbf{x}^{(t)})\| > \epsilon$  **do**
  - 7:    $t \leftarrow t + 1$
  - 8:    $s(t) = -\nabla f(\mathbf{x}^{(t)}) + \frac{\|\nabla f(\mathbf{x}^{(t)})\|^2}{\|\nabla f(\mathbf{x}^{(t-1)})\|^2} s(t-1)$
  - 9:    $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} + \alpha s(t)$
  - 10: **end while**
  - 11: **return**  $\mathbf{x}^{(t)}$
- 



# Example I

- Minimize  $f(\mathbf{x}) = \frac{3}{2}x_1^2 + \frac{1}{2}x_2^2 - x_1x_2 - 2x_1$ . Assume that the initial point is  $\mathbf{x}^{(0)} = (-2, 4)$  and  $\alpha = 0.5$

## Solution:

- Gradient:

$$\nabla f(\mathbf{x}) = (3x_1 - x_2 - 2, x_2 - x_1) \quad (7)$$

- $\nabla f(\mathbf{x}^{(0)}) = (-12, 6)$   
 $s(0) = (12, -6)$   
 $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha s(0) = (4, 1)$
- $\nabla f(\mathbf{x}^{(1)}) = (9, -3)$   
 $s(1) = (-9, 3) + \frac{1}{2}(12, -6) = (-3, 0)$   
 $\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha s(1) = (2.5, 1.0)$



## Example II

$$\textcircled{3} \quad \nabla f(\mathbf{x}^{(2)}) = (4.500, -1.500)$$

$$s(2) = (-4.500, 1.500) + \frac{1}{4}(-3.000, 0.000) = (-5.250, 1.500)$$

$$\mathbf{x}^{(3)} = \mathbf{x}^{(2)} + \alpha s(2) = (-0.125, 1.750)$$

$$\textcircled{4} \quad \nabla f(\mathbf{x}^{(3)}) = (-4.125, 1.875)$$

$$s(3) = (4.125, -1.875) + \frac{657}{720}(-5.250, 1.500) = (-0.6666, -0.506)$$

$$\mathbf{x}^{(4)} = \mathbf{x}^{(3)} + \alpha s(3) = (-0.458, 1.497)$$



# Conjugate Gradient Descent: Variants I

The conjugate gradient direction can be rewritten as:

$$s(t) = -\nabla f(\mathbf{x}^{(t)}) + \beta s(t-1) \quad (8)$$

Some variants are listed below:

- **Fletcher-Reeves:**

$$\beta = \frac{\nabla f(\mathbf{x}^{(t)})^T \nabla f(\mathbf{x}^{(t)})}{\nabla f(\mathbf{x}^{(t-1)})^T \nabla f(\mathbf{x}^{(t-1)})} \quad (9)$$



# Conjugate Gradient Descent: Variants II

- **Polak–Ribière:**

$$\beta = \frac{\nabla f(\mathbf{x}^{(t)})^T (\nabla f(\mathbf{x}^{(t)}) - \nabla f(\mathbf{x}^{(t-1)}))}{\nabla f(\mathbf{x}^{(t-1)})^T \nabla f(\mathbf{x}^{(t)})} \quad (10)$$

- **Hestenes-Stiefel:**

$$\beta = \frac{\nabla f(\mathbf{x}^{(t)})^T (\nabla f(\mathbf{x}^{(t)}) - \nabla f(\mathbf{x}^{(t-1)}))}{s(t-1)^T (\nabla f(\mathbf{x}^{(t)}) - \nabla f(\mathbf{x}^{(t-1)}))} \quad (11)$$





# Quasi-Newton methods I

- Quasi-Newton methods, as the name suggest, are an alternative to the Newton's method
- They are specially useful when computing the Hessian matrix is computational prohibitive
- They rank among the most efficient methods available



# Quasi-Newton methods II

- The basic principle in quasi-Newton methods is that the direction of search is based on an  $n \times n$  direction matrix **S** that is generated from available data and is contrived to be an approximation of  $\mathbf{H}^{-1}$
- Quasi-Newton methods updates the “Hessian matrix” by analyzing successive gradient vectors



# Quasi-Newton methods III

- In general, a point is updated as follows:

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha \mathbf{S}^{(t)} \nabla f(\mathbf{x}^{(t)}) \quad (12)$$

where

$$\mathbf{S}^{(t)} = \begin{cases} \mathbf{I}_n & \text{for the steepest descent method} \\ \mathbf{H}_t^{-1} & \text{for the Newton method} \end{cases}$$

- The idea is, therefore, to choose some positive definite  $\mathbf{S}^{(t)}$  such that is equal to or, at least, approximately equal to  $\mathbf{H}_t^{-1}$



# Quasi-Newton methods IV

---

**Algorithm 4** Basic Quasi-Newton Method

---

- 1: Set  $t \leftarrow 0$
  - 2: Set an initial search point  $\mathbf{x}^{(t)}$
  - 3: Set a step size,  $\alpha \geq 0$
  - 4: Set  $\mathbf{S}^{(t)} \leftarrow \mathbf{I}_n$
  - 5: **while**  $\|\nabla f(\mathbf{x}^{(t)})\| > \epsilon$  **do**
  - 6:    $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha \mathbf{S}^{(t)} \nabla f(\mathbf{x}^{(t)})$
  - 7:   Update  $\mathbf{S}^{(t+1)}$  by following a proper strategy
  - 8:    $t \leftarrow t + 1$
  - 9: **end while**
  - 10: **return**  $\mathbf{x}^{(t)}$
- 



# Newton's Method vs Quasi-Newton Methods

Quasi-Newton Method	Newton's Method
Only need the function values and gradients	Need the function values, gradients and Hessians
$H_t^{-1}$ maintains positive definite for several updates	$H_t^{-1}$ is not sure to be positive definite
Need $\mathcal{O}(n^2)$ multiplications in each iteration	Need $\mathcal{O}(n^3)$ multiplications in each iteration



# How can we update the Hessian Matrix? I

There are several approaches to approximate the Hessian matrix

- **Symmetric Rank-One:**

$$H_{t+1} = H_t + \frac{(y_{(t)} - H_t s_{(t)}) (y_{(t)} - H_t s_{(t)})^T}{(y_{(t)} - H_t s_{(t)})^T s_{(t)}} \quad (13)$$

where:

$$\begin{aligned} y_{(t)} &= \nabla f(\mathbf{x}^{(t)}) - \nabla f(\mathbf{x}^{(t-1)}) \\ s_{(t)} &= \mathbf{x}^{(t)} - \mathbf{x}^{(t-1)} \end{aligned}$$



# How can we update the Hessian Matrix? II

- **Davidon–Fletcher–Powell (DFP):**

$$H_{t+1} = H_t + \frac{s_{(t)}s_{(t)}^T}{s_{(t)}^T y_{(t)}} - \frac{(H_t y_{(t)}) (H_t y_{(t)})^T}{y_{(t)}^T H_t y_{(t)}} \quad (14)$$

- **Broyden–Fletcher–Goldfarb–Shanno (BFGS):**

$$H_{t+1} = H_t + \frac{y_{(t)}y_{(t)}^T}{y_{(t)}^T s_{(t)}} - \frac{H_t s_{(t)} s_{(t)}^T H_t^T}{s_{(t)}^T H_t s_{(t)}} \quad (15)$$



# Example I

- Minimize  $f(\mathbf{x}) = \frac{3}{2}x_1^2 + \frac{1}{2}x_2^2 - x_1x_2 - 2x_1$ . Assume that the initial point is  $\mathbf{x}^{(0)} = (-2, 4)$  and  $\alpha = 0.5$

## Solution:

- Initial Hessian matrix

$$H_0 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

- The gradient vector:  $\nabla f(\mathbf{x}) = (3x_1 - x_2 - 2, x_2 - x_1)$
- $\nabla f(\mathbf{x}^{(0)}) = (-12, 6)$   
 $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \alpha H_0 \nabla f(\mathbf{x}^{(0)}) = (4, 1)$





## Example II

$$\textcircled{4} \quad \nabla f(\mathbf{x}^{(1)}) = (9.000, -3.000)$$

$$S_0 = \mathbf{x}^{(1)} - \mathbf{x}^{(0)} = (6.000, -3.000)$$

$$y_0 = \nabla f(\mathbf{x}^{(1)}) - \nabla f(\mathbf{x}^{(0)}) = (21.000, -9.000)$$

$$\begin{aligned} H_1 &= H_0 + \frac{(y_0 - H_0 s_0)(y_0 - H_0 s_0)^T}{(y_0 - H_0 s_0)^T s_0} \\ &= \begin{bmatrix} 3.083 & -0.833 \\ -0.833 & 1.333 \end{bmatrix} \end{aligned}$$

$$\textcircled{5} \quad \mathbf{x}^{(2)} = \mathbf{x}^{(1)} - \alpha H_1^{-1} \nabla f(\mathbf{x}^{(1)}) = (2.610, 1.256)$$



## Example III

$$\textcircled{6} \quad \nabla f(\mathbf{x}^{(2)}) = (4.573, -1.354)$$

$$S_0 = \mathbf{x}^{(2)} - \mathbf{x}^{(1)} = (-1.390, 0.256)$$

$$y_0 = \nabla f(\mathbf{x}^{(2)}) - \nabla f(\mathbf{x}^{(1)}) = (-4.427, 1.646)$$

$$\begin{aligned} H_2 &= H_1 + \frac{(y_1 - H_1 s_1)(y_1 - H_1 s_1)^T}{(y_1 - H_1 s_0)^T s_1} \\ &= \begin{bmatrix} 3.000 & -1.000 \\ -1.000 & 1.000 \end{bmatrix} \end{aligned}$$

$$\textcircled{7} \quad \mathbf{x}^{(3)} = \mathbf{x}^{(2)} - \alpha H_2^{-1} \nabla f(\mathbf{x}^{(2)}) = (1.805, 1.128)$$



# Quasi-Newton Methods

- Note, however, that the above approaches still require computing the inverse of the Hessian matrix
- Can we instead update the inverse of the Hessian matrix?
- The BFGS formula can be updated as follows:

$$H_{t+1}^{-1} = \left( \mathbf{I}_n - \frac{s_t y_t^T}{s_t^T y_t} \right) H_t^{-1} \left( \mathbf{I}_n - \frac{y_t s_t^T}{s_t^T y_t} \right) + \frac{s_t s_t^T}{s_t^T y_t} \quad (16)$$



# Questions?

