

# Non-Linear Optimization Methods - Part II

Optimization Methods in Management Science

Master in Management

HEC Lausanne

Dr. Rodrigue Ouevray

Fall 2019 Semester

# Non-Linear Optimization Methods

## **Unconstrained** optimization:

- Descent methods
- Linesearch
- Newton's method

# Context

- We now consider the more general problem described below:

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}),$$

where  $f$  is typically twice continuously differentiable

- We don't assume any more that  $f(\mathbf{x})$  is a quadratic function
- It can be **any** twice continuously differentiable function
- A **descent method** is an **iterative** optimization algorithm for finding a **local optimum** of such a function

# Descent Direction

- A direction  $\mathbf{d}_k$  for which  $\nabla f(\mathbf{x}_k)^T \mathbf{d}_k < 0$  is called a **descent** direction
- If  $\mathbf{d}_k$  is a **descent** direction, then we have the guarantee that

$$f(\mathbf{x}_k + \alpha \mathbf{d}_k) < f(\mathbf{x}_k)$$

if  $\alpha > 0$  is sufficiently small

- This just means that we can **decrease** the objective function if we move into the direction  $\mathbf{d}_k$
- $\nabla f(\mathbf{x}_k)^T \mathbf{d}_k < 0$  can be interpreted as the **directional derivative** of  $f$  at  $\mathbf{x}_k$  along a vector  $\mathbf{d}_k$

## Descent Direction (2)

### Proposition

*We assume that  $\mathbf{Q}$  is a positive definite matrix and that  $\nabla f(\mathbf{x}_k) \neq \mathbf{0}$ . Then any direction  $\mathbf{d}_k$  given by  $\mathbf{d}_k = -\mathbf{Q}\nabla f(\mathbf{x}_k)$  is a descent direction.*

**Proof :**  $-\nabla f(\mathbf{x}_k)^T \mathbf{Q} \nabla f(\mathbf{x}_k) < 0$ .

# Steepest Descent Direction

- The **steepest** descent is the direction given by the **opposite** of the gradient  $-\nabla f(\mathbf{x})$
- This is a descent direction that is optimal in the sense that

$$-\nabla f(\mathbf{x})^T \nabla f(\mathbf{x}) \leq \mathbf{d}^T \nabla f(\mathbf{x})$$

for any  $\mathbf{d}$  such that  $\|\mathbf{d}\| = \|\nabla f(\mathbf{x})\|$

# Descent Method Framework

**Descent methods** comprise the following steps:

- Find a direction  $\mathbf{d}_k$  such that  $\nabla f(\mathbf{x}_k)^T \mathbf{d}_k < 0$
- Find a **step**  $\alpha_k$  such that  $f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) \leq f(\mathbf{x}_k)$
- Compute  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$ ,  $k = k + 1$ , and repeat the process until a **stopping criterion** is satisfied

## Example

We would like to **minimize** the following function  $f$ :

$$f(x_1, x_2) = \frac{1}{2}(x_1^2 + 10x_2^2)$$

The **solution** is the point  $(0, 0)$ . The gradient of  $f$  is given by:

$$\nabla f(x_1, x_2) = (x_1 \ 10x_2)^T$$

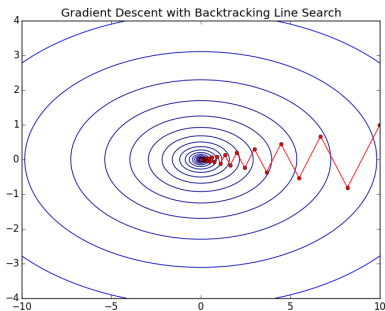
We would like to apply the **steepest descent** method with the initial point given by  $\mathbf{x} = (10, 1)$ . At this point, the descent direction is

$$\mathbf{d}_1 = -\nabla f(10, 1) = (-10 \ -10)^T$$



## Example (Cont'd)

Here is the sequence of iterates produced by the steepest descent method:



We can see that this sequence of iterates **converge** to the **optimal** solution by **zig-zaging**. We will see later in this presentation much performant algorithms than the steepest descent method

## Length of the Step

- Indeed, it is not necessary to solve the following problem:

$$\alpha_k = \operatorname{argmin}_{\alpha \geq 0} f(\mathbf{x}_k + \alpha \mathbf{d}_k)$$

- We just need to find a step  $\alpha_k$  that reduces  $f$  **sufficiently** in order to have an algorithm that **converges** to the **optimal** solution
- This step should be **easy** to compute
- The challenges in finding a good  $\alpha_k$  are both in avoiding that the step length is **too long**, or **too short**
- If the step is too long or too short, then this can **prevent** the convergence of the descent method to the **optimal** solution

## Example

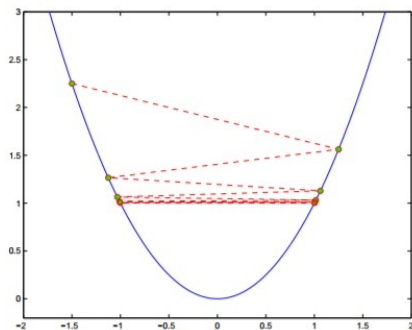
- We would like to **minimize**  $f(x) = x^2$  with a descent method
- The **minimum** is at  $x = 0$  and its value is 0
- The sequence of iterates is generated by the following formula:

$$x_{k+1} = x_k + \alpha_k d_k$$

- We would like that this sequence of iterates **converge** to  $x = 0$
- With a problem with one variable, there is only **two possible directions**  $d_k$ :  $d_k = -1$  (we move to the left from the current iterate) and  $d_k = 1$  (we move to the right)
- Let's assume that the **initial point** is given by  $x_0 = 2$

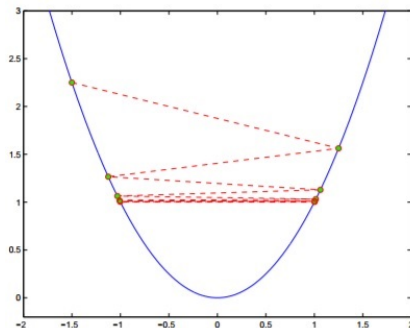
## Example: Steps Too Long...

- Let's consider the case where  $d_k = (-1)^{k+1}$  and  $\alpha_k = 2 + 3(2^{-k-1})$
- $k = 0, d_0 = -1, \alpha_0 = 3.5, k = 1, d_1 = 1, \alpha_1 = 2.75, k = 2, d_2 = -1, \alpha_2 = 2 + 3/8, \dots$
- Starting from  $x_0 = 2$ , the descent method generates the following sequence of iterates:  $x_0 = 2, x_1 = -1.5, x_2 = 1.25, \dots$



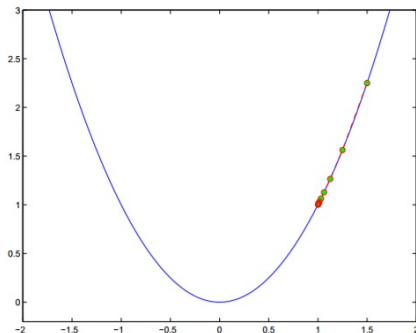
## Example: Steps Too Long... (Cont'd)

- The objective function decreases at each iteration but  $\{x_k\}$  never converge to 0 !
- When  $k$  is even (resp. odd), the sequence of iterates converge to 1 (resp. to -1). But the **optimum** is at  $x = 0$  !



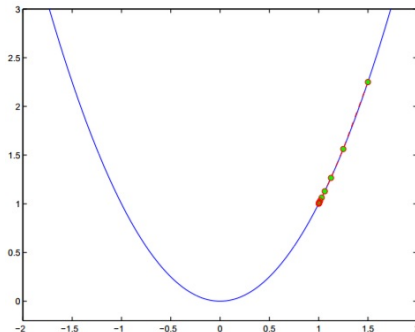
## Example: Steps Too Short...

- Same example as before but we consider the case where  $d_k = -1$  and  $\alpha_k = 2^{-k-1}$
- $k = 0, d_0 = -1, \alpha_0 = 0.5, k = 1, d_1 = -1, \alpha_0 = 0.25, k = 2, d_2 = -1, \alpha_0 = 0.125$
- The descent method generates the following sequence of iterates:  $x_0 = 2, x_1 = 1.5, x_2 = 1.25, \dots$



## Example: Steps Too Short...(Cont'd)

- The objective function decreases at each iteration but  $\{x_k\}$  converge to 1 and not to 0 !



- These examples just show that the steps  $\alpha_k$  must satisfy some conditions to guarantee that the algorithm converges to an optimal solution

# Wolfe Conditions

The conditions that  $f$  must satisfy in terms of decrease to converge to the optimum are called the **Wolfe conditions**

## Theorem

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a differentiable function, a point  $\mathbf{x} \in \mathbb{R}^n$ , a descent direction  $\mathbf{d}_k$  such that  $\nabla f(\mathbf{x}_k)^T \mathbf{d}_k < 0$  and a step  $\alpha_k > 0$ . We say that the step length  $\alpha_k$  satisfies the **Wolfe conditions** if

$$(1) \quad f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) \leq f(\mathbf{x}_k) + \alpha_k \beta_1 \nabla f(\mathbf{x}_k)^T \mathbf{d}_k \quad (\text{Armijo rule})$$

$$(2) \quad \nabla f(\mathbf{x}_k + \alpha_k \mathbf{d}_k)^T \mathbf{d}_k \geq \beta_2 \nabla f(\mathbf{x}_k)^T \mathbf{d}_k \quad (\text{curvature condition})$$

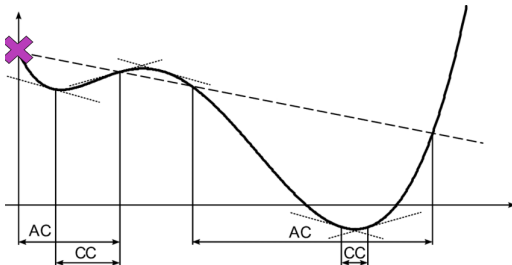
with  $0 < \beta_1 < \beta_2 < 1$

**Armijo rule** ensures that the step length  $\alpha_k$  decreases  $f$  sufficiently, and the curvature condition ensures that the slope has been reduced sufficiently



# Illustration of the Wolfe Conditions

- Let's define the function  $\phi$  as follows:  $\phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{d}_k)$



- $\phi(0) = f(\mathbf{x}_k)$  corresponds to the violet cross
- The dotted line starting from the violet cross is given by  $g(\alpha) = f(\mathbf{x}_k) + \alpha \beta_1 \nabla f(\mathbf{x}_k)^T \mathbf{d}_k$
- AC represents the set of  $\alpha$  satisfying the Armijo condition
- CC corresponds to set of  $\alpha$  satisfying the curvature condition
- The set of  $\alpha$  satisfying both conditions are given by the intersection between AC and CC

## Wolfe Conditions: Exercise

We consider the following function:

$$\begin{aligned} f &: \mathbb{R}^2 \rightarrow \mathbb{R} \\ (x, y) &\mapsto f(x, y) = x^4 + x^2 + y^2 \end{aligned}$$

and the initial data:

$$\mathbf{x}_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \beta_1 = 0.1 \quad \text{and} \quad \beta_2 = 0.5$$

- a) Compute  $\nabla f(\mathbf{x}_1)$
- b) Show that  $\mathbf{d}_1 = \begin{pmatrix} -3 \\ -1 \end{pmatrix}$  is a descent direction at  $\mathbf{x}_1$
- c) Determine if the steps  $\alpha_1 = 1$ ,  $\alpha_2 = 0.1$  and  $\alpha_3 = 0.5$  satisfy Wolfe's conditions

## Wolfe Conditions: Exercise (Cont'd)

a)  $\nabla f(\mathbf{x}) = \begin{pmatrix} 4x^3 + 2x \\ 2y \end{pmatrix} \quad \forall \mathbf{x} \in \mathbb{R}^2, \quad \nabla f(1, 1) = \begin{pmatrix} 6 \\ 2 \end{pmatrix}$

b)  $\mathbf{d}_1$  is a descent direction:

$$\nabla f(\mathbf{x}_1)^T \mathbf{d}_1 = \begin{pmatrix} 6 & 2 \end{pmatrix} \begin{pmatrix} -3 \\ -1 \end{pmatrix} = -20 < 0$$

## Wolfe Conditions: Exercise (Cont'd)

- c)  $f(x, y) = x^4 + x^2 + y^2$ ,  $\nabla f(\mathbf{x}) = (4x^3 + 2x, 2y)^T$ . For  $\alpha_1 = 1$ , let's check Wolfe conditions with  $\beta_1 = 0.1$  and  $\beta_2 = 0.5$

$$(1) \quad f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) \leq f(\mathbf{x}_k) + \alpha_k \beta_1 \nabla f(\mathbf{x}_k)^T \mathbf{d}_k \quad (\text{Armijo rule})$$

$$(2) \quad \nabla f(\mathbf{x}_k + \alpha_k \mathbf{d}_k)^T \mathbf{d}_k \geq \beta_2 \nabla f(\mathbf{x}_k)^T \mathbf{d}_k \quad (\text{curvature condition})$$

Application :  $\mathbf{x}_1 = (1, 1)$ ,  $\nabla f(1, 1) = (6, 2)^T$ ,  $\mathbf{d}_1 = (-3, -1)^T$ ,  
 $\mathbf{x}_2 = \mathbf{x}_1 + \alpha \mathbf{d}_1 = (-2, 0)$  and  $\nabla f(\mathbf{x}_2) = (-36, 0)^T$

- It satisfies the curvature condition but not the Armijo rule:

$$(1) \quad f(-2, 0) = 20 \not\leq 1 = 3 + 1 \cdot 0.1 \cdot \begin{pmatrix} 6 & 2 \end{pmatrix} \begin{pmatrix} -3 \\ -1 \end{pmatrix}$$

$$(2) \quad \begin{pmatrix} -36 & 0 \end{pmatrix} \begin{pmatrix} -3 \\ -1 \end{pmatrix} = 108 \geq -10 = 0.5 \cdot \begin{pmatrix} 6 & 2 \end{pmatrix} \begin{pmatrix} -3 \\ -1 \end{pmatrix}$$

## Wolfe Conditions: Exercise (Cont'd)

- c)  $f(x, y) = x^4 + x^2 + y^2$ ,  $\nabla f(\mathbf{x}) = (4x^3 + 2x, 2y)^T$ . For  $\alpha_2 = 0.1$ , let's check Wolfe conditions with  $\beta_1 = 0.1$  and  $\beta_2 = 0.5$

$$(1) \quad f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) \leq f(\mathbf{x}_k) + \alpha_k \beta_1 \nabla f(\mathbf{x}_k)^T \mathbf{d}_k \quad (\text{Armijo rule})$$

$$(2) \quad \nabla f(\mathbf{x}_k + \alpha_k \mathbf{d}_k)^T \mathbf{d}_k \geq \beta_2 \nabla f(\mathbf{x}_k)^T \mathbf{d}_k \quad (\text{curvature condition})$$

Application :  $\mathbf{x}_1 = (1, 1)$ ,  $\nabla f(1, 1) = (6, 2)^T$ ,  $\mathbf{d}_1 = (-3, -1)^T$ ,  
 $\mathbf{x}_2 = \mathbf{x}_1 + \alpha \mathbf{d}_1 = (0.7, 0.9)$  and  $\nabla f(\mathbf{x}_2) = (-2.772, 1.8)^T$

- It satisfies Armijo rule but not the curvature condition:

$$(1) \quad f(0.7, 0.9) = 1.54 \leq 2.8 = 3 + 0.1 \cdot 0.1 \cdot (-20)$$

$$(2) \quad \begin{pmatrix} -2.772 & 1.8 \end{pmatrix} \begin{pmatrix} -3 \\ -1 \end{pmatrix} = -10.116 \not\geq -10 = 0.5 \cdot \begin{pmatrix} 6 & 2 \end{pmatrix} \begin{pmatrix} -3 \\ -1 \end{pmatrix}$$

## Wolfe Conditions: Exercise (Cont'd)

- c)  $f(x, y) = x^4 + x^2 + y^2$ ,  $\nabla f(\mathbf{x}) = (4x^3 + 2x, 2y)^T$ . For  $\alpha_3 = 0.5$ , let's check Wolfe conditions with  $\beta_1 = 0.1$  and  $\beta_2 = 0.5$

$$(1) \quad f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) \leq f(\mathbf{x}_k) + \alpha_k \beta_1 \nabla f(\mathbf{x}_k)^T \mathbf{d}_k \quad (\text{Armijo rule})$$

$$(2) \quad \nabla f(\mathbf{x}_k + \alpha_k \mathbf{d}_k)^T \mathbf{d}_k \geq \beta_2 \nabla f(\mathbf{x}_k)^T \mathbf{d}_k \quad (\text{curvature condition})$$

Application :  $\mathbf{x}_1 = (1, 1)$ ,  $\nabla f(1, 1) = (6, 2)^T$ ,  $\mathbf{d}_1 = (-3, -1)^T$ ,  
 $\mathbf{x}_2 = \mathbf{x}_1 + \alpha \mathbf{d}_1 = (-0.5, 0.5)$  and  $\nabla f(\mathbf{x}_2) = (-1.5, 1)^T$

- It satisfies both conditions

$$(1) \quad f(-0.5, 0.5) = 0.5625 \leq 2 = 3 + 0.5 \cdot 0.1 \cdot (-20)$$

$$(2) \quad \begin{pmatrix} -1.5 & 1 \end{pmatrix} \begin{pmatrix} -3 \\ -1 \end{pmatrix} = 3.5 \geq -10 = 0.5 \cdot \begin{pmatrix} 6 & 2 \end{pmatrix} \begin{pmatrix} -3 \\ -1 \end{pmatrix}$$

# Linesearch

- A **linesearch** is a method that produces a step  $\alpha^*$  satisfying the Wolfe conditions **at each iteration**
- Let's assume that we are at **iteration**  $k$ . How to generate a step  $\alpha_k$  satisfying these conditions ?
  - ▶ We start from a initial step  $\alpha_0$
  - ▶ If this step violates the Armijo rule, then the step is too long and we reduce it
  - ▶ If this step violates the curvature rule, then the step is too short and we increase it

# Backtracking Linesearch

- In practice, we use a **backtracking linesearch**
- The **backtracking linesearch** starts with a **large** step  $\alpha_0$  and iteratively shrinks it to satisfy the Armijo rule if it is necessary
- The curvature condition can be **dispensed**
- This makes this algorithm very **efficient** to find a step satisfying the Wolfe conditions



# Newton's Method

We consider the following problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}),$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is twice continuously differentiable

## Newton's method:

- We **approximate**  $f$  by a **quadratic** function given by its **Taylor series expansion**
- The **minimization** of this quadratic function generates the **next** iterate
- We repeat these steps until a stopping criterion is reached

## Newton's Method (Cont'd)

### Important Remark

The **quadratic** approximation  $q(\mathbf{x})$  of  $f(\mathbf{x})$  around the current iterate  $\mathbf{x}_k$  is given by its Taylor series expansion:

$$f(\mathbf{x}) \approx q(\mathbf{x}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T (\mathbf{x} - \mathbf{x}_k) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_k)^T \nabla^2 f(\mathbf{x}_k) (\mathbf{x} - \mathbf{x}_k)$$

The first order condition to minimize the quadratic function is given by  $\nabla q(\mathbf{x}) = 0$ , i.e.:

$$\nabla f(\mathbf{x}_k) + \nabla^2 f(\mathbf{x}_k) (\mathbf{x} - \mathbf{x}_k) = 0$$

## Newton's Method (Cont'd)

- If the matrix  $\nabla^2 f(\mathbf{x}_k)$  is invertible, then the solution  $\mathbf{x}$  of this equation is given by

$$\mathbf{x} = \mathbf{x}_k - \nabla^2 f(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k)$$

- The next iterate  $\mathbf{x}_{k+1}$ :

$$\mathbf{x}_{k+1} = \mathbf{x} = \mathbf{x}_k - \nabla^2 f(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k)$$

- We repeat these steps until a **stopping criterion** is reached

# Newton's Method: Algorithm

- **Input:** a given point  $\mathbf{x}_0 \in \mathbb{R}^n$  and a convergence parameter  $\epsilon > 0$
- **Output:** an approximation  $\mathbf{x}^*$  of the solution
- **Initialization:**  $k = 0$
- **Iterations:**
  - ▶  $\mathbf{x}_{k+1} = \mathbf{x}_k - \nabla^2 f(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k)$
  - ▶  $k = k + 1$
- **Stopping criterion:** if  $\|\nabla f(\mathbf{x})\| \leq \epsilon$ , then  $\mathbf{x}^* = \mathbf{x}_k$

# Newton's Method: Direction

## Important Remark

- With the **steepest descent** method, the direction  $\mathbf{d}_k$  is given by

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

Then we determine a step  $\alpha_k$  and  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$

- With the **Newton's method**, the direction  $\mathbf{d}_k$  is given by

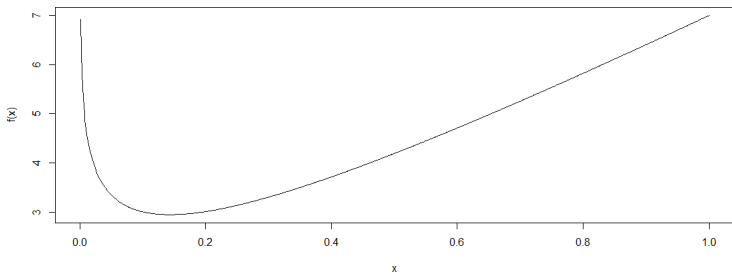
$$\mathbf{d}_k = -\nabla^2 f(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k)$$

and  $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{d}_k$

- When the method converges to the optimal solution, the Newton's direction is often a much better direction than the one provided by the steepest descent

## Newton's Method: Example

- We would like to minimize the function given by  $f(x) = 7x - \ln(x)$
- The domain of  $f$  is  $x > 0$
- Its first derivative is given by  $f'(x) = 7 - \frac{1}{x}$  and its second derivative by  $f''(x) = \frac{1}{x^2}$
- It is not difficult to check that  $x^* = \frac{1}{7}$  is the **unique global minimizer**



## Newton's Method: Example (Cont'd)

- The Newton's direction at  $x$  is

$$d = -\nabla^2 f(x)^{-1} \nabla f(x) = -\frac{f'(x)}{f''(x)} = -x^2(7 - \frac{1}{x}) = x - 7x^2$$

- So, the Newton's method generates the sequence of iterates  $\{x_k\}$  with

$$x_{k+1} = x_k + d_k = x_k + (x_k - 7x_k^2) = 2x_k - 7x_k^2$$

## Newton's Method: Example (Cont'd)

- The table below gives some examples of the sequences generated by this method for **different** starting points  $x_0 = 1, y_0 = 0.1, z_0 = 0.01$ :

$k$	$x_k$	$y_k$	$z_k$
0	1	0.1	0.01
1	-5	0.13	0.0193
2		0.1417	0.0359925
3		0.14284777	0.062916884
4		0.142857142	0.098124028
5		0.142857143	0.128849782
6			0.1414837
7			0.142843938
8			0.142857142
9			0.142857143
10			0.142857143

- When we start at  $x_0 = 1$ , the next iterate is  $x_1 = -5$ . But  $f(x)$  is not defined for  $x \leq 0$  ! The algorithm has failed !
- The sequences  $\{y_k\}$  and  $\{z_k\}$  have converged to  $1/7$  but not  $\{x_k\}$



# Issues with the Newton's Method

- When the method converges, the iterates of Newton's method are attracted to **critical points**. Indeed, the method is just trying to solve the system of equations  $\nabla f(\mathbf{x}) = 0$
- There is **no guarantee** that the sequence of iterates will converge to the optimum
- The **hessian** is assumed to be **nonsingular** at each iteration. Indeed, even if  $\nabla^2 f(\mathbf{x}_k)$  is nonsingular, it may converge to a non-singular matrix
- $\mathbf{d}_k$  is not guaranteed to be a **descent** direction, unless  $\nabla^2 f(\mathbf{x}_k)$  is positive definite
- It may be **too expensive** to compute second order partial derivatives in particular when the dimension of the system is high