

# CS257 Linear and Convex Optimization

## Lecture 10

Bo Jiang

John Hopcroft Center for Computer Science  
Shanghai Jiao Tong University

November 9, 2020

# Recap

**Strong convexity.**  $f$  is  $m$ -strongly convex if

- $f(x) - \frac{m}{2}\|x\|^2$  is convex
- first-order condition

$$f(y) \geq f(x) + \nabla f(x)^T(y - x) + \frac{m}{2}\|y - x\|^2$$

- second-order condition

$$\nabla^2 f(x) \succeq mI \iff \lambda_{\min}(\nabla^2 f(x)) \geq m$$

**Convergence.** For  $m$ -strongly convex and  $L$ -smooth  $f$  with minimum  $x^*$ , gradient descent with constant step size  $t \in (0, \frac{1}{L}]$  satisfies

$$f(x_k) - f(x^*) \leq \frac{L(1 - mt)^k}{m} [f(x_0) - f(x^*)]$$

**Condition number.** For  $Q \succ O$ ,

$$\kappa(Q) = \frac{\lambda_{\max}(Q)}{\lambda_{\min}(Q)}$$

**Well-/Ill-conditioned** if  $\kappa(Q)$  is small/large  $\implies$  fast/slow convergence.

# Today

- exact line search
- backtracking line search
- Newton's method

# Step Size

## Gradient descent

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

- constant step size:  $t_k = t$  for all  $k$
- exact line search: optimal  $t_k$  for each step

$$t_k = \arg \min_s f(\mathbf{x}_k - s \nabla f(\mathbf{x}_k))$$

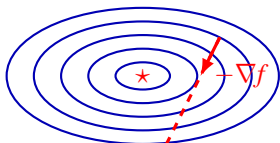
- backtracking line search (Armijo's rule):  $t_k$  satisfies

$$f(\mathbf{x}_k) - f(\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k)) \geq \alpha t_k \|\nabla f(\mathbf{x}_k)\|_2^2$$

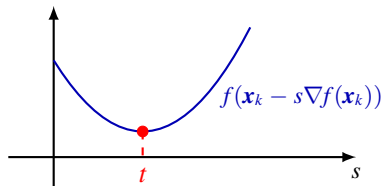
for some given  $\alpha \in (0, 1)$ .

# Exact Line Search

- 1: initialization  $\mathbf{x} \leftarrow \mathbf{x}_0 \in \mathbb{R}^n$
- 2: **while**  $\|\nabla f(\mathbf{x})\| > \delta$  **do**
- 3:      $t \leftarrow \arg \min_s f(\mathbf{x} - s \nabla f(\mathbf{x}))$
- 4:      $\mathbf{x} \leftarrow \mathbf{x} - t \nabla f(\mathbf{x})$
- 5: **end while**
- 6: **return**  $\mathbf{x}$



level curves of  $f(x_1, x_2) = \frac{x_1^2}{4} + x_2^2$



**Note.** Often impractical; used only if the inner minimization is cheap.

# Exact Line Search for Quadratic Functions

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{b}^T \mathbf{x}, \quad \mathbf{Q} \succ \mathbf{0}$$

- gradient at  $\mathbf{x}_k$  is  $\mathbf{g}_k = \nabla f(\mathbf{x}_k) = \mathbf{Q} \mathbf{x}_k + \mathbf{b}$
- second-order Taylor expansion is exact for quadratic functions,

$$\begin{aligned} h(t) &= f(\mathbf{x}_k - t \mathbf{g}_k) \\ &= f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T (-t \mathbf{g}_k) + \frac{1}{2} (-t \mathbf{g}_k)^T \nabla^2 f(\mathbf{x}_k) (-t \mathbf{g}_k) \\ &= \left( \frac{1}{2} \mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k \right) t^2 - \mathbf{g}_k^T \mathbf{g}_k t + f(\mathbf{x}_k) \end{aligned}$$

- minimizing  $h(t)$  yields best step size

$$t_k = \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k}$$

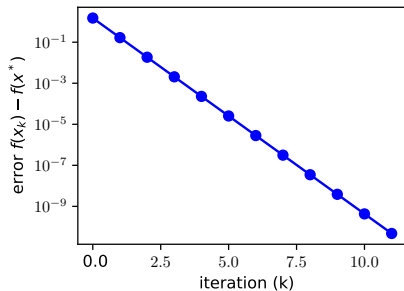
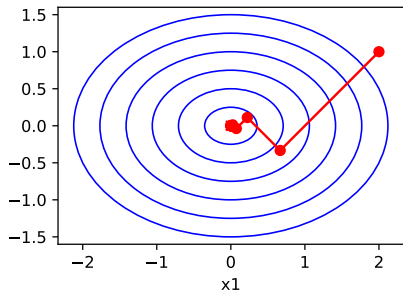
- update step

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t_k \mathbf{g}_k = \mathbf{x}_k - \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{Q} \mathbf{g}_k} \mathbf{g}_k$$

## Example

$$f(x_1, x_2) = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} = \frac{\gamma}{2} x_1^2 + \frac{1}{2} x_2^2, \quad \mathbf{Q} = \text{diag}\{\gamma, 1\}$$

Well-conditioned.  $\gamma = 0.5$ ,  $\mathbf{x}_0 = (2, 1)^T$



Fast convergence.

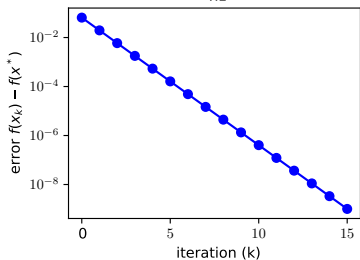
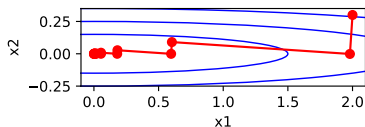
**Note.** Successive gradient directions are always orthogonal, as

$$0 = h'(t_k) = -\nabla f(\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k))^T \nabla f(\mathbf{x}_k) = -\nabla f(\mathbf{x}_{k+1})^T \nabla f(\mathbf{x}_k)$$

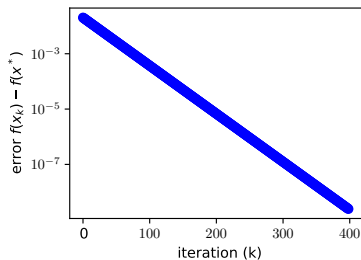
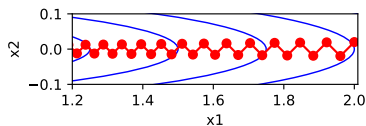
## Example (cont'd)

$$f(x_1, x_2) = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} = \frac{\gamma}{2} x_1^2 + \frac{1}{2} x_2^2, \quad \mathbf{Q} = \text{diag}\{\gamma, 1\}$$

Ill-conditioned.  $\gamma = 0.01$ , convergence rate depends on initial point



$\mathbf{x}_0 = (2, 0.3)$ , fast convergence



$\mathbf{x}_0 = (2, 0.02)$ , slow convergence



# Convergence Analysis

**Theorem.** If  $f$  is  $m$ -strongly convex and  $L$ -smooth, and  $\mathbf{x}^*$  is a minimum of  $f$ , then the sequence  $\{\mathbf{x}_k\}$  produced by gradient descent with exact line search satisfies

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \left(1 - \frac{m}{L}\right)^k [f(\mathbf{x}_0) - f(\mathbf{x}^*)]$$

## Notes.

- $0 \leq 1 - \frac{m}{L} < 1$ , so  $\mathbf{x}_k \rightarrow \mathbf{x}^*$  and  $f(\mathbf{x}_k) \rightarrow f(\mathbf{x}^*)$  exponentially fast
- The number of iterations to reach  $f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \epsilon$  is  $O(\log \frac{1}{\epsilon})$ . For  $\epsilon = 10^{-p}$ ,  $k = O(p)$ , linear in the number of significant digits.
- The convergence rate depends on the condition number  $L/m$  and can be slow if  $L/m$  is large. When close to  $\mathbf{x}^*$ , we can estimate  $L/m$  by  $\kappa(\nabla f^2(\mathbf{x}^*))$ .

## Proof

1. By the quadratic upper bound for  $L$ -smooth functions,

$$f(\mathbf{x}_k - t\nabla f(\mathbf{x}_k)) \leq f(\mathbf{x}_k) - t\|\nabla f(\mathbf{x}_k)\|^2 + \frac{Lt^2}{2}\|\nabla f(\mathbf{x}_k)\|^2 \triangleq q(t)$$

2. Minimizing over  $t$  in step 1,

$$f(\mathbf{x}_{k+1}) = \min_t f(\mathbf{x}_k - t\nabla f(\mathbf{x}_k)) \leq \min_t q(t) = q\left(\frac{1}{L}\right) = f(\mathbf{x}_k) - \frac{1}{2L}\|\nabla f(\mathbf{x}_k)\|^2$$

3. By  $m$ -strong convexity,

$$f(\mathbf{x}) \geq f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T(\mathbf{x} - \mathbf{x}_k) + \frac{m}{2}\|\mathbf{x} - \mathbf{x}_k\|^2 \triangleq \hat{f}(\mathbf{x})$$

4. Minimizing over  $\mathbf{x}$  in step 3,

$$f(\mathbf{x}^*) = \min_{\mathbf{x}} f(\mathbf{x}) \geq \min_{\mathbf{x}} \hat{f}(\mathbf{x}) = \hat{f}\left(\mathbf{x}_k - \frac{1}{m}\nabla f(\mathbf{x}_k)\right) = f(\mathbf{x}_k) - \frac{1}{2m}\|\nabla f(\mathbf{x}_k)\|^2$$

5. By 4,  $\|\nabla f(\mathbf{x}_k)\|^2 \geq 2m[f(\mathbf{x}_k) - f(\mathbf{x}^*)]$ . Plugging into 2,

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq \left(1 - \frac{m}{L}\right) [f(\mathbf{x}_k) - f(\mathbf{x}^*)]$$

# Backtracking Line Search

Exact line search is often expensive and not worth it. Suffices to find a good enough step size. One way to do so is to use **backtracking line search**, aka **Armijo's rule**.

## Gradient descent with backtracking line search

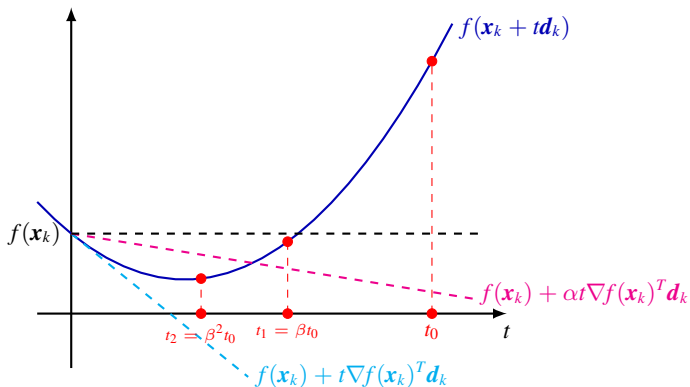
```
1: initialization  $\mathbf{x} \leftarrow \mathbf{x}_0 \in \mathbb{R}^n$ 
2: while  $\|\nabla f(\mathbf{x})\| > \delta$  do
3:    $t \leftarrow t_0$ 
4:   while  $f(\mathbf{x} - t\nabla f(\mathbf{x})) > f(\mathbf{x}) - \alpha t \|\nabla f(\mathbf{x})\|_2^2$  do
5:      $t \leftarrow \beta t$ 
6:   end while
7:    $\mathbf{x} \leftarrow \mathbf{x} - t\nabla f(\mathbf{x})$ 
8: end while
9: return  $\mathbf{x}$ 
```

$\alpha \in (0, 1)$  and  $\beta \in (0, 1)$  are constants. Armijo used  $\alpha = \beta = 0.5$

Values suggested in [BV]:  $\alpha \in [0.01, 0.3]$ ,  $\beta \in [0.1, 0.8]$

**Note.** For general  $\mathbf{d}$ , use condition  $f(\mathbf{x} + t\mathbf{d}) > f(\mathbf{x}) + \alpha t \nabla f(\mathbf{x})^T \mathbf{d}$

# Backtracking Line Search (cont'd)



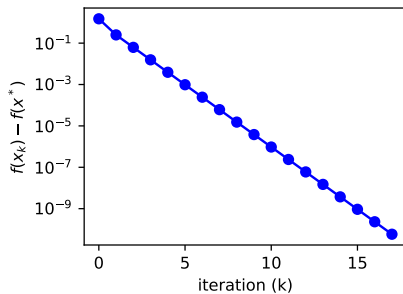
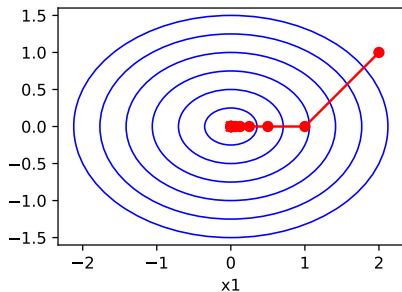
- $\nabla f(\mathbf{x}_k)^T \mathbf{d}_k < 0$  for descent direction  $\mathbf{d}_k$
- start from some “large” step size  $t_0$  ([BV] uses  $t_0 = 1$ )
- reduce step size geometrically until decrease is “large enough”

$$\underbrace{f(\mathbf{x}_k) - f(\mathbf{x}_k + t\mathbf{d}_k)}_{\text{actual decrease in function value}} \geq \alpha \times \underbrace{t |\nabla f(\mathbf{x}_k)^T \mathbf{d}_k|}_{\text{decrease along tangent line}}$$

## Example

$$f(x_1, x_2) = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} = \frac{\gamma}{2} x_1^2 + \frac{1}{2} x_2^2, \quad \mathbf{Q} = \text{diag}\{\gamma, 1\}$$

Well-conditioned.  $\gamma = 0.5$ ,  $\mathbf{x}_0 = (2, 1)^T$

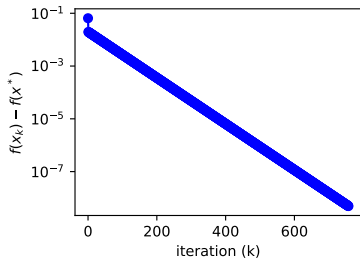
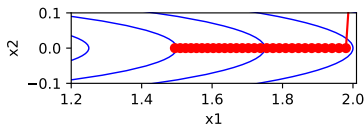


Fast convergence.

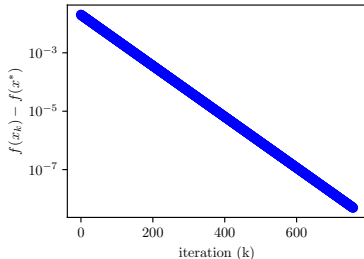
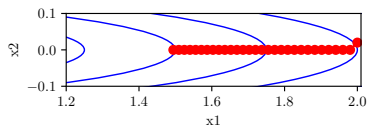
## Example (cont'd)

$$f(x_1, x_2) = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} = \frac{\gamma}{2} x_1^2 + \frac{1}{2} x_2^2, \quad \mathbf{Q} = \text{diag}\{\gamma, 1\}$$

Ill-conditioned.  $\gamma = 0.01$



$\mathbf{x}_0 = (2, 0.3)$ , slow convergence



$\mathbf{x}_0 = (2, 0.02)$ , slow convergence

# Convergence Analysis

**Theorem.** If  $f$  is  $m$ -strongly convex and  $L$ -smooth, and  $\mathbf{x}^*$  is a minimum of  $f$ , then the sequence  $\{\mathbf{x}_k\}$  produced by gradient descent with backtracking line search satisfies

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq c^k [f(\mathbf{x}_0) - f(\mathbf{x}^*)]$$

where

$$c = 1 - \min \left\{ 2m\alpha t_0, \frac{4m\beta\alpha(1-\alpha)}{L} \right\}$$

**Notes.**

- $c \in (0, 1)$ , as

$$\frac{4m\beta\alpha(1-\alpha)}{L} \leq \frac{\beta m}{L} \leq \beta < 1$$

so  $\mathbf{x}_k \rightarrow \mathbf{x}^*$  and  $f(\mathbf{x}_k) \rightarrow f(\mathbf{x}^*)$  exponentially fast

- Number of iterations to reach  $f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \epsilon$  is  $O(\log \frac{1}{\epsilon})$ . For  $\epsilon = 10^{-p}$ ,  $k = O(p)$ , linear in the number of significant digits.

## Proof

The inner loop terminates with a step size bounded from below.

1. By the quadratic upper bound for  $L$ -smooth functions,

$$f(\mathbf{x}_k - t\nabla f(\mathbf{x}_k)) \leq f(\mathbf{x}_k) - t(1 - \frac{Lt}{2})\|\nabla f(\mathbf{x}_k)\|^2$$

2. The inner loop terminates for sure if

$$-t(1 - \frac{Lt}{2})\|\nabla f(\mathbf{x}_k)\|^2 \leq -\alpha t\|\nabla f(\mathbf{x}_k)\|^2 \implies t \leq \frac{2(1 - \alpha)}{L}$$

3. The step size in backtracking line search satisfies

$$t_k \geq \eta \triangleq \min \left\{ t_0, \frac{2\beta(1 - \alpha)}{L} \right\}$$

►  $t_k = t_0$  if Armijo's condition is satisfied by  $t_0$

otherwise,  $\frac{t_k}{\beta} > \frac{2(1-\alpha)}{L}$ , since the inner loop did not terminate at  $\frac{t_k}{\beta}$





## Proof (cont'd)

Now we look at the outer loop

4. By Armijo's condition in the inner loop,

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

5. By 3 and 4,

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq f(\mathbf{x}_k) - f(\mathbf{x}^*) - \alpha \eta \|\nabla f(\mathbf{x}_k)\|^2$$

6. By step 4 of slide 9,

$$\|\nabla f(\mathbf{x}_k)\|^2 \geq 2m[f(\mathbf{x}_k) - f(\mathbf{x}^*)]$$

7. By 5 and 6,

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq (1 - 2m\alpha\eta)[f(\mathbf{x}_k) - f(\mathbf{x}^*)] = c[f(\mathbf{x}_k) - f(\mathbf{x}^*)]$$

so

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq c^k[f(\mathbf{x}_0) - f(\mathbf{x}^*)]$$

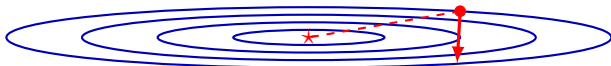
## Better Descent Direction

Gradient descent uses first-order information (i.e. gradient),

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

Locally  $-\nabla f(\mathbf{x}_k)$  is the max-rate descending direction, but globally it may not be the “right” direction.

**Example.** For  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{Q}\mathbf{x}$  with  $\mathbf{Q} = \text{diag}\{0.01, 1\}$ , optimum is  $\mathbf{x}^* = \mathbf{0}$ .



The negative gradient is

$$-\nabla f(\mathbf{x}) = -\mathbf{Q}\mathbf{x} = -(0.01x_1, x_2)^T$$

quite different from the “right” descent direction  $\mathbf{d} = -\mathbf{x}$ . Note

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

With second-order information (i.e. Hessian), we hope to do better.

# Newton's Method

By second-order Taylor expansion,

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

Minimizing quadratic approximation  $\hat{f}$ ,

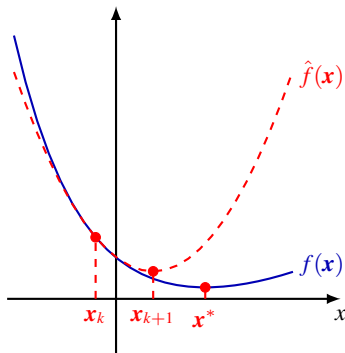
$$\nabla \hat{f}(\mathbf{x}) = \nabla^2 f(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k) + \nabla f(\mathbf{x}_k) = \mathbf{0}$$

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

provided  $\nabla^2 f(\mathbf{x}_k) \succ \mathbf{0}$ .

Newton step

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



**Note.** If  $f$  is quadratic, then  $f = \hat{f}$ , and Newton's method gets to the optimum in a single step starting from any  $\mathbf{x}_0$ .

## Newton's Method (cont'd)

- 1: initialization  $\mathbf{x} \leftarrow \mathbf{x}_0 \in \mathbb{R}^n$
- 2: **while**  $\|\nabla f(\mathbf{x})\| > \delta$  **do**
- 3:      $\mathbf{x} \leftarrow \mathbf{x} - [\nabla^2 f(\mathbf{x})]^{-1} \nabla f(\mathbf{x})$
- 4: **end while**
- 5: **return**  $\mathbf{x}$

**Note.** As in the case of gradient descent, other stopping criteria can be used. [BV] uses  $\nabla f(\mathbf{x})[\nabla^2 f(\mathbf{x})]^{-1} \nabla f(\mathbf{x}) > \delta$ .

The Newton step is a special case of  $\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k$  with

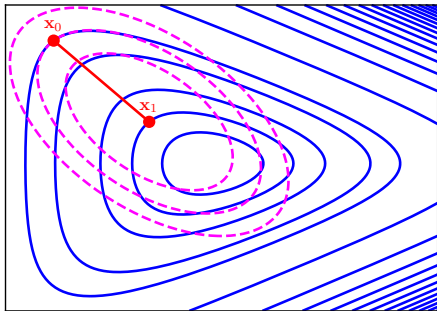
- **Newton direction**  $\mathbf{d}_k = -[\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$
- constant step size  $t_k = 1$

For  $\nabla^2 f(\mathbf{x}_k) \succ \mathbf{O}$ , the Newton direction is a descent direction

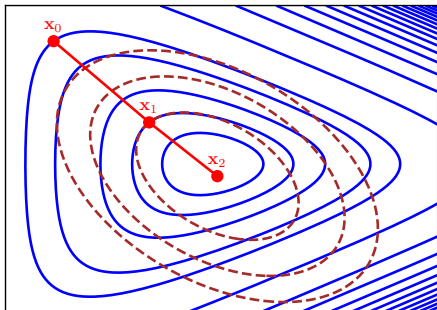
$$\nabla f(\mathbf{x}_k)^T \mathbf{d}_k = -\nabla f(\mathbf{x}_k)^T [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k) < 0 \quad \text{if } \nabla f(\mathbf{x}_k) \neq \mathbf{0}$$

## Newton's Method (cont'd)

The magenta curves are the level curves of the quadratic approximation of  $f$  at  $x_0$



The brown curves are the level curves of the quadratic approximation of  $f$  at  $x_1$ .



## Example

$$f(x_1, x_2) = e^{x_1+3x_2-0.1} + e^{x_1-3x_2-0.1} + e^{-x_1-0.1}$$

Newton step at  $\mathbf{x}_0 = (-2, 1)^T$ .

- gradient

$$\nabla f(\mathbf{x}_0) = e^{-0.1} \begin{pmatrix} e^{x_1+3x_2} + e^{x_1-3x_2} - e^{-x_1} \\ 3e^{x_1+3x_2} - 3e^{x_1-3x_2} \end{pmatrix} \Big|_{\mathbf{x}=\mathbf{x}_0} = \begin{pmatrix} -4.22019458 \\ 7.36051909 \end{pmatrix}$$

- Hessian

$$\begin{aligned} \nabla^2 f(\mathbf{x}_0) &= e^{-0.1} \begin{pmatrix} e^{x_1+3x_2} + e^{x_1-3x_2} + e^{-x_1} & 3e^{x_1+3x_2} - 3e^{x_1-3x_2} \\ 3e^{x_1+3x_2} - 3e^{x_1-3x_2} & 9e^{x_1+3x_2} + 9e^{x_1-3x_2} \end{pmatrix} \Big|_{\mathbf{x}=\mathbf{x}_0} \\ &= \begin{pmatrix} 9.1515943 & 7.36051909 \\ 7.36051909 & 22.19129872 \end{pmatrix} \end{aligned}$$

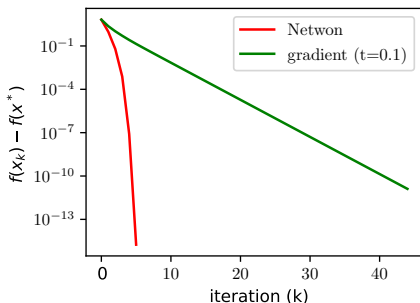
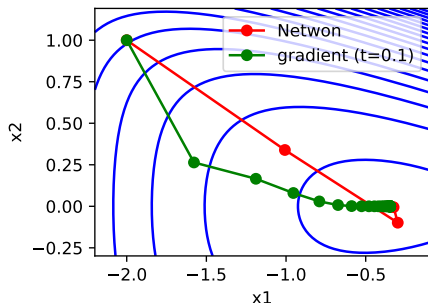
- Newton step

$$\mathbf{x}_1 = \mathbf{x}_0 - [\nabla^2 f(\mathbf{x}_0)]^{-1} \nabla f(\mathbf{x}_0) = \begin{pmatrix} -1.00725064 \\ 0.33903509 \end{pmatrix}$$

## Example (cont'd)

$$f(x_1, x_2) = e^{x_1+3x_2-0.1} + e^{x_1-3x_2-0.1} + e^{-x_1-0.1}$$

Solution using Newton's method and gradient descent with constant step size 0.1. Initial point  $\mathbf{x}_0 = (-2, 1)^T$ .



- Newton's method takes a more “direct” path
- Newton's method requires much fewer iterations, but each iteration is more expensive

## Connection to Root Finding

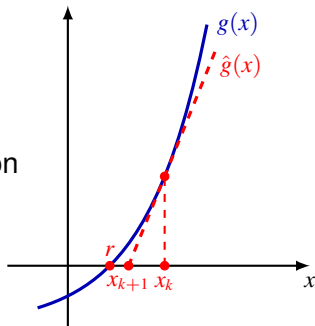
Newton's method is originally an algorithm for solving  $g(x) = 0$ .

By the first-order Taylor expansion,

$$g(x) \approx \hat{g}(x) \triangleq g(x_k) + g'(x_k)(x - x_k)$$

Use the root of  $\hat{g}(x)$  as the next approximation

$$x_{k+1} = x_k - \frac{g(x_k)}{g'(x_k)}$$



**Example (computing  $\sqrt{C}$ ).**  $\sqrt{C}$  is a root of  $g(x) = x^2 - C$ . Newton's method yields

$$x_{k+1} = x_k - \frac{x_k^2 - C}{2x_k} = \frac{1}{2} \left( x_k + \frac{C}{x_k} \right)$$

For  $x_0 > 0$ ,  $x_k$  converges to  $\sqrt{C}$ .



## Connection to Root Finding (cont'd)

Back to the optimization problem,

$$\min_x f(x)$$

The optimal solution  $x^*$  satisfies

$$f'(x^*) = 0$$

Letting  $g = f'$  in Newton's root finding algorithm,

$$x_{k+1} = x_k - \frac{f'(x_k)}{f''(x_k)} = x_k - [f''(x_k)]^{-1}f'(x_k)$$

In  $n$ -dimension,  $f' \rightarrow \nabla f$ ,  $f'' \rightarrow \nabla^2 f$ . We want to solve

$$\nabla f(\mathbf{x}^*) = \mathbf{0}$$

Newton's algorithm becomes

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