

## Convergence of Steepest Descent Method: Quadratic case

Consider the problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} - \mathbf{c}^T \mathbf{x}$$

where  $\mathbf{H}$  is a symmetric positive-definite matrix.

- $\mathbf{g}(\mathbf{x}) = \mathbf{H}\mathbf{x} - \mathbf{c}$ .  $\therefore \mathbf{x}^* = \mathbf{H}^{-1}\mathbf{c}$ .
- How does steepest descent method perform, when applied to  $f(\mathbf{x})$ ?
- Assume that *exact line search* is used in each iteration

**What is the step length  $\alpha^k$  at iteration  $k$ ?**

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{H} \mathbf{x} - \mathbf{c}^T \mathbf{x}. \therefore \mathbf{g}^k = \mathbf{g}(\mathbf{x}^k) = \mathbf{H} \mathbf{x}^k - \mathbf{c}$$

Define  $\phi(\alpha) = f(\mathbf{x}^k + \alpha \mathbf{d}^k) = f(\mathbf{x}^k - \alpha \mathbf{g}^k)$ .

**Exact line search:**

$$\begin{aligned}\alpha^k &= \arg \min_{\alpha > 0} \phi(\alpha) \\ \phi'(\alpha) = 0 &\Rightarrow \nabla f(\mathbf{x}^k - \alpha \mathbf{g}^k)^T (-\mathbf{g}^k) = 0 \\ &\Rightarrow (\mathbf{H} \mathbf{x}^k - \alpha \mathbf{H} \mathbf{g}^k - \mathbf{c})^T \mathbf{g}^k = 0 \\ &\Rightarrow (\mathbf{g}^k - \alpha \mathbf{H} \mathbf{g}^k)^T \mathbf{g}^k = 0\end{aligned}$$

Therefore,

$$\begin{aligned}\alpha^k &= \frac{\mathbf{g}^{kT} \mathbf{g}^k}{\mathbf{g}^{kT} \mathbf{H} \mathbf{g}^k} \\ \therefore \mathbf{x}^{k+1} &= \mathbf{x}^k - \left( \frac{\mathbf{g}^{kT} \mathbf{g}^k}{\mathbf{g}^{kT} \mathbf{H} \mathbf{g}^k} \right) \mathbf{g}^k\end{aligned}$$

**At what rate does  $\{\mathbf{x}^k\}$  converge?**

Define

$$E(\mathbf{x}^k) = \frac{1}{2}(\mathbf{x}^k - \mathbf{x}^*)^T \mathbf{H}(\mathbf{x}^k - \mathbf{x}^*). \quad (E(\mathbf{x}^k) > 0, \text{ if } \mathbf{x}^k \neq \mathbf{x}^*)$$

Note that  $E(\mathbf{x}^k) = f(\mathbf{x}^k) + \underbrace{\frac{1}{2}\mathbf{x}^{*T} \mathbf{H} \mathbf{x}^*}_{\text{constant}}.$

Define  $\mathbf{y}^k = \mathbf{x}^k - \mathbf{x}^*. \therefore \mathbf{H}\mathbf{y}^k = \mathbf{g}^k.$

Using

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \left( \frac{\mathbf{g}^{kT} \mathbf{g}^k}{\mathbf{g}^{kT} \mathbf{H} \mathbf{g}^k} \right) \mathbf{g}^k,$$

Relative decrease in  $E$ ,

$$\begin{aligned} & \frac{E(\mathbf{x}^k) - E(\mathbf{x}^{k+1})}{E(\mathbf{x}^k)} \\ = & \frac{(\mathbf{x}^k - \mathbf{x}^*)^T \mathbf{H}(\mathbf{x}^k - \mathbf{x}^*) - (\mathbf{x}^{k+1} - \mathbf{x}^*)^T \mathbf{H}(\mathbf{x}^{k+1} - \mathbf{x}^*)}{\mathbf{y}^{kT} \mathbf{H} \mathbf{y}^k} \end{aligned}$$

$$\begin{aligned}
& \frac{E(\mathbf{x}^k) - E(\mathbf{x}^{k+1})}{E(\mathbf{x}^k)} \\
= & \frac{(\mathbf{x}^k - \mathbf{x}^*)^T \mathbf{H}(\mathbf{x}^k - \mathbf{x}^*) - (\mathbf{x}^{k+1} - \mathbf{x}^*)^T \mathbf{H}(\mathbf{x}^{k+1} - \mathbf{x}^*)}{\mathbf{y}^k{}^T \mathbf{H} \mathbf{y}^k} \\
= & \frac{2\alpha^k \mathbf{g}^k{}^T \mathbf{g}^k - \alpha^{k2} \mathbf{g}^k{}^T \mathbf{H} \mathbf{g}^k}{\mathbf{y}^k{}^T \mathbf{H} \mathbf{y}^k}
\end{aligned}$$

Substituting  $\alpha^k = \frac{\mathbf{g}^k{}^T \mathbf{g}^k}{\mathbf{g}^k{}^T \mathbf{H} \mathbf{g}^k}$ , we get

$$\frac{E(\mathbf{x}^k) - E(\mathbf{x}^{k+1})}{E(\mathbf{x}^k)} = \frac{(\mathbf{g}^k{}^T \mathbf{g}^k)^2}{(\mathbf{g}^k{}^T \mathbf{H} \mathbf{g}^k)(\mathbf{g}^k{}^T \mathbf{H}^{-1} \mathbf{g}^k)}$$

## Kantorovich inequality

Let  $\mathbf{H} \in \mathbb{R}^{n \times n}$  be a symmetric positive definite matrix. Let  $\lambda_1$  and  $\lambda_n$  be respectively the smallest and largest eigenvalues of  $\mathbf{H}$ . Then, for any  $\mathbf{x} \neq \mathbf{0}$ ,

$$\frac{(\mathbf{x}^T \mathbf{x})^2}{(\mathbf{x}^T \mathbf{H} \mathbf{x})(\mathbf{x}^T \mathbf{H}^{-1} \mathbf{x})} \geq \frac{4\lambda_1 \lambda_n}{(\lambda_1 + \lambda_n)^2}$$

Using this inequality,

$$\begin{aligned} \frac{E(\mathbf{x}^k) - E(\mathbf{x}^{k+1})}{E(\mathbf{x}^k)} &= \frac{(\mathbf{g}^{kT} \mathbf{g}^k)^2}{(\mathbf{g}^{kT} \mathbf{H} \mathbf{g}^k)(\mathbf{g}^{kT} \mathbf{H}^{-1} \mathbf{g}^k)} \\ &\geq \frac{4\lambda_1 \lambda_n}{(\lambda_1 + \lambda_n)^2} \end{aligned}$$

Therefore,

$$E(\mathbf{x}^{k+1}) \leq \left( \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^2 E(\mathbf{x}^k)$$

$$E(\mathbf{x}^{k+1}) \leq \left( \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^2 E(\mathbf{x}^k)$$

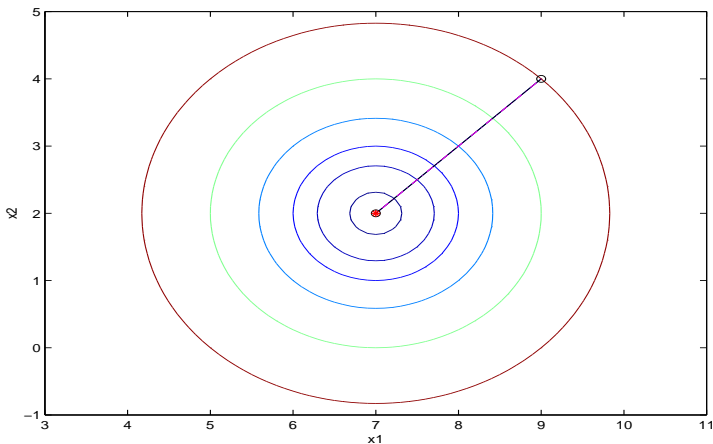
Therefore,  $E(\mathbf{x}^k) \rightarrow 0$  and  $\mathbf{x}^k \rightarrow \mathbf{x}^*$  ( $\mathbf{H}$  is positive definite).

With respect to  $E$ , the steepest descent method

- converges linearly with convergence rate no greater than  $\left( \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^2$
- Actual convergence rate depends upon  $\mathbf{x}^0$
- Define the *condition number* of  $\mathbf{H}$ ,  $r = \frac{\lambda_n}{\lambda_1}$
- Convergence rate of the steepest descent method depends on the condition number of  $\mathbf{H}$ 
  - $r = 1$  (circular contours)  $\Rightarrow$  convergence in one iteration
  - $r \gg 1$  (elliptical contours)  $\Rightarrow$  convergence is slow
- For nonquadratic functions, rate of convergence to  $\mathbf{x}^*$  depends on the condition number of  $\mathbf{H}(\mathbf{x}^*)$

Example:

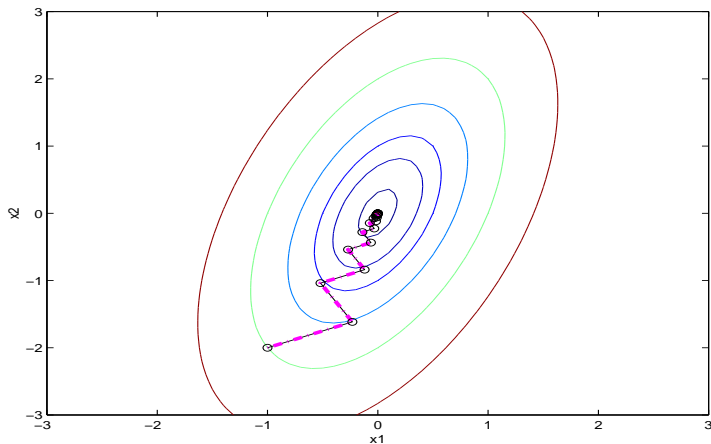
$$\min f(\mathbf{x}) \triangleq (x_1 - 7)^2 + (x_2 - 2)^2$$



Steepest descent algorithm (with exact line search) applied to  $f(\mathbf{x})$  converges in **one iteration** from **any starting point**

Example:

$$\min f(\mathbf{x}) \triangleq 4x_1^2 + x_2^2 - 2x_1x_2$$



Steepest descent algorithm (with exact line search) applied to  $f(\mathbf{x})$  requires many iterations before it converges