

Numerical Optimization

Unconstrained Optimization

Shirish Shevade

Computer Science and Automation
Indian Institute of Science
Bangalore 560 012, India.

NPTEL Course on Numerical Optimization

Steepest Descent Method

- Uses the steepest descent direction, $\mathbf{d}_{SD}^k = -\mathbf{g}^k$

Steepest Descent Algorithm

- (1) Initialize \mathbf{x}^0 and ϵ , set $k := 0$.
- (2) **while** $\|\mathbf{g}^k\| > \epsilon$
 - (a) $\mathbf{d}^k = -\mathbf{g}^k$
 - (b) Find $\alpha^k (> 0)$ along \mathbf{d}^k such that
 - (i) $f(\mathbf{x}^k + \alpha^k \mathbf{d}^k) < f(\mathbf{x}^k)$
 - (ii) α^k satisfies Armijo-Wolfe conditions
 - (c) $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha^k \mathbf{d}^k$
 - (d) $k := k + 1$

endwhile

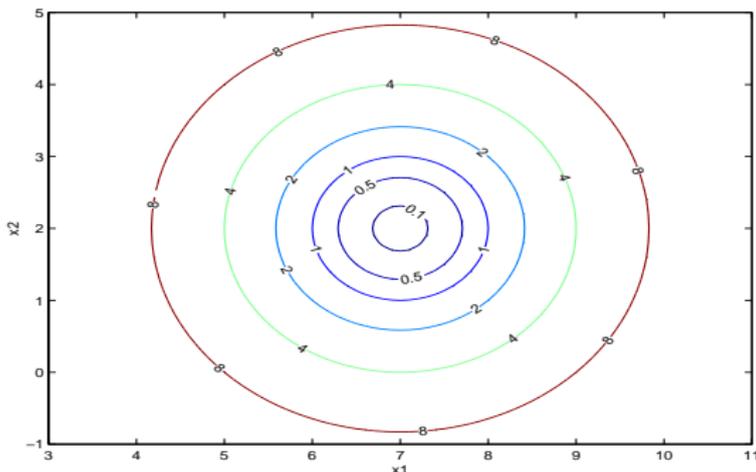
Output : $\mathbf{x}^* = \mathbf{x}^k$, a stationary point of $f(\mathbf{x})$.

-
- Exact or Backtracking line search can be used in step 2(b)

Example:

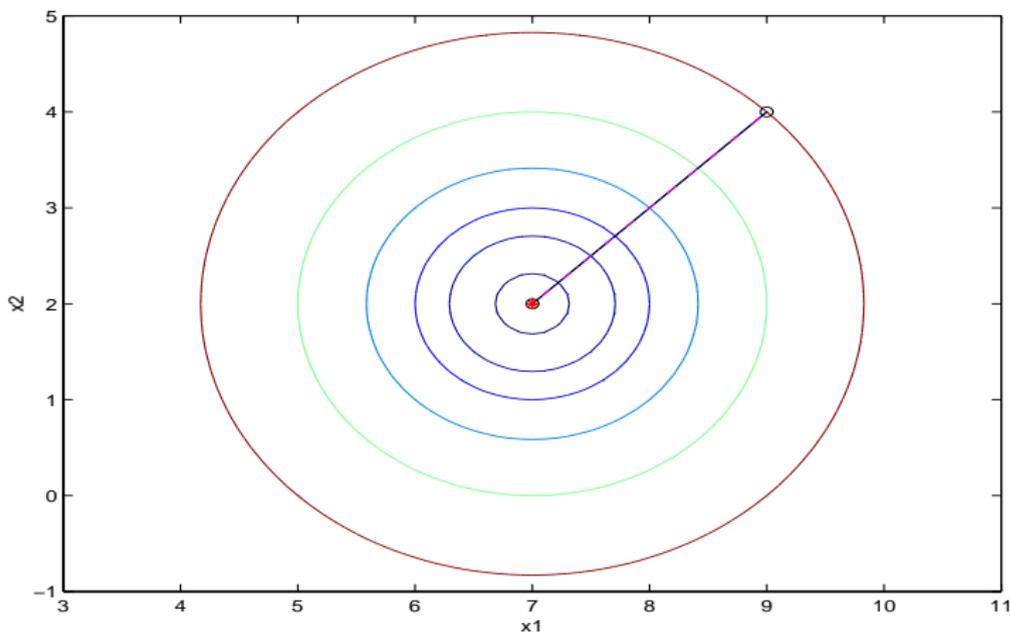
$$\min f(\mathbf{x}) \triangleq (x_1 - 7)^2 + (x_2 - 2)^2$$

- $\mathbf{g}(\mathbf{x}) = \begin{pmatrix} 2(x_1 - 7) \\ 2(x_2 - 2) \end{pmatrix}$, $\mathbf{H}(\mathbf{x}) = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$.
- $\mathbf{x}^* = \begin{pmatrix} 7 \\ 2 \end{pmatrix}$



Example:

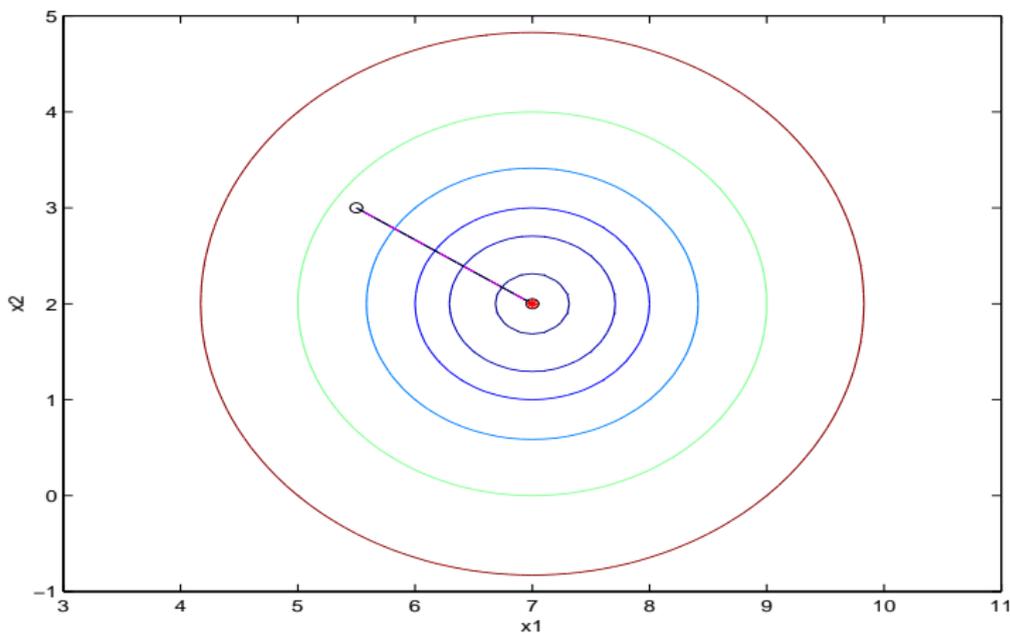
$$\min f(\mathbf{x}) \triangleq (x_1 - 7)^2 + (x_2 - 2)^2$$



Behaviour of the steepest descent algorithm (with exact line search) applied to $f(\mathbf{x})$ using $\mathbf{x}^0 = (9, 4)^T$

Example:

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

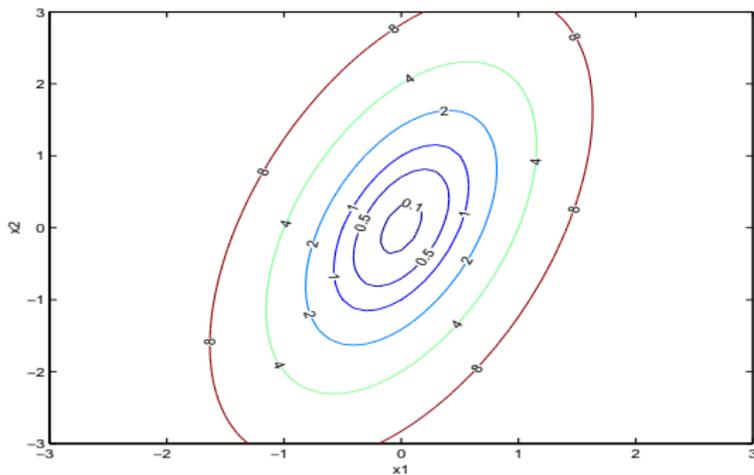


Behaviour of the steepest descent algorithm (with exact line search) applied to $f(\mathbf{x})$ using $\mathbf{x}^0 = (5.5, 3)^T$

Example:

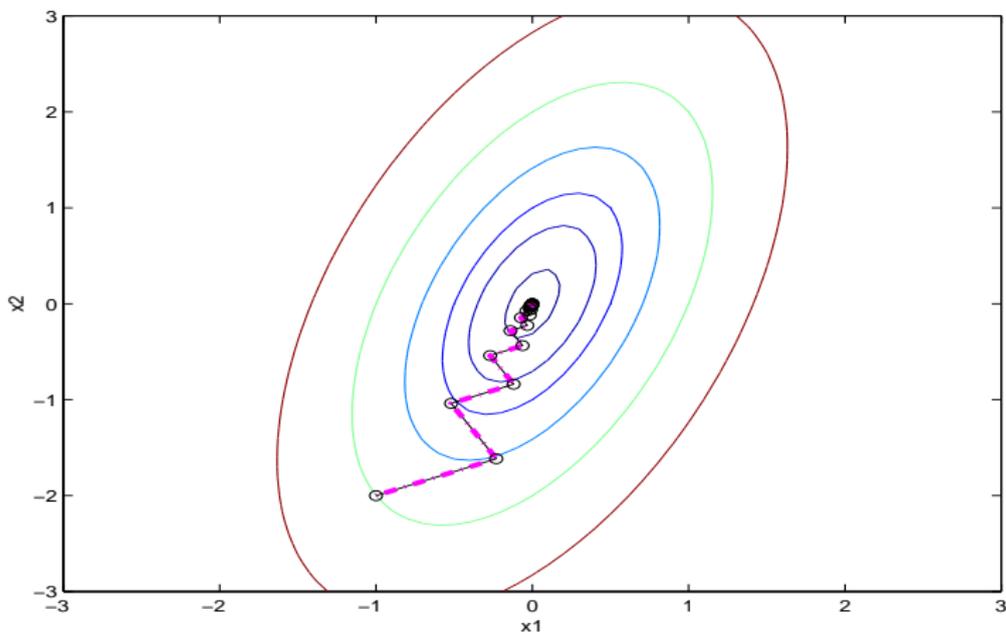
$$\min f(\mathbf{x}) \triangleq 4x_1^2 + x_2^2 - 2x_1x_2$$

- $\mathbf{g}(\mathbf{x}) = \begin{pmatrix} 8x_1 - 2x_2 \\ 2x_2 - 2x_1 \end{pmatrix}$, $\mathbf{H}(\mathbf{x}) = \begin{pmatrix} 8 & -2 \\ -2 & 2 \end{pmatrix}$.
- $\mathbf{x}^* = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$



Example:

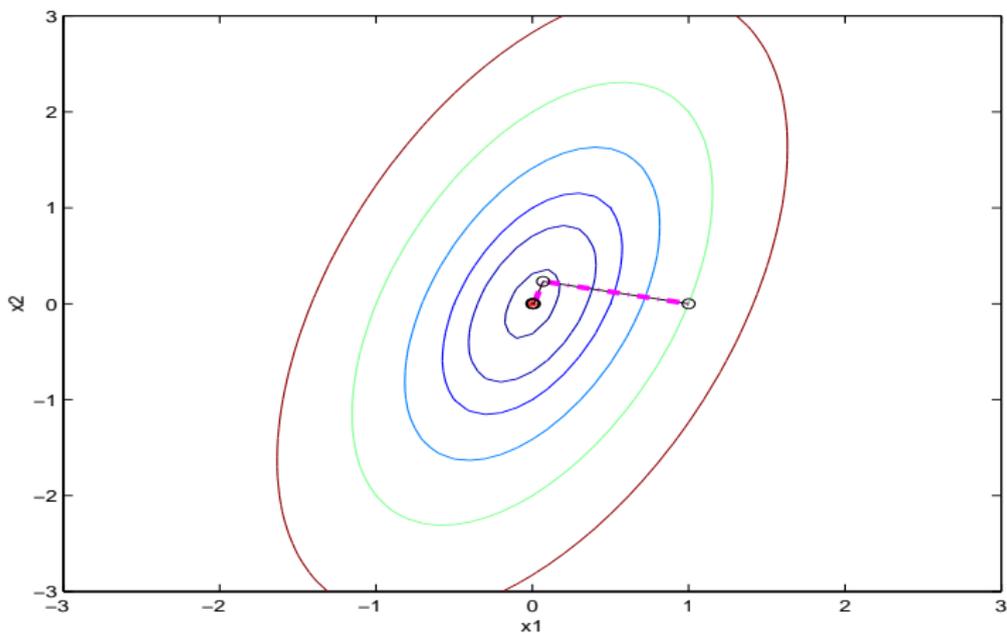
$$\min f(\mathbf{x}) \triangleq 4x_1^2 + x_2^2 - 2x_1x_2$$



Behaviour of the steepest descent algorithm (with exact line search) applied to $f(\mathbf{x})$ using $\mathbf{x}^0 = (-1, -2)^T$

Example:

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

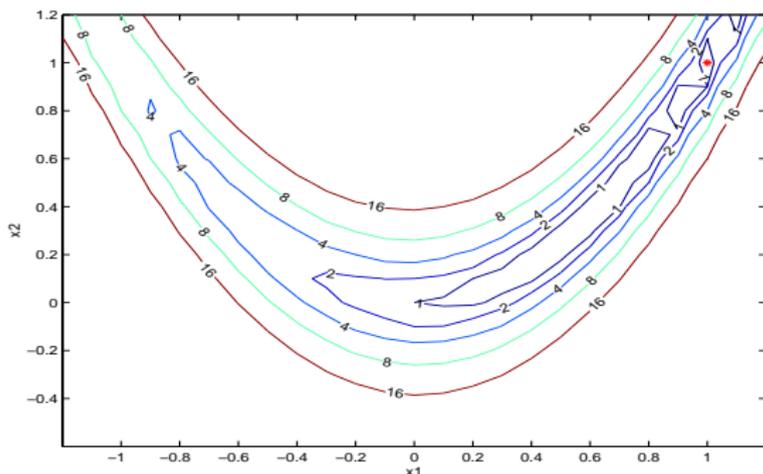


Behaviour of the steepest descent algorithm (with exact line search) applied to $f(\mathbf{x})$ using $\mathbf{x}^0 = (1, 0)^T$

Example (Rosenbrock function):

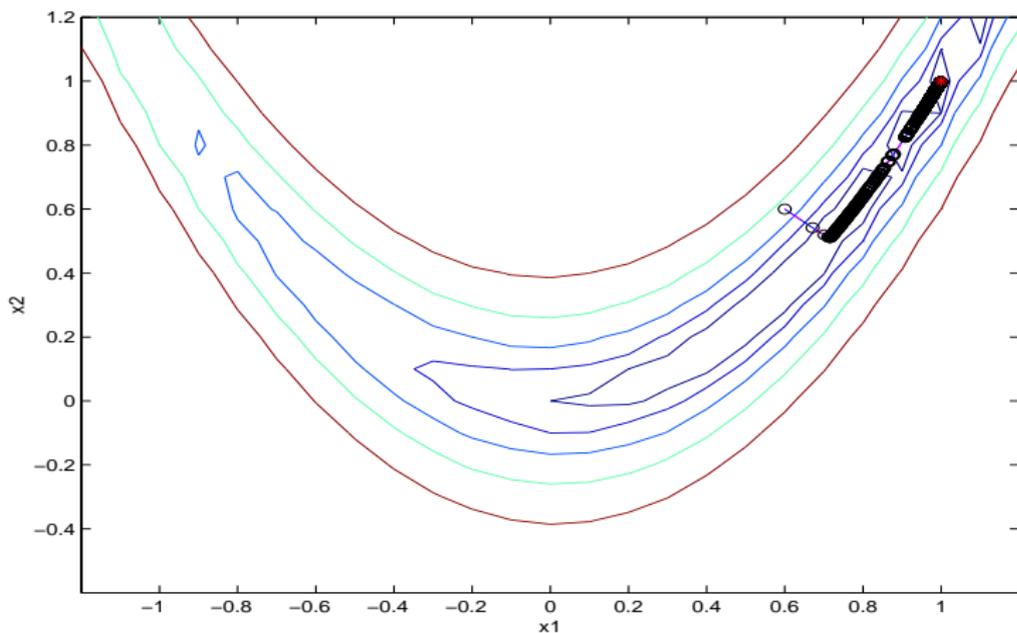
$$\min f(\mathbf{x}) \triangleq 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

• $\mathbf{x}^* = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, f(\mathbf{x}^*) = 0$



Example (Rosenbrock function):

$$\min f(\mathbf{x}) \triangleq 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$



Behaviour of the steepest descent algorithm (with backtracking line search) applied to $f(\mathbf{x})$ using $\mathbf{x}^0 = (0.6, 0.6)^T$

Example (Rosenbrock function):

$$\min f(\mathbf{x}) \triangleq 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

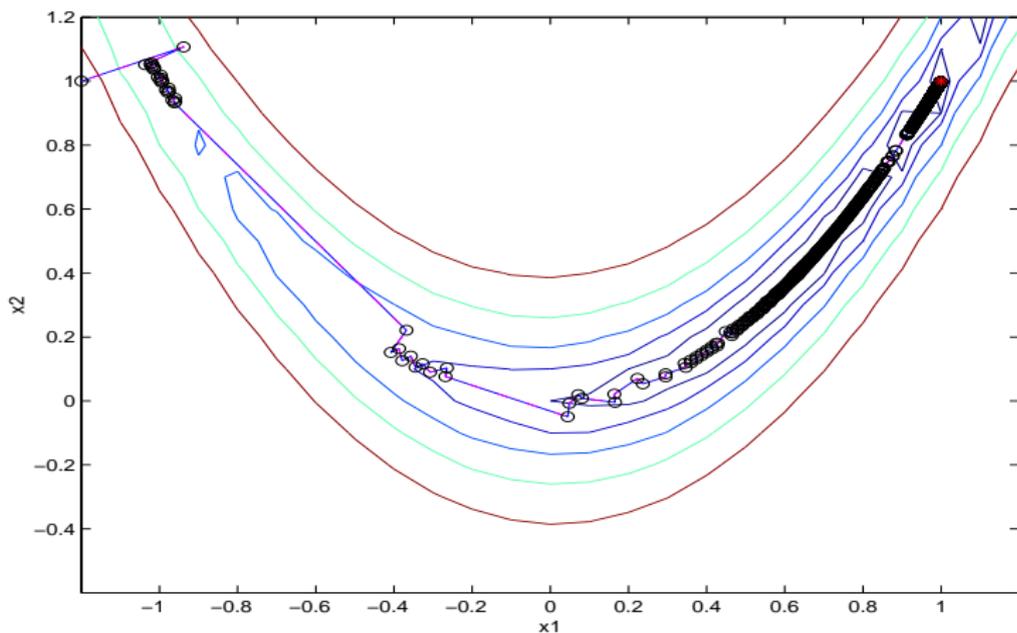
• $\mathbf{x}^* = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, f(\mathbf{x}^*) = 0$

k	\mathbf{x}_1^k	\mathbf{x}_2^k	$f(\mathbf{x}^k)$	$\ \mathbf{x}^k - \mathbf{x}^*\ $	$\ \mathbf{g}^k\ $
0	0.6	0.6	5.92	0.5657	75.59
10	0.72	0.52	0.0792	0.5601	0.3938
100	0.78	0.61	0.0465	0.4414	0.2451
1000	0.9914	0.9828	7.45×10^{-5}	0.0192	0.0069
2028	0.9989	0.9978	1.81×10^{-6}	0.0024	9.97×10^{-4}

Table: Steepest descent method (with backtracking line search) applied to Rosenbrock function, using $\mathbf{x}^0 = (0.6, 0.6)^T, \hat{\alpha} = .5, \rho = .3$ and $c_1 = 1.0 \times 10^{-4}$.

Example (Rosenbrock function):

$$\min f(\mathbf{x}) \triangleq 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$



Behaviour of the steepest descent algorithm (with backtracking line search) applied to $f(\mathbf{x})$ using $\mathbf{x}^0 = (-1.2, 1)^T$

Example (Rosenbrock function):

$$\min f(\mathbf{x}) \triangleq 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

• $\mathbf{x}^* = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, f(\mathbf{x}^*) = 0$

k	\mathbf{x}_1^k	\mathbf{x}_2^k	$f(\mathbf{x}^k)$	$\ \mathbf{x}^k - \mathbf{x}^*\ $	$\ \mathbf{g}^k\ $
0	-1.2	1.0	24.2	2.2	232.87
10	-1.00	1.01	4.02	2.0042	7.69
100	0.57	0.32	0.1867	0.80	0.84
1000	0.99	0.97	1.99×10^{-4}	0.0314	0.014
2300	0.9989	0.9979	1.11×10^{-6}	0.0024	9.63×10^{-4}

Table: Steepest descent method (with backtracking lines search) applied to Rosenbrock function, using $\mathbf{x}^0 = (-1.2, 1)^T, \hat{\alpha} = .5, \rho = .3$ and $c_1 = 1.0 \times 10^{-4}$.

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 α^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 ,

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

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

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

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

Kantorovich inequality

Let $\mathbf{H} \in \mathbb{R}^{n \times n}$ be a symmetric positive definite matrix. Let λ_1 and λ_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,

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

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

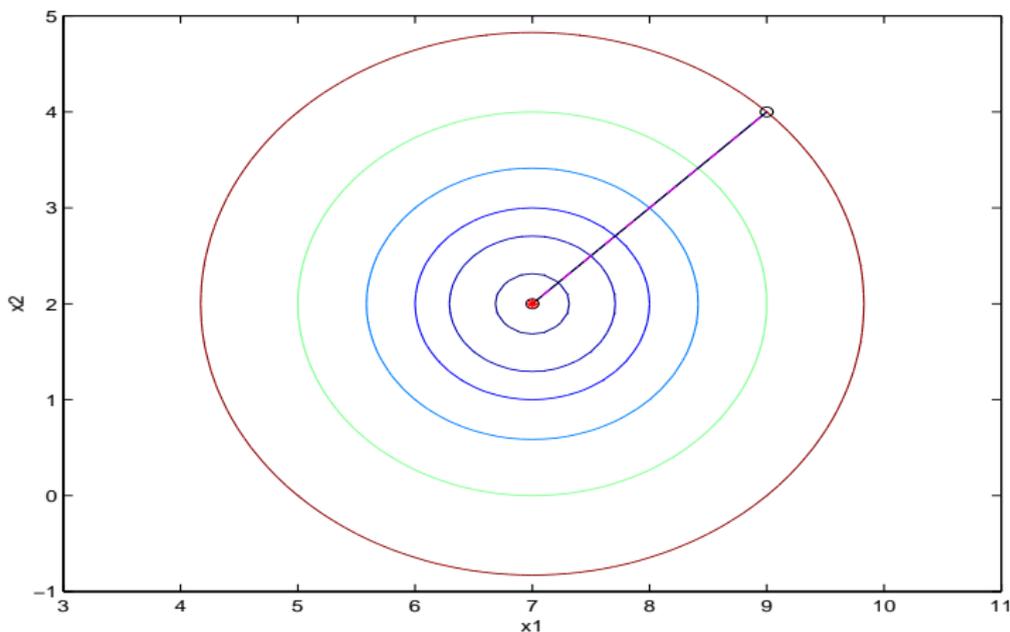
Therefore, $E(\mathbf{x}^k) \rightarrow 0$ and $\mathbf{x}^k \rightarrow \mathbf{x}^*$ (\mathbf{H} is positive definite).

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:

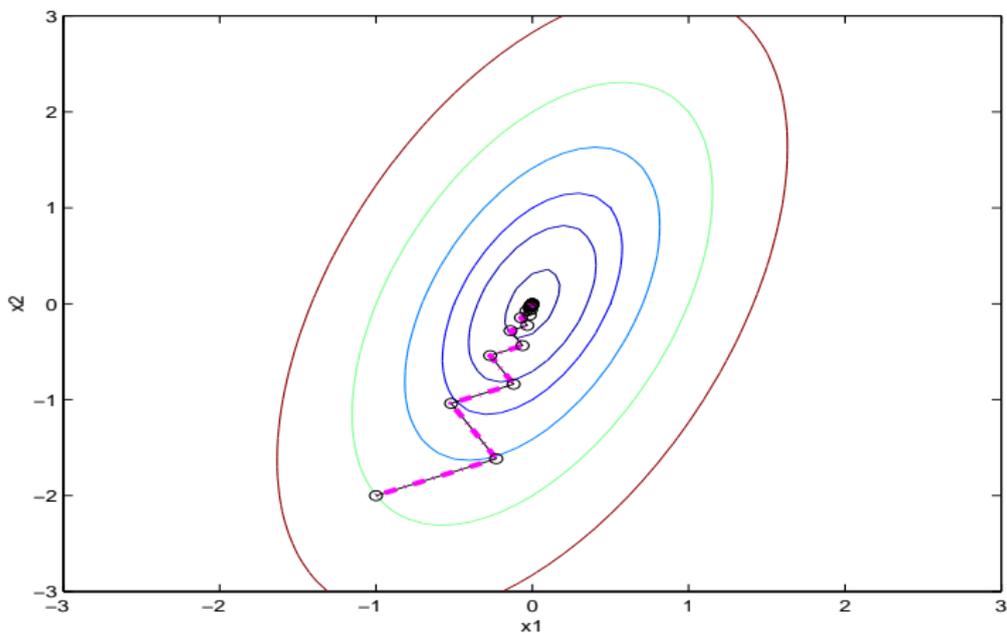
$$\min f(\mathbf{x}) \triangleq (x_1 - 7)^2 + (x_2 - 2)^2$$



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

Consider the problem to minimize

$$\min f(\mathbf{x}) \triangleq \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} - \mathbf{c}^T \mathbf{x}$$

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

- *Condition number* of the Hessian matrix controls the convergence rate of steepest descent method.
- Faster convergence if the Hessian matrix is \mathbf{I}
- Let $\mathbf{H} = \mathbf{L}\mathbf{L}^T$ be the Cholesky decomposition of \mathbf{H}
- Define $\mathbf{y} = \mathbf{L}^T \mathbf{x}$. Therefore, the function $f(\mathbf{x})$ is transformed to the function $h(\mathbf{y})$.

$$h(\mathbf{y}) \triangleq f(\mathbf{L}^{-T} \mathbf{y})$$

$$\begin{aligned}
h(\mathbf{y}) &= f(\mathbf{L}^{-T}\mathbf{y}) \\
&= \frac{1}{2}\mathbf{y}^T\mathbf{L}^{-1}\mathbf{H}\mathbf{L}^{-T}\mathbf{y} - \mathbf{c}^T\mathbf{L}^{-T}\mathbf{y} \\
&= \frac{1}{2}\mathbf{y}^T\mathbf{L}^{-1}\mathbf{L}\mathbf{L}^T\mathbf{L}^{-T}\mathbf{y} - \mathbf{c}^T\mathbf{L}^{-T}\mathbf{y} \\
&= \frac{1}{2}\mathbf{y}^T\mathbf{y} - \mathbf{c}^T\mathbf{L}^{-T}\mathbf{y}
\end{aligned}$$

- The Hessian matrix of $h(\mathbf{y})$ is \mathbf{I}
- Let us apply steepest descent method in \mathbf{y} -space

$$\begin{aligned}
\mathbf{y}^{k+1} &= \mathbf{y}^k - \nabla h(\mathbf{y}^k) \\
&= \mathbf{y}^k - \mathbf{L}^{-1}\nabla f(\mathbf{L}^{-T}\mathbf{y}^k) \\
\therefore \mathbf{L}^{-T}\mathbf{y}^{k+1} &= \mathbf{L}^{-T}\mathbf{y}^k - \mathbf{L}^{-T}\mathbf{L}^{-1}\nabla f(\mathbf{L}^{-T}\mathbf{y}^k) \\
\therefore \mathbf{x}^{k+1} &= \mathbf{x}^k - \mathbf{H}^{-1}\nabla f(\mathbf{x}^k)
\end{aligned}$$