

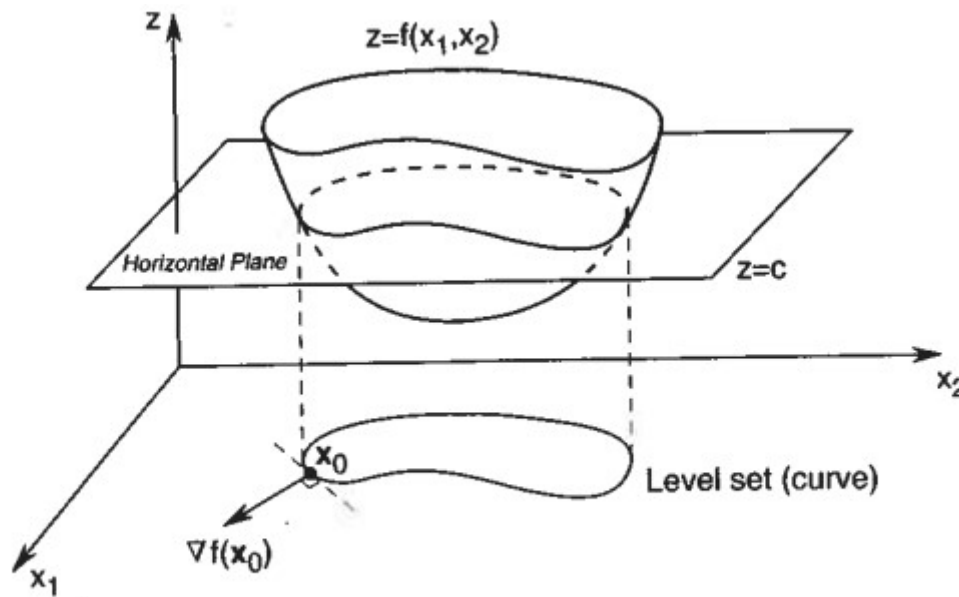
# Optimization – Gradient Methods

Henrik Karstoft

@

Aarhus University, Department of Engineering

# The gradient



**Figure 8.1** Constructing a level set corresponding to level  $c$  for  $f$ .

Gradient of  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is a column vector

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

First derivative—row vector

$$Df(\mathbf{x}) = \nabla f(\mathbf{x})^\top = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \cdots & \frac{\partial f}{\partial x_n} \end{bmatrix}$$

The directional derivative of  $f$  at  $\mathbf{x} \in \mathbb{R}^n$  in the direction  $\mathbf{d}$  is denoted

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{d}}$$

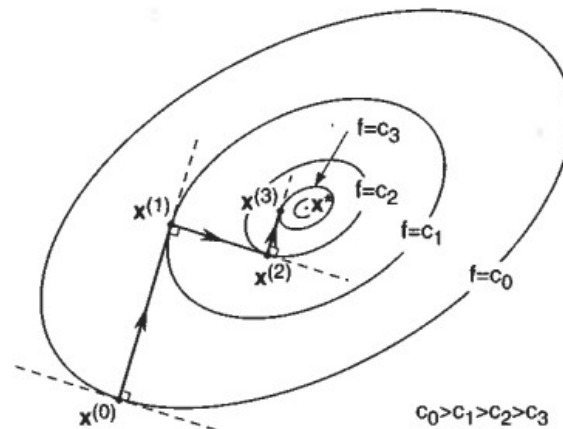
$$\begin{aligned} \frac{\partial f(\mathbf{x})}{\partial \mathbf{d}} &= \lim_{\alpha \rightarrow 0} \frac{f(\mathbf{x} + \alpha \mathbf{d}) - f(\mathbf{x})}{\alpha} \\ &= \left. \frac{d}{d\alpha} f(\mathbf{x} + \alpha \mathbf{d}) \right|_{\alpha=0} \\ &= \begin{bmatrix} Df(\mathbf{x}) \end{bmatrix} \begin{bmatrix} \mathbf{d} \end{bmatrix} \end{aligned}$$

# The Method of Steepest Descent

Finding best  $\alpha$ :

$$\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)})).$$

**Proposition 8.1** *If  $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$  is a steepest descent sequence for a given function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , then for each  $k$  the vector  $\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$  is orthogonal to the vector  $\mathbf{x}^{(k+2)} - \mathbf{x}^{(k+1)}$ .  $\square$*



**Figure 8.2** Typical sequence resulting from the method of steepest descent.

# The Method of Steepest Descent

**Example 2.2.** This is a 2-dimensional minimization example, illustrated on the front page. A tourist has lost his way in a hilly country. It is a foggy day so he cannot see far and he has no map. He knows that his rescue is at the bottom of a nearby valley. As tools he has a compass and his sense of balance, which can tell him about the local slope of the ground.

In order not to walk in circles he decides to use straight strides, i.e. with constant compass bearing. From what his feet tell him about the local slope he chooses a direction and walks in that direction until his feet tell him that he is on an uphill slope.

Now he has to decide on a new direction and he starts his next stride. Let us hope that he reaches his goal in the end. ■

**Example 2.3.** Let us return to our tourist who is lost in the fog in a hilly country. By experimenting with his compass he can find out that “half” the compass bearings give strides that start uphill and that the “other half” gives strides that start downhill. Between the two halves are two strides which start off going neither uphill or downhill. These form the tangent to the level curve corresponding to his position. ■

# The Method of Steepest Descent

**Proposition 8.2** *If  $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$  is the steepest descent sequence for  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  and if  $\nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$ , then  $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$ .  $\square$*

Stopping criteria:

$$\nabla f(\mathbf{x}^{(k+1)}) = \mathbf{0},$$

$$|f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)})| < \varepsilon,$$

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| < \varepsilon.$$

$$\frac{|f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)})|}{|f(\mathbf{x}^{(k)})|} < \varepsilon$$

$$\frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|}{\|\mathbf{x}^{(k)}\|} < \varepsilon.$$

# The Method of Steepest Descent

**Example 8.1** We use the method of steepest descent to find the minimizer of

$$f(x_1, x_2, x_3) = (x_1 - 4)^4 + (x_2 - 3)^2 + 4(x_3 + 5)^4.$$

The initial point is  $\mathbf{x}^{(0)} = [4, 2, -1]^T$ . We perform three iterations.

We find that

$$\nabla f(\mathbf{x}) = [4(x_1 - 4)^3, 2(x_2 - 3), 16(x_3 + 5)^3]^T.$$

Hence,

$$\nabla f(\mathbf{x}^{(0)}) = [0, -2, 1024]^T.$$

To compute  $\mathbf{x}^{(1)}$ , we need

$$\begin{aligned} \alpha_0 &= \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})) \\ &= \arg \min_{\alpha \geq 0} (0 + (2 + 2\alpha - 3)^2 + 4(-1 - 1024\alpha + 5)^4) \\ &= \arg \min_{\alpha \geq 0} \phi_0(\alpha). \end{aligned}$$

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \alpha_0 \nabla f(\mathbf{x}^{(0)}) = [4.000, 2.008, -5.062]^T.$$

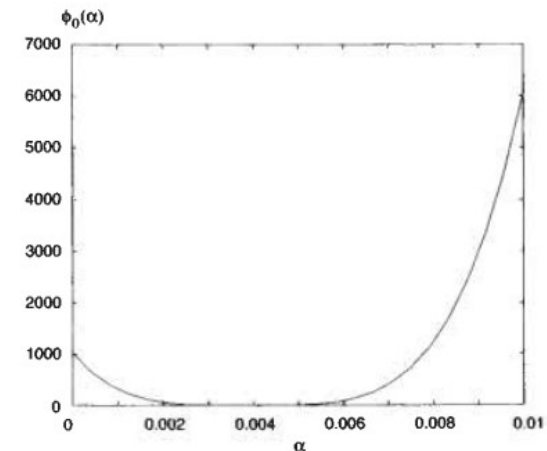


Figure 8.3 Plot of  $\phi_0(\alpha)$  versus  $\alpha$ .

# Review of eigenvalues quadratic forms

## Quadratic form

$$Q(x) = x^T A x$$

where  $x \in R^n$ ,  $A$  is symmetric  $n \times n$  matrix

**Eigenvalues** for a  $A$   $n \times n$  matrix are roots in characteristic the polynomial  $\chi(\lambda)$  of degree  $n$ :

$$\chi(\lambda) = \det(A - \lambda I)$$

where "det" is the determinat of  $A$ , and  $I$  is the identity matrix

# The Method of Steepest Descent, quadratic

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{b}^\top \mathbf{x},$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)},$$

$$\begin{aligned} \alpha_k &= \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)}) \\ &= \arg \min_{\alpha \geq 0} \left( \frac{1}{2} (\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)})^\top \mathbf{Q} (\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)}) - (\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)})^\top \mathbf{b} \right). \end{aligned}$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \frac{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}} \mathbf{g}^{(k)},$$

$$\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)}) = \mathbf{Q} \mathbf{x}^{(k)} - \mathbf{b}.$$



# Analysis of gradient methods, convergence

**Theorem 8.1** Let  $\{\mathbf{x}^{(k)}\}$  be the sequence resulting from a gradient algorithm  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)}$ . Let  $\gamma_k$  be as defined in Lemma 8.1, and suppose that  $\gamma_k > 0$  for all  $k$ . Then,  $\{\mathbf{x}^{(k)}\}$  converges to  $\mathbf{x}^*$  for any initial condition  $\mathbf{x}^{(0)}$  if and only if

$$\sum_{k=0}^{\infty} \gamma_k = \infty.$$

□

$$\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)}) = \mathbf{Q}\mathbf{x}^{(k)} - \mathbf{b}. \quad \gamma_k = \alpha_k \frac{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{Q}^{-1} \mathbf{g}^{(k)}} \left( 2 \frac{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}} - \alpha_k \right).$$

**Theorem 8.2** In the steepest descent algorithm, we have  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$  for any  $\mathbf{x}^{(0)}$ . □

# Analysis of gradient methods, convergence

**Theorem 8.3** For the fixed-step-size gradient algorithm,  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$  for any  $\mathbf{x}^{(0)}$  if and only if

$$0 < \alpha < \frac{2}{\lambda_{\max}(\mathbf{Q})}.$$

□

**Example 8.4** Let the function  $f$  be given by

$$f(\mathbf{x}) = \mathbf{x}^\top \begin{bmatrix} 4 & 2\sqrt{2} \\ 0 & 5 \end{bmatrix} \mathbf{x} + \mathbf{x}^\top \begin{bmatrix} 3 \\ 6 \end{bmatrix} + 24.$$

We wish to find the minimizer of  $f$  using a fixed-step-size gradient algorithm

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)}),$$

where  $\alpha \in \mathbb{R}$  is a fixed step size.

# Gradient methods, rate of convergens

**Theorem 8.4** *In the method of steepest descent applied to the quadratic function, at every step  $k$  we have*

$$V(\mathbf{x}^{(k+1)}) \leq \frac{\lambda_{\max}(Q) - \lambda_{\min}(Q)}{\lambda_{\max}(Q)} V(\mathbf{x}^{(k)}).$$

□

$$r = \frac{\lambda_{\max}(Q)}{\lambda_{\min}(Q)} = \|Q\| \|Q^{-1}\|,$$

called the *condition number* of  $Q$ . Then, it follows from Theorem 8.4 that

$$V(\mathbf{x}^{(k+1)}) \leq \left(1 - \frac{1}{r}\right) V(\mathbf{x}^{(k)}).$$