PURDUE UNIVERSITY · ECE 58000

OPTIMIZATION METHODS
PROF. ŻAK, PROF. CHONG

LECTURE NOTE 08
*Linhui Xie*
*February 1, 2023*

# 8  Gradient Method

- Gradient algorithm:
$$\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} - \alpha_k \nabla f\left(\boldsymbol{x}^{(k)}\right).$$

- Define $\boldsymbol{g}^{(k)} := \nabla f\left(\boldsymbol{x}^{(k)}\right)$ and set descent direction to $\boldsymbol{d}^{(k)} = -\boldsymbol{g}^{(k)}$.

- Step size $\alpha_k$ can be chosen in many different ways.

- For sufficiently small step size, the gradient algorithm has *descent* property.
  Define $\phi(\alpha) := f\left(\boldsymbol{x}^{(k)} - \alpha\boldsymbol{g}^{(k)}\right)$, then $\phi$ has Taylor expansion:

$$f\left(\boldsymbol{x}^{(k)} - \alpha\boldsymbol{g}^{(k)}\right) = f\left(\boldsymbol{x}^{(k)}\right) - \alpha \left\|\boldsymbol{g}^{(k)}\right\|^2 + o(\alpha)$$

  For $\alpha$ sufficiently small, we have

$$f\left(\boldsymbol{x}^{(k)} - \alpha\boldsymbol{g}^{(k)}\right) \leq f\left(\boldsymbol{x}^{(k)}\right)$$

- ▶ 𝒫𝒭𝒪𝒫𝒪𝒮𝒥𝒯𝒥𝒪𝒩 Suppose $\boldsymbol{g}^{(k)} = \nabla f\left(\boldsymbol{x}^{(k)}\right) \neq \boldsymbol{0}$. There exists $\bar{\alpha} > 0$ such that for all $\alpha_k \in (0, \bar{\alpha})$, we have
$$f\left(\boldsymbol{x}^{(k+1)}\right) < f\left(\boldsymbol{x}^{(k)}\right)$$

- Remark: if $\boldsymbol{g}^{(k)} = \boldsymbol{0}$, the FONC holds.
  Short proof:
  By chain rule, we have
$$\phi'(0) = f\left(\boldsymbol{x}^{(k)}\right) = -\left\|\boldsymbol{g}^{(k)}\right\|^2 < 0.$$

  Gradient is negative thus function value is decreasing. Hence, there exists $\bar{\alpha} > 0$ such that for all $\alpha_k \in (0, \bar{\alpha})$, we have

$$\phi\left(\alpha_k\right) < \phi(0).$$

  Rewriting, we obtain
$$f\left(\boldsymbol{x}^{(k+1)}\right) < f\left(\boldsymbol{x}^{(k)}\right).$$

- A variety of options exist for selecting $\alpha_k$.

- If $\alpha_k$ too small, an increased number of iterations may be needed to get optimal solution $\boldsymbol{x}^*$.

- If $\alpha_k$ too big, algorithm may lead to an oscillated track (zig-zag) around $\boldsymbol{x}^*$ (overshoot).

- General approach is to set a constant $\alpha_k = \alpha$ for all $k$.

- Steepest approach change $\alpha_k$ with each successive iteration.

## 8.1 Steepest descent algorithm

- Greedy scheme to pick for $\alpha_k$.

$$\alpha_k = \arg\min_{\alpha \geq 0} f\left(\boldsymbol{x}^{(k)} - \alpha \boldsymbol{g}^{(k)}\right).$$

- The steepest descent algorithm has the descent property.

▶ $\mathcal{PROPOSITION}$ 8.1 Let $\left\{\boldsymbol{x}^{(k)}\right\}$ be obtained by steepest descent method,

$$\left(\boldsymbol{x}^{(k+2)} - \boldsymbol{x}^{(k+1)}\right)^{\top} \left(\boldsymbol{x}^{(k+1)} - \boldsymbol{x}^{(k)}\right) = 0.$$

Short Proof:

Based on definition,
$$\boldsymbol{x}^{(k+2)} = \boldsymbol{x}^{(k+1)} - \alpha_{k+1}\boldsymbol{g}^{(k+1)},$$
$$\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} - \alpha_{k}\boldsymbol{g}^{(k)}.$$

Let $\phi(\alpha) = f\left(\boldsymbol{x}^{(k)} - \alpha \boldsymbol{g}^{(k)}\right) = f\left(\boldsymbol{x}^{(k+1)}\right)$. Since $\alpha_k = \arg\min \phi(\alpha)$, by FONC, we have

$$\phi'\left(\alpha_k\right) = 0.$$

Hence,

$$\phi'\left(\alpha_k\right) = \nabla f\left(\boldsymbol{x}^{(k)} - \alpha_k \boldsymbol{g}^{(k)}\right)^{\top} \boldsymbol{g}^{(k)} = \nabla f\left(\boldsymbol{x}^{(k+1)}\right)^{\top} \boldsymbol{g}^{(k)} \stackrel{\boldsymbol{g}^{(k)} = \nabla f\left(\boldsymbol{x}^{(k)}\right)}{=\!=\!=\!=\!=} \boldsymbol{g}^{(k+1)\top} \boldsymbol{g}^{(k)} = 0.$$

Therefore,
$$\left(\boldsymbol{x}^{(k+2)} - \boldsymbol{x}^{(k+1)}\right)^{\top} \left(\boldsymbol{x}^{(k+1)} - \boldsymbol{x}^{(k)}\right) = \alpha_{k+1}\alpha_k \boldsymbol{g}^{(k+1)\top} \boldsymbol{g}^{(k)} = 0.$$

- For a prescribed $\epsilon > 0$, terminate the iteration if one of the followings is met:

  ○ $\left\|\boldsymbol{g}^{(k)}\right\| < \epsilon$;
  ○ $\left\|\boldsymbol{x}^{(k+1)} - \boldsymbol{x}^{(k)}\right\| < \epsilon$;
  ○ $\left|f\left(\boldsymbol{x}^{(k+1)}\right) - f\left(\boldsymbol{x}^{(k)}\right)\right| < \epsilon$.

- More <u>preferable choices</u> using "relative change", because they are "scale-free".

  ○ $\left|f\left(\boldsymbol{x}^{(k+1)}\right) - f\left(\boldsymbol{x}^{(k)}\right)\right| / \left|f\left(\boldsymbol{x}^{(k)}\right)\right| < \epsilon$;
  ○ $\left\|\boldsymbol{x}^{(k+1)} - \boldsymbol{x}^{(k)}\right\| / \left\|\boldsymbol{x}^{(k)}\right\| < \epsilon$.

## 8.2 Analysis of optimization algorithms

- <u>Globally convergent</u>: starting from any initial point $\boldsymbol{x}^{(0)}$, an algorithm that generates sequence $\boldsymbol{x}^{(k)} \to \boldsymbol{x}^*$, where $\boldsymbol{x}^*$ satisfying the FONC.

- <u>Locally convergent</u>: starting from an initial point $\boldsymbol{x}^{(0)}$ is sufficiently close to $\boldsymbol{x}^*$, an algorithm generates sequences converges to $\boldsymbol{x}^*$.

- <u>Rate of convergence</u>: how fast an algorithm converges.