

09. GRADIENT DESCENT

- Descent methods
- Descent direction and step size
- Gradient descent.

Descent direction

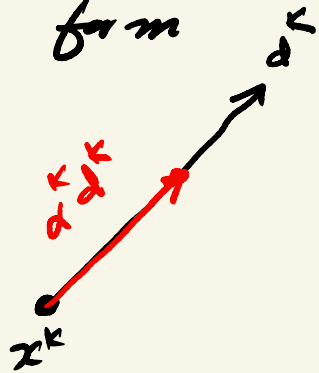
- Unconstrained optimization (non-linear)

$$\min_{x \in \mathbb{R}^n} f(x), \quad f: \mathbb{R}^n \rightarrow \mathbb{R} \quad \text{cont. diff.}$$

- We will consider iterative algorithms of the form

$$x^{k+1} = x^k + \alpha^k d^k$$

$$\alpha^k := \text{step size} \quad d^k := \text{descent direction}$$



- (Definition) A search direction is a descent direction for f at $x \in \mathbb{R}^n$ if the directional derivative of f at x is negative, i.e.:

$$f'(x; d) := \nabla f(x)^T d < 0.$$

Descent property

If f is cont. diff. and d is a descent direction,

then $\exists \varepsilon > 0$ s.t.

$$f(x+ad) < f(x)$$

for $d \in (0, \varepsilon]$.

$$\nabla f(x)^T d < 0$$

proof:

• Because $f'(x; d) < 0$ we have:

$$\lim_{d \rightarrow 0^+} \frac{f(x+ad) - f(x)}{d} < 0$$

• Using definition of limits, $\exists \varepsilon > 0$ s.t.

$$\frac{f(x+ad) - f(x)}{d} < 0 \quad \text{for } d \in (0, \varepsilon].$$

$\lim_{x \rightarrow a} f(x) = L$ use (ε, δ) definition of limit.

Gen. Method for Descent Direction.

Initialization: choose $x \in \mathbb{R}^n$.

For $k = 0, 1, 2, \dots$

(a) compute descent direction $d^k \in \mathbb{R}^n$.

(b) compute a stepsize α^k s.t. $f(x^k + \alpha^k d^k) < f(x^k)$

(c) update $x^{k+1} = x^k + \alpha^k d^k$

(d) Stopping criteria.

Question:

- How to choose d^k ?
- what are advantages/disadvantages of different d^k ?
- How do we compute α^k ?
- stopping criteria?

Stepsize selection d^k .

- Constant stepsize: $d^k = \bar{d} \in \mathbb{R}$ for all k .
- Exact line search: choose d^k that minimizes the 1-dimensional optimization problem
$$d^k := \arg \min_{d \geq 0} f(x^k + d d^k)$$

- Diminishing stepsize: choose d^k that satisfies
$$d^k \rightarrow 0 \quad \text{and} \quad \sum_{k=1}^{\infty} d^k = \infty$$

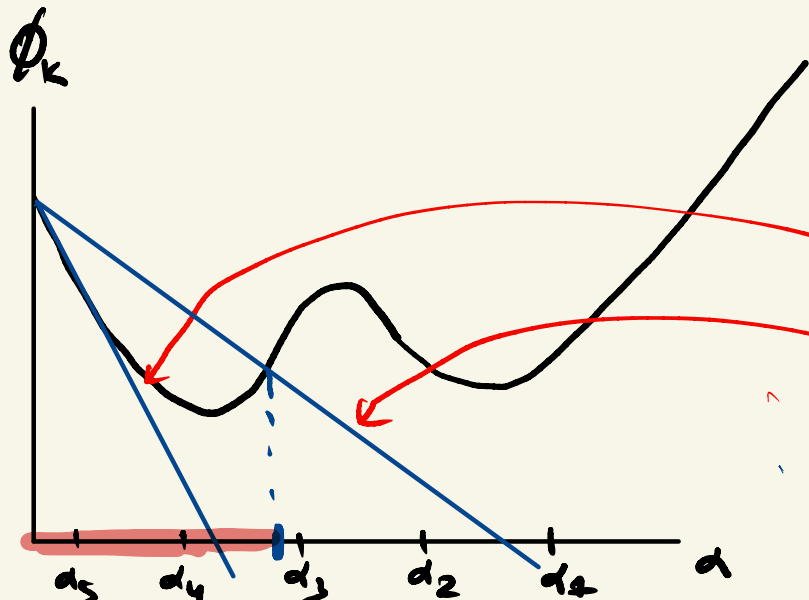
eg: $d^k = \frac{1}{k}, \frac{1}{\sqrt{k}}$

- Backtracking "Armijo's" linesearch: For some parameter $\mu \in (0, 1)$ reduce stepsize α (eg: $\alpha \leftarrow \alpha/2$ starting at $\alpha=1$) until

$$f(x^k) - f(x^k + \alpha d^k) \geq -\mu \alpha \nabla f(x^k)^T d^k.$$

$$\phi_k(\alpha) = f(x^k + \alpha d^k)$$

$$\phi_k(0) = \nabla f(x^k)^T d^k.$$



$$f(x^k) + \alpha \nabla f(x^k)^T d^k.$$

$$f(x^k) + \alpha \mu \nabla f(x^k)^T d^k.$$

α_4, α_5 satisfies the sufficient decrease condition.

Sufficient decrease condition.

- The sufficient decrease condition is satisfied for small enough α^k .

- suppose $d \neq 0$, $d \in \mathbb{R}^n$ is a descent direction of f at x . let $\mu \in (0, 1)$. Then exists $\varepsilon > 0$ s.t.

$$f(x) - f(x + \alpha d) \geq -\mu \alpha \nabla f(x)^T d,$$

for some $\alpha \in (0, \varepsilon]$

Exact line search for quadratic function.

An exact line search may not be possible for general functions but is possible for quadratic functions.

$$f(x) = \frac{1}{2} x^T A x + b^T x + c, \quad A \succeq 0, \quad A = A^T$$

Exact line search solve 1-d optimization problem

$$\min_{d \geq 0} f(x+ad). \quad \nabla f(x) = Ax + b.$$

Derivation: $f(x+ad) = \frac{1}{2} (x+ad)^T A (x+ad) + b^T (x+ad) + c$

$$= \frac{1}{2} x^T A x + a x^T A d + \frac{a^2}{2} d^T A d + b^T x + a b^T d + c$$

$$\Rightarrow \frac{d}{da} f(x+ad) = x^T A d + a d^T A d + b^T d = a d^T A d + d^T \nabla f(x)$$

$$\text{so, } \frac{d}{da} f(x+ad) = 0 \Rightarrow d = \frac{-d^T \nabla f(x)}{d^T A d} \succ 0$$

Gradient descent

$$d^k = -g_k \quad g_k \equiv \nabla f(x^k)$$

- The negative gradient direction $-g_k$ provides a descent direction:

$$f'(x^k, -g_k) = -\nabla f(x^k)^T g_k = -\|g_k\|_2^2 < 0$$

if x^k is not a stationary point

- The negative gradient $d \equiv -\nabla f(x)$ is the steepest direction of descent i.e.

$$\min_d \{ f'(x; d) \mid \|d\|_2 = 1 \}$$

proof: $f'(x; d) = \nabla f(x)^T d$
 $\geq -\|\nabla f(x)\| \|d\|$ Cauchy Schwartz
 $= -\|\nabla f(x)\|$ $|\nabla f^T d| \leq \|\nabla f\| \|d\|$

The lower is achieved by $d = \frac{-\nabla f(x)}{\|\nabla f(x)\|}$

Input : $\epsilon > 0$ tolerance
 x_0 starting point.

For $k=0, 1, 2, \dots$

- compute gradient $g^k = \nabla f(x^k)$
- choose a step size α^k that satisfies
 $f(x^k - \alpha^k g^k) \leq f(x^k)$. $\phi_k(x^k - \alpha g_k)$
- $x^{k+1} = x^k - \alpha^k g^k$
- stop if $\|\nabla f(x^{k+1})\| \leq \epsilon$.

