

CPSC 406
Computational Optimization
Dept of Computer Science
University of British Columbia

GRADIENT DESCENT

DESCENT DIRECTIONS

- Unconstrained nonlinear optimization:

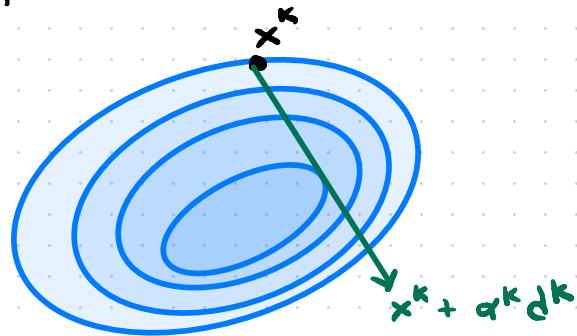
$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x), \quad f: \mathbb{R}^n \rightarrow \mathbb{R}, \quad \text{continuously differentiable}$$

- We will consider iterative algorithms of the form

$$x^{k+1} = x^k + \alpha^k d^k, \quad k = 0, 1, 2, \dots$$

where

- $d^k \equiv$ search direction
- $\alpha^k \equiv$ step length



- A search direction $d \neq 0$ is a descent dir for f at x if the directional derivative is negative, i.e.,

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

DESCENT PROPERTY

If $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable, and $d \in \mathbb{R}^n$ is a descent direction at x , then for some $\varepsilon > 0$,

$$f(x + \alpha d) < f(x) \quad \forall \alpha \in (0, \varepsilon] \quad (\text{Descent})$$

Proof:

- Because $f'(x; d) < 0$

$$\lim_{\alpha \rightarrow 0^+} \frac{f(x + \alpha d) - f(x)}{\alpha} = f'(x; d) < 0$$

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

$$\frac{f(x + \alpha d) - f(x)}{\alpha} < 0 \quad \forall \alpha \in (0, \varepsilon]$$

which implies (Descent)

GENERIC DESCENT METHOD — conceptual algorithm

Initialization: choose $x_0 \in \mathbb{R}^n$

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

(a) compute descent direction d^k

(b) compute stepsize α^k st $f(x^k + \alpha^k d^k) < f(x^k)$

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

(d) check stopping criteria

Questions

- How to determine a starting point?
- What are advantages / disadvantages of different dirs d^k ?
- How to compute a step length α^k ?
- When to stop?

STEP SIZE SELECTION α^k

[descent direction computation later]

STEP SIZE SELECTION (α^k)

These are the selection rules most used in practice:

1. Constant stepsize: $\alpha^k = \bar{\alpha} \quad \forall k$ (needs additional conditions on f to be reliable)

2. Exact linesearch: choose α^k to minimize f along ray $x^k + \alpha d^k$:

$$\alpha^k \in \operatorname{argmin}_{\alpha \geq 0} f(x^k + \alpha d^k)$$

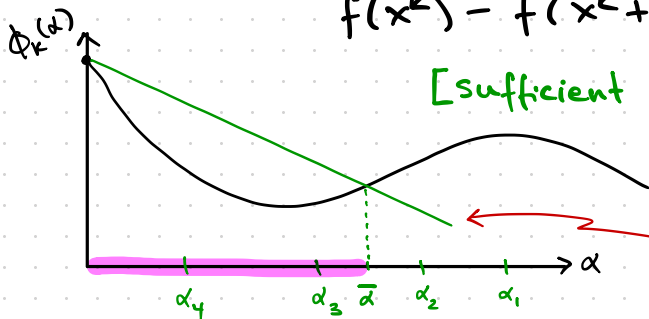
3. Backtracking "Armijo" Linesearch: for some parameter $\mu \in (0, 1)$ reduce α (eg, $\alpha \leftarrow \alpha/2$ beginning with $\alpha=1$) until

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

[sufficient decrease property - see $\alpha_3 \nless \alpha_4$ below]

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

$$f(x^k) + \underbrace{\mu \alpha \nabla f(x^k)^T d^k}_{\geq 0 \text{ by def'n of } d^k}$$



EXACT LINESEARCH FOR QUADRATIC FUNCTIONS

An exact linesearch is typically only possible for quadratic func's:

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

Exact linesearch solves the 1-dimensional optimization problem

$$\min_{\alpha \geq 0} f(x + \alpha d) \quad (\text{where } d \text{ is descent dir})$$

Derivation of solution [details in class]:

$$f(x + \alpha d) = \frac{1}{2} (x + \alpha d)^T A (x + \alpha d) + b^T (x + \alpha d) + c$$

$$\frac{d}{d\alpha} f(x + \alpha d) = \alpha d^T A d + x^T A d + b^T d = \alpha d^T A d + \nabla f(x)^T d$$

$$\frac{d}{d\alpha} f(x + \alpha d) = 0 \iff \left\{ \alpha = \frac{-\nabla f(x)^T d}{d^T A d} > 0 \right\}$$

$$\iff \left\{ \begin{array}{l} d^T A d > 0 \quad (A \succ 0) \\ \nabla f(x)^T d < 0 \quad (d \text{ is descent dir}) \end{array} \right\}$$

SEARCH DIRECTIONS d^k

GRADIENT DESCENT

$$d^k := -g^k \quad \text{where} \quad g^k := \nabla f(x^k)$$

- The negative gradient direction $(-g_k)$ provides descent:

$$f'(x^k; -g^k) = -g_k^T g_k = -\|g_k\|^2 < 0$$

if $g_k \neq 0$, ie, x^k is not already stationary.

- The negative gradient $g \equiv -\nabla f(x)$ is the steepest descent direction of f at x , ie, it solves

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

[set $g \equiv \nabla f(x)$]

Proof: $f'(x; d) \equiv g^T d \geq -\|g\| \cdot \|d\|$
 $\geq -\|g\|$

[Cauchy-Schwartz Ineq]

[$\|d\| = 1$]

Lower bound is achieved by setting $d = -g / \|g\|$

GRADIENT METHOD

Input : $\epsilon > 0$ (tolerance)
 x_0 (starting iterate)

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

- evaluate gradient $g^k = \nabla f(x^k)$

- choose step length α^k based on reducing the function

$$\phi(\alpha) = f(x^k - \alpha g^k)$$

[see stepsize selection slide]

- $x^{k+1} = x^k - \alpha^k g^k$

- STOP if $\|\nabla f(x^{k+1})\| < \epsilon$

[DEMO on function $f(x,y) = x^2 + 2y^2$]

"ZIG-ZAG" OF GRADIENT METHOD

Let x_1, x_2, x_3, \dots be the iterates generated by the gradient method. Then

$$(x_{k+2} - x_{k+1})^\top (x_{k+1} - x_k) = 0$$

Proof By definition of the gradient update:

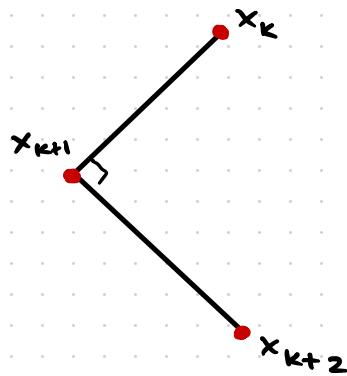
$$\left. \begin{aligned} x_{k+1} &= x_k - \alpha_k \nabla f_k \\ x_{k+2} &= x_{k+1} - \alpha_{k+1} \nabla f_{k+1} \end{aligned} \right\} \Rightarrow \begin{aligned} x_{k+1} - x_k &= -\alpha_k \nabla f_k \\ x_{k+2} - x_{k+1} &= -\alpha_{k+1} \nabla f_{k+1} \end{aligned}$$

$$\Rightarrow (x_{k+2} - x_{k+1})^\top (x_{k+1} - x_k) = 0 \Leftrightarrow \nabla f_k^\top \nabla f_{k+1} = 0.$$

Because $\alpha^k \in \arg \min \{ \phi(\alpha) := f(x_k - \alpha \nabla f_k) \}$

$$0 = \phi'(\alpha_k) = -\nabla f^\top \underbrace{x_k - \alpha \nabla f_k}_{\equiv x_{k+1}} = -\nabla f_k^\top \nabla f(x_{k+1})$$

$$\Leftrightarrow \nabla f(x_k)^\top \nabla f(x_{k+1}) = 0 //$$



"Zig-Zag" behavior often the reason why the gradient method is slow. 10

GRADIENT METHOD WITH CONSTANT STEPSIZE

- Constant stepsize sets $\alpha^k = \bar{\alpha}$ for all k .
- How to choose $\bar{\alpha}$?
 - $\bar{\alpha}$ too small \Rightarrow gradient method slow
 - $\bar{\alpha}$ too large \Rightarrow gradient method diverges

[Demo on $\min x^2 + 2y^2$ with
different values of constant stepsize]

- Must choose steplength $\bar{\alpha} \in (0, \alpha_{\max})$ for method to converge.
- α_{\max} depends on a property of $\nabla f(x)$ called Lipschitz continuity.

LIPSCHITZ CONTINUITY OF GRADIENT

A continuously differentiable function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ has a Lipschitz continuous gradient with parameter L if

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\| \quad [2\text{-norm throughout}]$$

for all vectors x, y and some $L > 0$ constant.

EXAMPLE $f(x) = \frac{1}{2} x^T A x + b^T x + c$ (quadratic, $A = A^T$)

$$\nabla f(x) = Ax + b$$

$$\begin{aligned} \|\nabla f(x) - \nabla f(y)\| &= \|(Ax + b) - (Ay + b)\| \\ &= \|Ax - Ay\| \\ &= \|A(x - y)\| \leq \|A\| \cdot \|x - y\| \\ &\quad \uparrow \|A\|_2 = \lambda_{\max}(A) \end{aligned}$$

EXAMPLE: Take $A = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$ from previous slide. $\|A\|_2 = 2$.

CONSTANT STEP-SIZE THRESHOLD

- If $f: \mathbb{R}^n \rightarrow \mathbb{R}$ has an L -Lipschitz continuous gradient and a minimizer exists, then the gradient method with constant stepsize converges if

$$\bar{\alpha} \in (0, 2/L).$$

- Previous quadratic example

- $f(x) = \frac{1}{2}x^T A x + b^T x + c$ with $A = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$

- $L = \lambda_{\max}(A) = 2$

- gradient method converges for all $\bar{\alpha} \in (0, 1)$

CONVERGENCE OF THE GRADIENT METHOD

For the minimization of $f: \mathbb{R}^n \rightarrow \mathbb{R}$ bnd below with L -Lipschitz gradient and one of the linesearches

(1) constant step size $\bar{\alpha} \in (0, 2/L)$

(2) exact linesearch

(3) backtracking linesearch with $\mu \in (0, 1)$

Then

(a) $f(x_{k+1}) < f(x_k)$ for all $k=0, 1, 2, \dots$ unless $\nabla f(x_k) = 0$

(Decreasing)

(b) $\|\nabla f(x_k)\| \rightarrow 0$ as $k \rightarrow \infty$.

(convergence to stationary point)

(Backtracking line search demo)

CONDITION NUMBER OF A MATRIX

The condition number of a $n \times n$ positive definite matrix A is defined by

$$K(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)} \geq 1$$

- ill-conditioned matrices have $K(A)$ large
- Condition number of the Hessian at the solution influences the speed at which the gradient method converges

$$H \equiv \nabla^2 f(x^*)$$

$K(H)$ small \Rightarrow GM typically converges quickly

$K(H)$ large \Rightarrow GM " " slowly

EXAMPLE: ROSENBRock FUNCTION

$$f(x_1, x_2) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

$$\nabla f(x_1, x_2) = \begin{bmatrix} -400x_1(x_2 - x_1^2) - 2(1 - x_1) \\ 200(x_2 - x_1^2) \end{bmatrix}$$

$$\nabla^2 f(x_1, x_2) = \begin{bmatrix} -400x_2 + 1200x_1^2 + 2 & -400x_1 \\ -400x_1 & 200 \end{bmatrix}$$

Solution $(x_1, x_2) = (1, 1)$ is unique. [Verify $\nabla f(1, 1) = 0$]

$$\nabla^2 f(1, 1) = \begin{bmatrix} 802 & -400 \\ -400 & 200 \end{bmatrix}$$

[Backtracking Demo]

SCALED GRADIENT METHOD

$$(P) \quad \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) \quad f: \mathbb{R}^n \rightarrow \mathbb{R}$$

- Make a linear change of variables with

$$S \text{ nonsingular } n \times n, \quad x = Sy \quad \text{ie} \quad y := S^{-1}x$$

$$(P_{\text{scaled}}) \quad \underset{y \in \mathbb{R}^n}{\text{minimize}} \quad g(y) := f(Sy)$$

- Apply gradient method to scaled problem:

$$y_{k+1} = y_k - \alpha_k \nabla g(y_k) \quad \text{with} \quad \nabla g(y) = S^T \nabla f(Sy)$$

- Multiply on left by S :

$$x_{k+1} = x_k - \alpha_k S S^T \nabla f(x_k)$$

- Scaled gradient method: with $D = S S^T$,

$$x_{k+1} = x_k - \alpha_k D \nabla f(x_k)$$

SCALED DESCENT

The scaled gradient $-D\nabla f(x)$ is a descent direction:

$$f'(x; -D\nabla f(x)) = -\nabla f(x)^T D \nabla f(x) < 0$$

because $D = SS^T \succ 0$ (S nonsingular)

Scaled Gradient Method

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

- choose scaling matrix D_k
- compute scaled gradient $d_k = D_k \nabla f(x_k)$
- compute steplength α_k by linesearch on the func'n
$$\phi(\alpha) = f(x_k - \alpha d_k)$$
- $x_{k+1} = x_k - \alpha_k d_k$
- STOP if $\|\nabla f(x_{k+1})\| \leq \text{tol}$

GAUSS-NEWTON METHOD

$$(NLS) \text{ minimize}_{x \in \mathbb{R}^n} f(x) = \frac{1}{2} \|\mathbf{r}(x)\|^2$$

$$\begin{aligned} \mathbf{r}_i: \mathbb{R}^n &\rightarrow \mathbb{R} \\ \text{cont. diff}^1 \\ i &= 1, \dots, m \end{aligned}$$

Gauss-Newton Method:

$$x_{k+1} = \underset{x}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{r}_k + \mathbf{A}_k(x - x_k)\|^2$$

[Linearization of \mathbf{r} at x_k]

$$\begin{aligned} \mathbf{r}_k &\equiv \mathbf{r}(x_k) \\ \mathbf{A}_k &= \begin{bmatrix} \nabla \mathbf{r}_1(x_k)^T \\ \vdots \\ \nabla \mathbf{r}_m(x_k)^T \end{bmatrix} \end{aligned}$$

$$= (\mathbf{A}_k^T \mathbf{A}_k)^{-1} \mathbf{A}_k^T (\mathbf{A}_k x_k - \mathbf{r}_k)$$

$$= x_k - (\mathbf{A}_k^T \mathbf{A}_k)^{-1} \mathbf{A}_k^T \mathbf{r}_k$$

$$= x_k - (\mathbf{A}_k^T \mathbf{A}_k)^{-1} \nabla f(x_k)$$

$$\nabla f(x_k) = \mathbf{A}_k^T \mathbf{r}_k$$

Thus, we see that the Gauss-Newton method is a scaled gradient method with the scaling matrix

$$\mathbf{D}_k = (\mathbf{A}_k^T \mathbf{A}_k)^{-1}$$