

Lecture 2: Gradient Methods

- Gradient Methods - Motivation
- Principal Gradient Methods
- Gradient Methods - Choices of Direction
- Asymptotic Convergence
- Local Convergence Rate

Motivation

Consider the minimization of a differentiable function f

$$\begin{array}{ll} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in \mathbb{R}^n \end{array}$$

Suppose we are at a point \mathbf{x} .

- If $\nabla f(\mathbf{x}) = \mathbf{0}$, then \mathbf{x} is a candidate solution. Done.
- If $\nabla f(\mathbf{x}) \neq \mathbf{0}$, there is an interval $(0, \delta)$ of stepsizes such that

$$f(\mathbf{x} - \alpha \nabla f(\mathbf{x})) < f(\mathbf{x}), \quad \text{for all } \alpha \in (0, \delta).$$

- Moreover, if the direction \mathbf{d} makes an angle with $\nabla f(\mathbf{x})$ that is greater than 90 degrees,

$$\nabla f(\mathbf{x})' \mathbf{d} < 0,$$

there is an interval $(0, \delta)$ of stepsizes such that

$$f(\mathbf{x} + \alpha \mathbf{d}) < f(\mathbf{x}) \quad \text{for all } \alpha \in (0, \delta).$$

Iterative Descent Methods

$$\mathbf{x}^{r+1} = \mathbf{x}^r + \alpha_r \mathbf{d}^r, \quad r = 0, 1, \dots$$

where, if $\nabla f(\mathbf{x}^r) \neq \mathbf{0}$, the direction \mathbf{d}^r satisfies $\nabla f(\mathbf{x}^r)' \mathbf{d}^r < 0$, and α_r is a positive stepsize.

- Principal example:

$$\mathbf{x}^{r+1} = \mathbf{x}^r - \alpha_r \mathbf{D}^r \nabla f(\mathbf{x}^r), \quad r = 0, 1, \dots$$

where \mathbf{D}^r is a positive definite symmetric matrix

- Simplest method: Steepest descent

$$\mathbf{x}^{r+1} = \mathbf{x}^r - \alpha_r \nabla f(\mathbf{x}^r), \quad r = 0, 1, \dots$$

- Most sophisticated method: Newton's method

$$\mathbf{x}^{r+1} = \mathbf{x}^r - \alpha_r \nabla^2 f(\mathbf{x}^r)^{-1} \nabla f(\mathbf{x}^r), \quad r = 0, 1, \dots$$

Discussion

- Steepest descent: no memory, slow convergence
 - ★ zigzagging
 - ★ consider the convex quadratic case; convergence related to the condition number. Plot.
- Newton's method: no memory, fast convergence (with $\alpha_r = 1$);
 - ★ Given \mathbf{x}^r , the method obtains \mathbf{x}^{r+1} as the minimum of a quadratic approximation of f based on a second order Taylor expansion around \mathbf{x}^r . Plot.
 - ★ insensitive to condition number
 - ★ how many Newton iterations does it take to minimize a quadratic function f ?

Other Choices of Direction

- Diagonally Scaled Steepest Descent:

$$\mathbf{D}^r = \text{Diagonal approximation to } (\nabla^2 f(\mathbf{x}^r))^{-1}.$$

- Modified Newton's Method:

$$\mathbf{D}^r = (\nabla^2 f(\mathbf{x}^0))^{-1}, \quad r = 0, 1, \dots$$

- Discretized Newton's Method:

$$\mathbf{D}^r = (\mathbf{H}(\mathbf{x}^r))^{-1}, \quad r = 0, 1, \dots,$$

where $\mathbf{H}(\mathbf{x}^r)$ is a finite-difference based approximation of $\nabla^2 f(\mathbf{x}^r)$.

- Gauss-Newton method for least squares problems $\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{g}(\mathbf{x})\|^2$. Here \mathbf{D}^r is an approximation of the true Hessian inverse $(\nabla^2 \mathbf{g}(\mathbf{x}^r))^{-1}$

$$\mathbf{D}^r = (\nabla \mathbf{g}(\mathbf{x}^r) \nabla \mathbf{g}(\mathbf{x}^r)')^{-1}, \quad r = 0, 1, \dots$$

where $\nabla \mathbf{g}(\mathbf{x}^r)$ is the Jacobian matrix of \mathbf{g} .

Choices of Stepsize

- **Minimization Rule:** α_r is such that

$$f(\mathbf{x}^r + \alpha_r \mathbf{d}^r) = \min_{\alpha \geq 0} f(\mathbf{x}^r + \alpha \mathbf{d}^r).$$

- **Limited Minimization Rule:** Fix some $s > 0$. Choose α_r such that

$$f(\mathbf{x}^r + \alpha_r \mathbf{d}^r) = \min_{\alpha \in [0, s]} f(\mathbf{x}^r + \alpha \mathbf{d}^r).$$

- **Constant stepsize:** α_r is such that $\alpha_r = s$ is a constant
- **Diminishing stepsize:** $\alpha_r \rightarrow 0$ but satisfies the infinite travel condition

$$\sum_{r=0}^{\infty} \alpha_r = \infty$$

- **Armijo rule:** Let $\sigma \in (0, \frac{1}{2})$. Start with s and continue with $\beta s, \beta^2 s, \dots$, until $\beta^m s$ falls within the set of α with

$$f(\mathbf{x}^r) - f(\mathbf{x}^r + \alpha \mathbf{d}^r) \geq -\sigma \alpha \nabla f(\mathbf{x}^r)' \mathbf{d}^r.$$

Claim: if $\mathbf{d}^r = -\mathbf{D}^r \nabla f(\mathbf{x}^r) \neq \mathbf{0}$ with $\mathbf{D}^r \succ \mathbf{0}$, then m is finite.

Convergence of Iterative Methods

- **Convergence to stationary points**
 - ★ Sanity check
 - ★ Minimal requirement of any reasonable algorithm
 - ★ Does not give global efficiency of the algorithm
- **Asymptotic convergence rate:** local analysis, assuming already close to a solution, let $\#$ of iterations go to infinity
 - ★ **1) linear, 2) supperlinear, 3) sublinear**
- **Iteration complexity analysis:**
 - ★ Measures the number of iterations required to get an ϵ optimal solution (e.g., $f(\mathbf{x}^r) - f^* \leq \epsilon$)
 - ★ Current analysis is all for the worst case and requires convexity
 - ★ Gives global behavior of the algorithm

Convergence of Gradient Descent Methods

- Only convergence to stationary points can be guaranteed
- Even convergence to a single limit may be hard to guarantee (capture theorem)
- Danger of nonconvergence if directions \mathbf{d}^r tend to be orthogonal to $\nabla f(\mathbf{x}^r)$
- **Gradient related condition:** For any subsequence $\{\mathbf{x}^r\}_{r \in \mathcal{K}}$ that converges to a nonstationary point, the corresponding subsequence $\{\mathbf{d}^r\}_{r \in \mathcal{K}}$ is bounded and satisfies

$$\limsup_{r \rightarrow \infty, r \in \mathcal{K}} \nabla f(\mathbf{x}^r)' \mathbf{d}^r < 0.$$

- The condition is satisfied if $\mathbf{d}^r = -\mathbf{D}^r \nabla f(\mathbf{x}^r)$ and the eigenvalues of \mathbf{D}^r are bounded above and bounded away from zero

Convergence Issues

- Let $\{\mathbf{x}^r\}$ be a sequence generated by a gradient method

$$\mathbf{x}^{r+1} = \mathbf{x}^r + \alpha_r \mathbf{d}^r,$$

where $\{\mathbf{d}^r\}$ is gradient related. Assume that for some constant $L > 0$, we have

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$$

Assume that either (a) there exists a scalar ϵ such that for all r

$$0 < \epsilon \leq \alpha_r \leq \frac{(2 - \epsilon)|\nabla f(\mathbf{x}^r)' \mathbf{d}^r|}{L\|\mathbf{d}^r\|^2} \quad (1)$$

or (b) $\alpha_r \rightarrow 0$ and

$$\sum_{r=1}^{\infty} \alpha_r = \infty.$$

Then either $f(\mathbf{x}^r) \rightarrow -\infty$ or else $\{f(\mathbf{x}^r)\}$ converges to a finite value and $\nabla f(\mathbf{x}^r) \rightarrow \mathbf{0}$.

- Specialize to the steepest descent: $\mathbf{d}^r = -\nabla f(\mathbf{x}^r)$.

Convergence Analysis

- Given \mathbf{x}^r and the descent direction \mathbf{d}^r , the Lipschitz assumption implies

$$f(\mathbf{x}^r + \alpha \mathbf{d}^r) - f(\mathbf{x}^r) \leq \alpha \nabla f(\mathbf{x}^r)' \mathbf{d}^r + \frac{1}{2} \alpha^2 L \|\mathbf{d}^r\|^2$$

- Minimization of this function over α yields the stepsize

$$\bar{\alpha}_r = -\frac{\nabla f(\mathbf{x}^r)' \mathbf{d}^r}{L \|\mathbf{d}^r\|^2}$$

- In case of steepest descent, $\bar{\alpha}_r = 1/L$.
- The constant stepsize (1) reduces the cost function f at each iteration.

$$f(\mathbf{x}^r + \alpha_r \mathbf{d}^r) - f(\mathbf{x}^r) \leq \alpha_r \mu \nabla f(\mathbf{x}^r)' \mathbf{d}^r \quad (2)$$

for some constant $\mu > 0$.

Convergence Analysis, Continued

- Assume $\bar{\mathbf{x}}$ is a nonstationary limit point. Then $f(\mathbf{x}^r) \downarrow f(\bar{\mathbf{x}})$, it follows from (2) that

$$\alpha_r \nabla f(\mathbf{x}^r)' \mathbf{d}^r \rightarrow 0.$$

- If $\{\mathbf{x}^r\}_{r \in \mathcal{K}} \rightarrow \bar{\mathbf{x}}$, then by the gradient relatedness

$$\limsup_{r \rightarrow \infty, r \in \mathcal{K}} \nabla f(\mathbf{x}^r)' \mathbf{d}^r < 0.$$

- This implies $\{\alpha_r\}_{r \in \mathcal{K}} \rightarrow 0$, contradicting stepsize (1).
- For the diminishing stepsize rule (b), assume the whole sequence converges to a non-stationary point $\bar{\mathbf{x}}$. Then the gradient relatedness implies

$$\limsup_{r \rightarrow \infty} \nabla f(\mathbf{x}^r)' \mathbf{d}^r < 0.$$

Sum (2) over all r shows $f(\mathbf{x}^r) \rightarrow -\infty$, contradiction.

- Similar proof for the Armijo rule. For large $r \in \mathcal{K}$

$$f(\mathbf{x}^r) - f(\mathbf{x}^r + (\alpha_r/\beta)\mathbf{d}^r) < -\sigma \frac{\alpha_r}{\beta} \nabla f(\mathbf{x}^r)' \mathbf{d}^r.$$

- Defining $\mathbf{p}^r = \mathbf{d}^r / \|\mathbf{d}^r\|$ and $\bar{\alpha}_r = \alpha_r \|\mathbf{d}^r\| / \beta$, we have

$$\frac{f(\mathbf{x}^r) - f(\mathbf{x}^r + \bar{\alpha}_r \mathbf{p}^r)}{\bar{\alpha}_r} < -\sigma \nabla f(\mathbf{x}^r)' \mathbf{p}^r.$$

- Letting $k \rightarrow \infty$, we get

$$-\nabla f(\bar{\mathbf{x}})' \bar{\mathbf{p}} \leq -\sigma \nabla f(\bar{\mathbf{x}})' \bar{\mathbf{p}},$$

where $\bar{\mathbf{p}}$ is a limit point of $\{\mathbf{p}^r\}_{r \in \mathcal{K}}$.

- This is a contradiction since $\nabla f(\bar{\mathbf{x}})' \bar{\mathbf{p}} < 0$.

Other Stepsize Rules

- **Non-monotone step size rule:** Fix some integer $K > 0$. Reduce α by a factor of β until

$$f(\mathbf{x}^r + \alpha \mathbf{d}^r) - \max_{k=1,2,\dots,K} f(\mathbf{x}^{r-k}) \leq \sigma \alpha \nabla f(\mathbf{x}^r)' \mathbf{d}^r.$$

The resulting sequence $\{f(\mathbf{x}^r)\}$ is no longer monotonically decreasing. But the convergence to stationary point can still be established.

- **BB (Barzilai-Borwein) stepsize rule:** consider minimizing a strongly convex quadratic function

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{2} \mathbf{x}' \mathbf{A} \mathbf{x} + \mathbf{b}' \mathbf{x}$$

The steepest descent method with exact line minimization stepsize rule is

$$\alpha_r = \frac{\|\nabla f(\mathbf{x}^r)\|^2}{(\nabla f(\mathbf{x}^r))' \mathbf{A} \nabla f(\mathbf{x}^r)}.$$

The BB stepsize rule uses

$$\alpha_r^{\text{BB}} = \frac{\|\nabla f(\mathbf{x}^{r-1})\|^2}{(\nabla f(\mathbf{x}^{r-1}))' \mathbf{A} \nabla f(\mathbf{x}^{r-1})}.$$

Practical performance of BB stepsize rule is quite good (compared to the steepest descent) in some applications, but the theoretical understanding of BB method remains limited (more later).

Local Convergence Rate Analysis

- Restrict attention to sequences $\{\mathbf{x}^r\}$ converging to a local min \mathbf{x}^* .
- Measure progress in terms of an error function $e(\mathbf{x})$ with $e(\mathbf{x}^*) = 0$, such as

$$e(\mathbf{x}) = \|\mathbf{x} - \mathbf{x}^*\|, \quad e(\mathbf{x}) = f(\mathbf{x}) - f(\mathbf{x}^*)$$

- Compare the tail of the sequence $e(\mathbf{x}^r)$ with the tail of standard sequences.
- Geometric or linear convergence [if $e(\mathbf{x}^r) \leq q\beta^r$ for some $q > 0$ and $\beta \in [0, 1)$, and for all r]. Holds if

$$\limsup_{r \rightarrow \infty} \frac{e(\mathbf{x}^{r+1})}{e(\mathbf{x}^r)} < \beta$$

- Superlinear convergence [if $e(\mathbf{x}^r) \leq q\beta^{p^r}$ for some $q > 0$, $p > 1$ and $\beta \in [0, 1)$, and for all r].

Quadratic Model Analysis

- Focus on the quadratic function $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}'\mathbf{Q}\mathbf{x}$, with $\mathbf{Q} \succ \mathbf{0}$. Optimal solution $\mathbf{x}^* = ?$
- Analysis also applies to nonquadratic problems in the neighborhood of a nonsingular local min. Consider steepest descent

$$\begin{aligned}\mathbf{x}^{r+1} &= \mathbf{x}^r - \alpha_r \nabla f(\mathbf{x}^r) = (\mathbf{I} - \alpha_r \mathbf{Q})\mathbf{x}^r \\ \Rightarrow \|\mathbf{x}^{r+1}\|^2 &= (\mathbf{x}^r)'(\mathbf{I} - \alpha_r \mathbf{Q})^2 \mathbf{x}^r \leq \max \text{eig.}(\mathbf{I} - \alpha_r \mathbf{Q})^2 \|\mathbf{x}^r\|^2\end{aligned}$$

- The eigenvalues of $(\mathbf{I} - \alpha_r \mathbf{Q})^2$ are equal to $(1 - \alpha_r \lambda_i)^2$, where λ_i are the eigenvalues of \mathbf{Q} , so

$$\max \text{eig of } (\mathbf{I} - \alpha_r \mathbf{Q})^2 = \max\{(1 - \alpha_r m)^2, (1 - \alpha_r M)^2\}$$

where m, M are the smallest and largest eigenvalues of \mathbf{Q} . Thus

$$\frac{\|\mathbf{x}^{r+1}\|}{\|\mathbf{x}^r\|} \leq \max\{|1 - \alpha_r m|, |1 - \alpha_r M|\}$$

Optimal Convergence Rate

- The value of α_r that minimizes the bound is $\alpha^* = 2/(M + m)$, in which case

$$\frac{\|\mathbf{x}^{r+1}\|}{\|\mathbf{x}^r\|} \leq \frac{M - m}{M + m}$$

- Convergence rate of f for the minimization stepsize

$$\frac{f(\mathbf{x}^{r+1})}{f(\mathbf{x}^r)} \leq \left(\frac{M - m}{M + m} \right)^2$$

- The ratio $\kappa := M/m$ is called the condition number of Q , and problems with M/m : large are called ill-conditioned.
- Convergence rate is independent of problem dimension.
- **Iteration complexity:** find an ϵ -optimal solution, we need

$$O(\kappa \ln D/\epsilon)$$

iterations, where $D := \|\mathbf{x}^0 - \mathbf{x}^*\|^2$.

- Converges in **one** step if $M = m$.

Causes of Ill-Conditioning

- Ill-conditioning usually arises when optimization variables are scaled incoherently.
- Example:
 - ★ x_1 : gas flow, lbs/hr, nominal value 11,000
 - ★ x_2 : waste buildup: lbs, nominal value 6×10^{-4} .
 - ★ x_3 : steam thermal resistance, $\text{BTU}/(\text{hr} \cdot \text{ft}^2 \cdot ^\circ\text{F})^{-1}$; nominal value 100
- Should rescale the units so that the nominal values are balanced.

Scaled Steepest Descent

- View the more general method

$$\mathbf{x}^{r+1} = \mathbf{x}^r - \alpha_r \mathbf{D}^r \nabla f(\mathbf{x}^r)$$

as a scaled version of steepest descent.

- Consider a change of variables $\mathbf{x} = \mathbf{S}\mathbf{y}$ with $\mathbf{S} = (\mathbf{D}^r)^{1/2}$. In the space of \mathbf{y} , the problem is

$$\begin{array}{ll} \text{minimize} & h(\mathbf{y}) := f(\mathbf{S}\mathbf{y}) \\ \text{subject to} & \mathbf{y} \in \mathbb{R}^n \end{array}$$

- Apply steepest descent to this problem, multiply with \mathbf{S} , and pass back to the space of \mathbf{x} , using $\nabla h(\mathbf{y}^r) = \mathbf{S} \nabla f(\mathbf{x}^r)$,

$$\begin{aligned} \mathbf{y}^{r+1} &= \mathbf{y}^r - \alpha_r \nabla h(\mathbf{y}^r) \\ \mathbf{S}\mathbf{y}^{r+1} &= \mathbf{S}\mathbf{y}^r - \alpha_r \mathbf{S} \nabla h(\mathbf{y}^r) \\ \mathbf{x}^{r+1} &= \mathbf{x}^r - \alpha_r \mathbf{D}^r \nabla f(\mathbf{x}^r) \end{aligned}$$

Diagonal Scaling

- Apply the results for steepest descent to the scaled iteration $\mathbf{y}^{r+1} = \mathbf{y}^r - \alpha_r \nabla h(\mathbf{y}^r)$:

$$\frac{\|\mathbf{y}^{r+1}\|}{\|\mathbf{y}^r\|} \leq \max\{|1 - \alpha_r m_r|, |1 - \alpha_r M_r|\}$$

$$\frac{f(\mathbf{x}^{r+1})}{f(\mathbf{x}^r)} = \frac{h(\mathbf{y}^{r+1})}{h(\mathbf{y}^r)} \leq \left(\frac{M_r - m_r}{M_r + m_r} \right)^2$$

where m_r and M_r are the smallest and largest eigenvalues of the Hessian of h , which is

$$\nabla^2 h(\mathbf{y}) = \mathbf{S} \nabla^2 f(\mathbf{x}) \mathbf{S} = (\mathbf{D}^r)^{1/2} \mathbf{Q} (\mathbf{D}^r)^{1/2}$$

- It is desirable to choose \mathbf{D}^r as close as possible to \mathbf{Q}^{-1} . Also if \mathbf{D}^r is so chosen, the stepsize $\alpha = 1$ is near the optimal $2/(M_r + m_r)$.
- Using as \mathbf{D}^r a diagonal approximation to \mathbf{Q}^{-1} is common and often very effective. Corrects for poor choice of units expressing the variables.

Non-quadratic Problems

- Rate of convergence to a nonsingular local minimum of a non-quadratic function is very similar to the quadratic case (linear convergence is typical).
- If $\mathbf{D}^r = (\nabla^2 f(\mathbf{x}^*))^{-1}$, we asymptotically obtain optimal scaling and superlinear convergence
- More generally, if the direction $\mathbf{d}^r = -\mathbf{D}^r \nabla f(\mathbf{x}^r)$ approaches asymptotically the Newton direction, i.e.,

$$\lim_{r \rightarrow \infty} \frac{\|\mathbf{d}^r + (\nabla^2 f(\mathbf{x}^*))^{-1} \nabla f(\mathbf{x}^r)\|}{\|\nabla f(\mathbf{x}^r)\|} = 0$$

and the Armijo rule is used with initial stepsize equal to one, the rate of convergence is superlinear.

- Convergence rate to a singular local min is typically sublinear (in effect, condition number $\kappa = \infty$)

Dealing with Nonnegative Constraints

$$\begin{array}{ll} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \geq \mathbf{0}. \end{array}$$

- Convert to an unconstrained optimization problem: $x_i = w_i^2$, $i = 1, 2, \dots, m$.
- Caution: this can create spurious solutions.
- Consider the example

$$\begin{array}{ll} \text{minimize} & (x_1 - 1)^2 + (x_2 + 1)^2 \\ \text{subject to} & x_1 \geq 0. \end{array}$$

- Optimal solution $\mathbf{x}^* = (1, -1)^T$. Conversion to an unconstrained problem:

$$\begin{array}{ll} \text{minimize} & (w_1^2 - 1)^2 + (w_2 + 1)^2 \\ \text{subject to} & \mathbf{w} \in \mathbb{R}^2. \end{array}$$

- We have

$$\nabla f = \begin{pmatrix} 4w_1(w_1^2 - 1) \\ 2(w_2 + 1) \end{pmatrix}, \quad \nabla^2 f = \begin{bmatrix} 4(3w_1^2 - 1) & 0 \\ 0 & 2 \end{bmatrix}$$

which introduces a saddle point $(0, -1)^T$.

- We will discuss how to deal with inequality constraints later.