

Unconstrained minimization: gradient descent

15.093: Optimization

Dimitris Bertsimas
Alexandre Jacquillat

Sloan School of Management
Massachusetts Institute of Technology



Unconstrained minimization

Formulation (Unconstrained minimization)

Let $f: \mathbb{R}^n \mapsto \mathbb{R}$ be a continuous (usually differentiable) function. The associated unconstrained non-linear optimization problem is given by

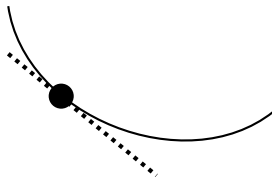
$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

- Notation: $z^* = \min f(\mathbf{x})$ and $\mathbf{x}^* \in \arg \min f(\mathbf{x})$
- Focus on the difficulties at the core of non-linear optimization
 - The techniques used for unconstrained non-linear minimization also lie at the core of algorithms for constrained non-linear minimization
- Sample applications of unconstrained non-linear minimization
 - Machine learning: linear regression, logistic regression, neural networks
 - Signal processing: interpolation, extrapolation, de-noising
 - Robotics: position, orientation, inverse kinematics

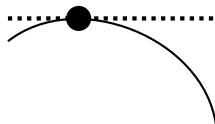
Necessary and sufficient optimality conditions

Some intuition in one dimension

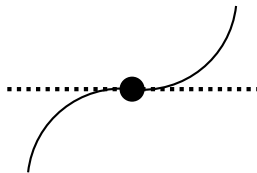
Case with $f'(x) \neq 0$



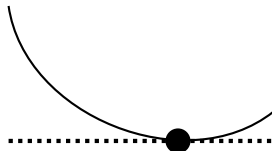
Case with $f'(x) = 0$ and $f''(x) < 0$



Case with $f'(x) = 0$ and $f''(x) = 0$



Case with $f'(x) = 0$ and $f''(x) > 0$



- Necessity of first-order condition at a local optimum: $\nabla f(\mathbf{x}^*) = 0$
- Necessity of additional condition to guarantee local optimality, e.g., second-order condition convexity

Necessary optimality conditions

Theorem

Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be a twice continuously differentiable function.

If $\mathbf{x}^* \in \mathbb{R}^n$ is a local minimum of $f(\cdot)$, then

$$\nabla f(\mathbf{x}^*) = \mathbf{0} \text{ and } \nabla^2 f(\mathbf{x}^*) \succeq \mathbf{0}$$

- Proof of $\nabla f(\mathbf{x}^*) = \mathbf{0}$ (first-order necessary condition)
 - For all $\mathbf{d} \in \mathbb{R}^n, \alpha > 0$, $\frac{f(\mathbf{x}^* + \alpha \mathbf{d}) - f(\mathbf{x}^*)}{\alpha} \geq 0$
 - Take limits as $\alpha \rightarrow 0$: for all $\mathbf{d} \in \mathbb{R}^n$: $\nabla f(\mathbf{x}^*)^\top \mathbf{d} \geq 0$
 - Apply the inequality to $\mathbf{e}_i \in \mathbb{R}^n$ and $-\mathbf{e}_i \in \mathbb{R}^n$: $\nabla_i f(\mathbf{x}^*) = 0$ for all i
- Proof of $\nabla^2 f(\mathbf{x}^*) \succeq \mathbf{0}$ (second-order necessary condition)
 - Taylor series expansion: $\rho(\cdot)$ with $\rho(\mathbf{y}) \rightarrow 0$ as $\mathbf{y} \rightarrow \mathbf{0}$, such that

$$\underbrace{f(\mathbf{x}^* + \alpha \mathbf{d}) - f(\mathbf{x}^*)}_{\geq 0} = \underbrace{\alpha \nabla f(\mathbf{x}^*)^\top \mathbf{d}}_{=0} + \frac{\alpha^2}{2} \mathbf{d}^\top \nabla^2 f(\mathbf{x}^*) \mathbf{d} + \alpha^2 \|\mathbf{d}\|^2 \rho(\alpha \mathbf{d})$$

$$\alpha \rightarrow 0 : \implies \mathbf{d}^\top \nabla^2 f(\mathbf{x}^*) \mathbf{d} \geq 0$$

Sufficient optimality conditions

Definition (Positive definite matrix)

A symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is positive definite, written $\mathbf{A} \succ \mathbf{0}$, if $\mathbf{u}^\top \mathbf{A} \mathbf{u} > 0$ for all $\mathbf{u} \neq \mathbf{0}$.

Equivalently, $\mathbf{A} \succ \mathbf{0}$ if and only if all its eigenvalues are positive.

Theorem

Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be twice continuously differentiable.

If $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$, then \mathbf{x}^* is a local minimum.

- Proof:

- Let $\lambda > 0$ be the smallest eigenvalue of $\nabla^2 f(\mathbf{x})$
- Taylor series expansion: $\rho(\cdot)$ with $\rho(\mathbf{y}) \rightarrow 0$ as $\mathbf{y} \rightarrow \mathbf{0}$, such that

$$\begin{aligned} f(\mathbf{x}^* + \mathbf{y}) - f(\mathbf{x}^*) &= \underbrace{\nabla f(\mathbf{x}^*)^\top}_{=0} \mathbf{y} + \frac{1}{2} \underbrace{\mathbf{y}^\top \nabla^2 f(\mathbf{x}^*) \mathbf{y}}_{\geq \lambda \|\mathbf{y}\|^2} + \|\mathbf{y}\|^2 \rho(\mathbf{y}) \\ &\geq \frac{\lambda}{2} \|\mathbf{y}\|^2 + \|\mathbf{y}\|^2 \rho(\mathbf{y}) \\ &\geq 0, \text{ for } \|\mathbf{y}\| \text{ small enough.} \end{aligned}$$

Example

- Consider the following two-dimensional function:

$$f(\mathbf{x}) = \frac{1}{2}x_1^2 + x_1x_2 + 2x_2^2 - 4x_1 - 4x_2 - x_2^3$$

- Compute its gradient and its Hessian:

$$\nabla f(\mathbf{x}) = (x_1 + x_2 - 4, x_1 + 4x_2 - 4 - 3x_2^2)^\top$$

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} 1 & 1 \\ 1 & 4 - 6x_2 \end{bmatrix}$$

→ Two candidates for local minima: $\mathbf{x}^* = (4, 0)$ and $\bar{\mathbf{x}} = (3, 1)$

- $\bar{\mathbf{x}}$ is not a local minimum from second-order information

$$\nabla^2 f(\bar{\mathbf{x}}) = \begin{bmatrix} 1 & 1 \\ 1 & -2 \end{bmatrix} \text{ is an indefinite matrix}$$

- \mathbf{x}^* is a local minimum:

$$\nabla^2 f(\mathbf{x}^*) = \begin{bmatrix} 1 & 1 \\ 1 & 4 \end{bmatrix} \succ \mathbf{0}$$

Another example

- Consider the following two-dimensional function:

$$f(\mathbf{x}) = x_1^3 + x_2^2$$

- Compute its gradient and its Hessian:

$$\nabla f(\mathbf{x}) = (3x_1^2, 2x_2)^\top$$

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} 6x_1 & 0 \\ 0 & 2 \end{bmatrix}$$

→ One candidate for local minimum: $\mathbf{x}^* = (0, 0)$

- \mathbf{x}^* is still a candidate a local minimum from second-order information

$$\nabla^2 f(\mathbf{x}^*) = \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix} \succeq \mathbf{0}$$

- However, \mathbf{x}^* is not a local minimum.

$$\tilde{\mathbf{x}} = (-\varepsilon, 0)^\top \implies f(\tilde{\mathbf{x}}) = -\varepsilon^3 < 0 = f(\mathbf{x}^*)$$

The case of convexity

Proposition (Reminder: first-order convexity condition)

Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be a differentiable function. It is convex if and only:

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}), \quad \forall \mathbf{x}, \mathbf{y}$$

Proposition (Reminder: second-order condition)

Let $f(\cdot)$ be a twice continuously differentiable function.

f is a convex function if and only if $H(\mathbf{x}) \succeq \mathbf{0}$ for all $\mathbf{x} \in \mathbb{R}^n$

Theorem

Let $f(\mathbf{x})$ be a continuously differentiable convex function.

Then \mathbf{x}^* is a global minimum of f if and only if $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

- Proof: If f is convex and $\nabla f(\mathbf{x}^*) = \mathbf{0}$, we have:

$$f(\mathbf{x}) - f(\mathbf{x}^*) \geq \nabla f(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*) = 0$$

Example: fitting a linear regression model

- Linear regression model with ℓ_2 -regularization: $\min f(\beta)$ with

$$f(\beta) = \frac{1}{2} \sum_{i=1}^n \left(y_i - \sum_{j=1}^m x_{ij} \beta_j \right)^2 + \frac{\lambda}{2} \sum_{j=1}^m \beta_j^2 = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \frac{\lambda}{2} \|\beta\|^2$$

- Differentiation of the loss function:

$$\nabla f(\beta) = -\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\beta) + \lambda\beta$$

$$\nabla^2 f(\beta) = \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}$$

- Importance of ℓ_2 -regularization in case \mathbf{X} does not have full rank
 - $\mathbf{X}^\top \mathbf{X}$ is positive semidefinite: $\mathbf{u}^\top \mathbf{X}^\top \mathbf{X} \mathbf{u} = \|\mathbf{X}\mathbf{u}\|^2 \geq 0, \forall \mathbf{u} \in \mathbb{R}^n$
 - $\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}$ is positive definite if $\lambda > 0$:
 $\mathbf{u}^\top (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}) \mathbf{u} = \|\mathbf{X}\mathbf{u}\|^2 + \lambda \|\mathbf{u}\|^2 > 0, \forall \mathbf{u} \neq \mathbf{0}$

→ A convex quadratic optimization problem: $\nabla^2 f(\beta) \succeq \mathbf{0}$

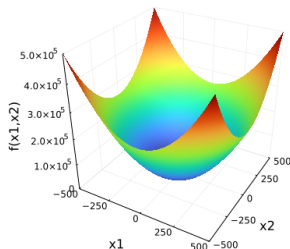
→ The problem admits a minimum at a stationary point

$$\beta^* = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

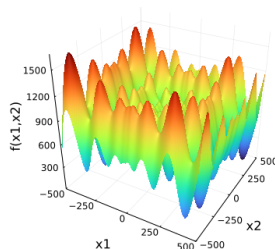
Summary and implications

- A local minimum can only occur in a stationary point: $\nabla f(\mathbf{x}^*) = \mathbf{0}$
- Convex optimization: any stationary point is a global minimum!
- Otherwise, local optimality depends on second-order conditions
 - If $\nabla^2 f(\mathbf{x}^*) \succ 0$, then \mathbf{x}^* is a local minimum
 - If $\nabla^2 f(\mathbf{x}^*) \succeq 0$, then \mathbf{x}^* *can be* a local minimum
 - If $\nabla^2 f(\mathbf{x}^*)$ is indefinite, then \mathbf{x}^* is not a local minimum
- In non-convex optimization, identifying global minima is very hard

Convex optimization



Non-convex optimization



Algorithms for unconstrained optimization

A generic descent method

- Seek a stationary point as a candidate for local minimum

$$\nabla f(\mathbf{x}) = \mathbf{0}$$

Algorithm

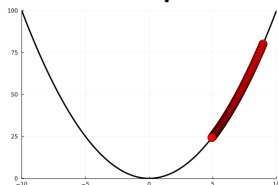
1. *Initialization: starting point $\mathbf{x}^0 \in \mathbb{R}^n$, and iteration counter $k = 0$*
 2. *Repeat, until termination criterion is reached (e.g., $\|\nabla f(\mathbf{x}^k)\| \leq \eta$)*
 - 2.1 *Update iteration counter: $k \leftarrow k + 1$*
 - 2.2 *Determine a descent direction \mathbf{d}^k , such that $\nabla f(\mathbf{x}^k)^\top \mathbf{d}^k < 0$*
 - 2.3 *Determine a step size $\alpha^k > 0$*
 - 2.4 *Update $\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k + \alpha^k \mathbf{d}^k$*
- A descent algorithm: improvement in the solution at each iteration as long as the step size is small enough and $\nabla f(\mathbf{x}^k)^\top \mathbf{d}^k < 0$

$$f(\mathbf{x}^{k+1}) \approx f(\mathbf{x}^k) + \alpha^k \nabla f(\mathbf{x}^k)^\top \mathbf{d} < f(\mathbf{x}^k), \quad \forall k = 0, 1, \dots$$

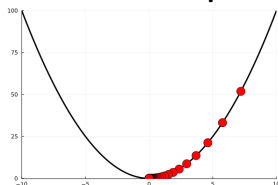
Impact of step size on algorithm convergence

- Core trade-off in setting step sizes
 - Small step sizes: slow convergence, with little progress per iteration
 - Large step sizes: unstable behavior
 - In-between, faster convergence can be achieved
- The appropriate step size depends on the shape of the function to minimize, the descent direction, and the progress of the algorithm

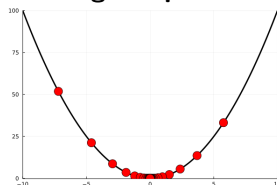
Small step size



Intermediate step size



Large step size



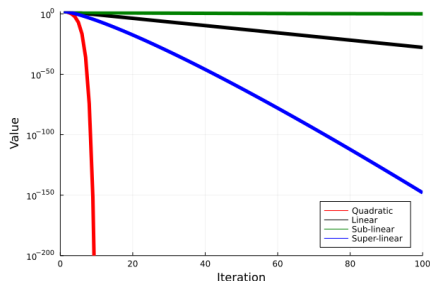
Rate of convergence of algorithms

Definition

A convergence sequence $z_1, \dots, z_n \rightarrow z$ has order of convergence p and rate of convergence β if

$$\lim_{k \rightarrow \infty} \frac{|z_{k+1} - z|}{|z_k - z|^p} = \beta$$

- Linear convergence: $p^* = 1$, $0 < \beta < 1$, e.g., $z_k = a^k$, $a < 1$
- Sub-linear convergence: $p^* = 1$, $\beta = 1$, e.g., $z_k = 1/k$
- Super-linear convergence: $p^* = 1$, $\beta = 0$, e.g., $z_k = (1/k)^k$
- Quadratic convergence: $p^* = 2$, e.g., $z_k = a^{2^k}$, $a < 1$



Gradient descent: convergence analysis

Gradient descent algorithm

Algorithm

1. *Initialization: starting point $\mathbf{x}^0 \in \mathbb{R}^n$, and iteration counter $k = 0$*
2. *Repeat, until stopping criterion is reached*
 - 2.1 *Update iteration counter: $k \leftarrow k + 1$*
 - 2.2 *Choose descent direction $\mathbf{d}^k = -\nabla f(\mathbf{x}^k)$*
 - 2.3 *Determine a step size $\alpha^k > 0$*
 - 2.4 *Update $\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k + \alpha^k \mathbf{d}^k$*

- Motivation: the gradient rule defines a valid descent direction

$$f(\mathbf{x}^{k+1}) \approx f(\mathbf{x}^k) - \alpha \|\nabla f(\mathbf{x}^k)\|^2 < f(\mathbf{x}^k)$$

→ Gradient descent defines a monotonically decreasing sequence

$$f(\mathbf{x}^0) > f(\mathbf{x}^1) > \dots > f(\mathbf{x}^k) > \dots$$

- Under assumptions, $\mathbf{x}^0, \mathbf{x}^1, \dots$, will converge to a stationary point

Further assumptions: smoothness and strong convexity

Definition (smoothness)

A function $f(\cdot)$ is M -smooth if its gradient is M -Lipschitz continuous:

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq M\|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y}$$

Definition (strong convexity)

A function $f(\cdot)$ is m -strongly convex if:

$$f(\mathbf{x}) - \frac{m}{2}\|\mathbf{x}\|^2 \text{ is convex.}$$

Proposition

If f is M -smooth, it satisfies:

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{M}{2}\|\mathbf{y} - \mathbf{x}\|^2, \quad \forall \mathbf{x}, \mathbf{y}$$

If f is m -strongly convex and continuously differentiable, it satisfies

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{m}{2}\|\mathbf{y} - \mathbf{x}\|^2, \quad \forall \mathbf{x}, \mathbf{y}$$

Gradient descent with constant step size: convergence

Theorem (Convergence with constant step size)

Assume that f is M -smooth and convex. The gradient descent algorithm with fixed step size $\alpha \leq 1/M$ converges in $\mathcal{O}(1/k)$ to a global minimum:

$$f(\mathbf{x}^k) - z^* \leq \frac{1}{2\alpha k} \|\mathbf{x}^0 - \mathbf{x}^*\|^2, \quad \forall k \geq 0$$

- Sub-linear convergence: $f(\mathbf{x}^k) - z^* \leq \frac{1}{2\alpha k} \|\mathbf{x}^0 - \mathbf{x}^*\|^2$
- Solution within ε of the optimum after at most

$$\frac{\|\mathbf{x}^0 - \mathbf{x}^*\|^2}{2\alpha\varepsilon} \text{ iterations}$$

→ Drivers of algorithm performance

- Initial condition: the weaker \mathbf{x}^0 , the slower the convergence
- Tolerance: the smaller ε , the slower the convergence
- Step size: convergence is faster when α is larger, yet not too large

Gradient descent with constant step size: proof

- Using $\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha \nabla f(\mathbf{x}^k)$ and M -smoothness, we have:

$$\begin{aligned}
 f(\mathbf{x}^{k+1}) &\leq f(\mathbf{x}^k) + \underbrace{\nabla f(\mathbf{x}^k)^\top (\mathbf{x}^{k+1} - \mathbf{x}^k)}_{-\alpha \|\nabla f(\mathbf{x}^k)\|^2} + \underbrace{\frac{M}{2}}_{\leq 1/(2\alpha)} \cdot \underbrace{\|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2}_{=\alpha^2 \|\nabla f(\mathbf{x}^k)\|^2} \\
 \implies f(\mathbf{x}^{k+1}) &\leq f(\mathbf{x}^k) - \frac{\alpha}{2} \|\nabla f(\mathbf{x}^k)\|^2 = f(\mathbf{x}^k) - \frac{1}{2\alpha} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2
 \end{aligned}$$

- By convexity: $f(\mathbf{x}^*) \geq f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^\top (\mathbf{x}^* - \mathbf{x}^k)$. We obtain

$$\begin{aligned}
 f(\mathbf{x}^{k+1}) - f(\mathbf{x}^*) &\leq -\nabla f(\mathbf{x}^k)^\top (\mathbf{x}^* - \mathbf{x}^k) - \frac{1}{2\alpha} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 \\
 &= \dots \text{ [after some algebra]} \\
 &= \frac{1}{2\alpha} \left(\|\mathbf{x}^k - \mathbf{x}^*\|^2 - \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 \right)
 \end{aligned}$$

- We conclude by telescoping the sum and exploiting monotonicity:

$$\underbrace{\sum_{i=1}^k (f(\mathbf{x}^i) - f(\mathbf{x}^*))}_{=k(f(\mathbf{x}^k) - f(\mathbf{x}^*))} \leq \sum_{i=1}^k (f(\mathbf{x}^i) - f(\mathbf{x}^*)) \leq \underbrace{\frac{1}{2\alpha} \left(\|\mathbf{x}^0 - \mathbf{x}^*\|^2 - \|\mathbf{x}^k - \mathbf{x}^*\|^2 \right)}_{\leq \frac{1}{2\alpha} \|\mathbf{x}^0 - \mathbf{x}^*\|^2}$$

Gradient descent with exact line search: convergence

- Gradient descent with exact line search: at each iteration k , choose step size α^k to maximize the one-step improvement:

$$\alpha^k \in \arg \min f(\mathbf{x}^k + \alpha^k \mathbf{d}^k)$$

Theorem (Convergence with exact line search)

Assume that f is M -smooth and m -strongly convex. The gradient descent algorithm with exact line search converges in $\mathcal{O}(c^k)$, with $c = 1 - m/M$

$$f(\mathbf{x}^k) - z^* \leq c^k (f(\mathbf{x}^{(0)}) - z^*), \quad \forall k \geq 0$$

- Stronger convergence with f strongly convex and optimized step sizes
 - Linear convergence: $f(\mathbf{x}^k) - z^* \leq c^k (f(\mathbf{x}^{(0)}) - z^*)$
 - Solution within ε of the optimum after at most

$$\frac{\log(f(\mathbf{x}^{(0)}) - z^*) - \log(\varepsilon)}{\log(1/c)} \text{ iterations}$$

Gradient descent with exact line search: proof

- Due to the M -smoothness of f :

$$f(\mathbf{x}^k - \alpha^k \nabla f(\mathbf{x}^k)) \leq f(\mathbf{x}^k) - \alpha^k \|\nabla f(\mathbf{x}^k)\|^2 + \frac{M(\alpha^k)^2}{2} \|\nabla f(\mathbf{x}^k)\|^2$$

- By optimizing over α^k and subtracting z^* on both sides, we obtain:

$$(f(\mathbf{x}^{k+1}) - z^*) \leq (f(\mathbf{x}^k) - z^*) - \frac{1}{2M} \|\nabla f(\mathbf{x}^k)\|^2$$

- Due to the strong convexity of f , we have

$$f(\mathbf{y}) \geq f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^\top (\mathbf{y} - \mathbf{x}^k) + \frac{m}{2} \|\mathbf{y} - \mathbf{x}^k\|^2, \quad \forall \mathbf{y}$$

- By minimizing over \mathbf{y} on both sides, we obtain:

$$z^* \geq f(\mathbf{x}^k) - \frac{1}{2m} \|\nabla f(\mathbf{x}^k)\|^2$$

- We conclude:

$$(f(\mathbf{x}^{k+1}) - z^*) \leq \left(1 - \frac{m}{M}\right) (f(\mathbf{x}^k) - z^*)$$

$$\implies (f(\mathbf{x}^k) - z^*) \leq \left(1 - \frac{m}{M}\right)^k (f(\mathbf{x}^{(0)}) - z^*)$$

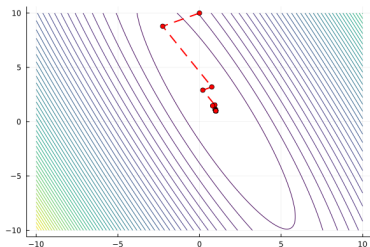
Example

$$\min f(x_1, x_2) = 5x_1^2 + x_2^2 + 4x_1x_2 - 14x_1 - 6x_2 + 20$$

$$\begin{pmatrix} d_1^k \\ d_2^k \end{pmatrix} = -\nabla f(x_1^k, x_2^k) = \begin{pmatrix} -10x_1^k - 4x_2^k + 14 \\ -2x_2^k - 4x_1^k + 6 \end{pmatrix}$$

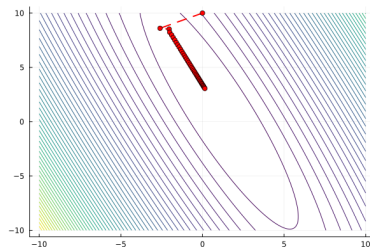
Exact line search

$$\alpha^k = \frac{(d_1^k)^2 + (d_2^k)^2}{2(5(d_1^k)^2 + (d_2^k)^2 + 4d_1^k d_2^k)}$$



Constant step size

$$\alpha^k = 0.1$$



Advanced topics

Role of the condition number in quadratic optimization

Definition (condition number)

The condition number of a non-singular matrix \mathbf{Q} is the ratio of its largest-to-smallest eigenvalues: $\kappa(\mathbf{Q}) = \frac{\lambda_{\max}}{\lambda_{\min}} \geq 1$

Theorem (Kantorovich, Nobel Prize in Economics, 1975)

If $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{c}^\top \mathbf{x} + b$ and $\mathbf{Q} \succ 0$, then

$$(f(\mathbf{x}^{k+1}) - z^*) \leq \left(\frac{\kappa(\mathbf{Q}) - 1}{\kappa(\mathbf{Q}) + 1} \right)^2 (f(\mathbf{x}^k) - z^*)$$

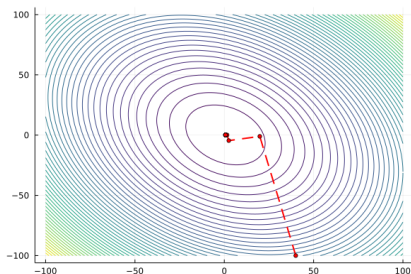
- The condition number plays a critical role in the convergence of gradient descent algorithms for quadratic functions
 - $\kappa(\mathbf{Q}) \approx 1$: “well-conditioned” matrix, fast convergence
 - $\kappa(\mathbf{Q}) \gg 1$: “ill-conditioned” matrix, slow convergence

Role of the condition number: illustration

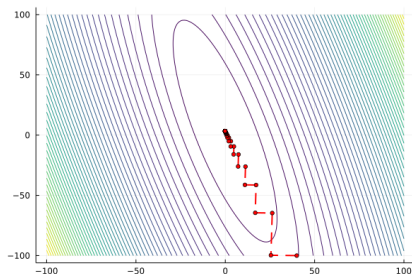
$$\min f(x_1, x_2) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{c}^\top \mathbf{x} + 10 \quad \mathbf{d}^k = -\nabla f(\mathbf{x}^k) = -\mathbf{Q} \mathbf{x}^k + \mathbf{c}$$

$$\text{Exact line search: } \alpha^k = \frac{(\mathbf{d}^k)^\top \mathbf{d}^k}{(\mathbf{d}^k)^\top \mathbf{Q} \mathbf{d}^k}$$

$$\mathbf{Q} = \begin{bmatrix} 20 & 5 \\ 5 & 16 \end{bmatrix} \rightarrow \kappa(\mathbf{Q}) = 1.85$$



$$\mathbf{Q} = \begin{bmatrix} 20 & 5 \\ 5 & 2 \end{bmatrix} \rightarrow \kappa(\mathbf{Q}) = 30.23$$



Beyond quadratic optimization

- Recall the drivers of the performance of the gradient descent algorithm
 - Initial condition: the weaker x^0 , the slower the convergence
 - Tolerance: the smaller ε , the slower the convergence
 - Function f : the higher the constant c , the slower the convergence
- The constant $c = 1 - \frac{m}{M}$ relates to the convex shape of the function f
- In fact, c depends on the condition number of the Hessian of f around the optimal value of the optimization problem

$$\kappa(\nabla^2 f(x^*)) \lesssim \frac{M}{m}$$

- Beyond quadratic optimization: the convergence rate is impacted by the condition number of the Hessian matrix
- Strong convergence when the Hessian is well-conditioned
 - Slow convergence when the Hessian is ill-conditioned

Extension: norm-based steepest gradient descent

Definition (steepest descent)

For any norm $\|\cdot\|$, the steepest descent direction is given as follows:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha^k \mathbf{d}^k, \text{ with } \mathbf{d}^k \in \arg \min \{ \nabla f(\mathbf{x}^k) \mathbf{v} : \|\mathbf{v}\| = 1 \}$$

- Motivation: finding the direction with the strongest improvement

$$f(\mathbf{x} + \mathbf{v}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^\top \mathbf{v}$$

- ℓ_2 -norm $\|\cdot\|_2 \implies \mathbf{d}^k = -\nabla f(\mathbf{x}^k)$
 - Cauchy-Schwarz inequality: $|\nabla f(\mathbf{x}^k) \mathbf{v}| \leq \|\nabla f(\mathbf{x}^k)\| \|\mathbf{v}\| = \|\nabla f(\mathbf{x}^k)\|$
 - Gradient descent is merely steepest descent with the ℓ_2 -norm
- Quadratic norm $\|\mathbf{z}\|_P = \sqrt{\mathbf{z}^\top \mathbf{P} \mathbf{z}} \implies \mathbf{d}^k = -\mathbf{P}^{-1} \nabla f \mathbf{x}$
 - Equivalent to gradient descent with re-scaling $\bar{\mathbf{x}} = \mathbf{P}^{1/2} \mathbf{x}$ to circumvent ill-conditioned problems
- ℓ_1 -norm $\|\cdot\|_1 \implies \mathbf{d}^k = -\frac{\partial f}{\partial x_i} \mathbf{e}_i$
 - Equivalent to coordinate descent: one variable at a time

Conclusion

Summary

Takeaway

Unconstrained minimization is the core problem in non-linear optimization, with applications in machine learning, signal processing, robotics, etc.

Takeaway

Descent methods can find local optima. They imply global optimality if the function is convex; global optimization is much harder otherwise.

Takeaway

Gradient descent is a simple algorithm that exhibits sub-linear convergence with constant step sizes and linear convergence with optimized step sizes.

Takeaway

Convergence is greatly impacted by the condition number of $\nabla^2 f(\mathbf{x}^)$.*