

Unconstrained minimization: Newton's method

15.093: Optimization

Dimitris Bertsimas
Alexandre Jacquillat

Sloan School of Management
Massachusetts Institute of Technology



Reminder: descent methods for unconstrained minimization

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

Algorithm

1. *Initialization: starting point $\mathbf{x}^0 \in \mathbb{R}^n$, and iteration counter $k = 0$*
 2. *Repeat, until termination criterion is reached*
 - 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 $\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k + \alpha^k \mathbf{d}^k$*
- Main design questions
 1. Initialization: how to determine the starting point \mathbf{x}^0 ?
 2. Descent: how to determine the descent direction \mathbf{d}^k ?
 3. Line search: how to choose the step size α^k ?
 4. Termination criterion; typically, $\|\nabla f(\mathbf{x}^k)\| \leq \eta$ for small $\eta > 0$

Newton's method

Motivation and intuition: a second-order view

- Taylor series expansion around x

$$f(\mathbf{y}) \approx \hat{f}(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^\top \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x})$$

Motivation and intuition: a second-order view

- Taylor series expansion around \mathbf{x}

$$f(\mathbf{y}) \approx \hat{f}(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^\top \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x})$$

- Approximation of the minimization problem: $\min f(\mathbf{y}) \rightarrow \min \hat{f}(\mathbf{y})$
- From first-order conditions, move in the Newton direction

$$\begin{aligned}\nabla \hat{f}(\mathbf{y}) = \mathbf{0} &\implies \nabla f(\mathbf{x}) + \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x}) = \mathbf{0} \\ &\implies \mathbf{y} = \mathbf{x} - (\nabla^2 f(\mathbf{x}))^{-1} \nabla f(\mathbf{x})\end{aligned}$$

- The Newton decrement approximates the algorithm's progress:

$$f(\mathbf{x}) - \min_{\mathbf{y}} \hat{f}(\mathbf{y}) = f(\mathbf{x}) - \hat{f}(\mathbf{x} + \mathbf{d}) = \frac{1}{2} \lambda(\mathbf{x})^2$$

Definition

- Newton direction: $\mathbf{d} = -(\nabla^2 f(\mathbf{x}))^{-1} \nabla f(\mathbf{x})$
- Newton decrement: $\lambda(\mathbf{x}) = (\nabla f(\mathbf{x})^\top \nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x}))^{1/2}$

Motivation and intuition: a first-order view

- We seek a stationary point: $\nabla f(\mathbf{y}) = \mathbf{0}$
- Taylor series expansion of $\nabla f(\mathbf{y})$ around \mathbf{x}

$$\nabla f(\mathbf{y}) \approx \hat{g}(\mathbf{y}) = \nabla f(\mathbf{x}) + \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x})$$

Motivation and intuition: a first-order view

- We seek a stationary point: $\nabla f(\mathbf{y}) = \mathbf{0}$
- Taylor series expansion of $\nabla f(\mathbf{y})$ around \mathbf{x}

$$\nabla f(\mathbf{y}) \approx \hat{g}(\mathbf{y}) = \nabla f(\mathbf{x}) + \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x})$$

- Approximation of the problem: $\nabla f(\mathbf{y}) = \mathbf{0} \rightarrow \hat{g}(\mathbf{y}) = \mathbf{0}$
- From the equation, move in the Newton direction

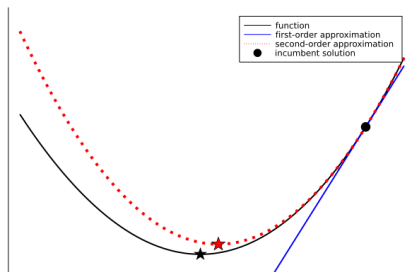
$$\begin{aligned}\hat{g}(\mathbf{y}) = \mathbf{0} &\implies \nabla f(\mathbf{x}) + \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x}) = \mathbf{0} \\ &\implies \mathbf{y} = \mathbf{x} - (\nabla^2 f(\mathbf{x}))^{-1} \nabla f(\mathbf{x})\end{aligned}$$

- History:
 - The Newton method was developed by Newton and Raphson in the 1600's for solving systems of equations
 - Extension to optimization by Simpson in the 1700's: $g(\mathbf{y}) = \nabla f(\mathbf{y})$

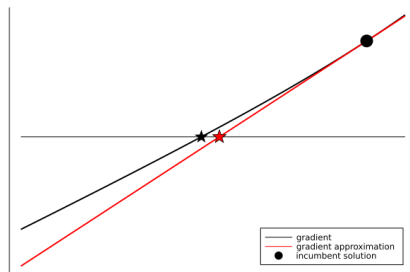
Visualization

- Two equivalent ways of interpreting Newton's method
 1. Second-order view: minimization of second-order Taylor approximation
 2. First-order view: root of first-order Taylor approximation of gradient
- By leveraging second-order (Hessian) information, we obtain a stronger approximation of the function, hence of the minimization problem

Second-order view



First-order view



Newton's method

Algorithm (Newton's method)

1. *Initialization: starting point $\mathbf{x}^0 \in \mathbb{R}^n$, and iteration counter $k = 0$*
2. *Repeat, until termination criterion is reached*
 - 2.1 *Update iteration counter: $k \leftarrow k + 1$*
 - 2.2 *Determine Newton's direction $\mathbf{d}^k = -(\nabla^2 f(\mathbf{x}^k))^{-1} \nabla f(\mathbf{x}^k)$*
 - 2.3 *Update $\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k + \mathbf{d}^k$*

Newton's method

Algorithm (Newton's method)

1. *Initialization: starting point $\mathbf{x}^0 \in \mathbb{R}^n$, and iteration counter $k = 0$*
2. *Repeat, until termination criterion is reached*
 - 2.1 *Update iteration counter: $k \leftarrow k + 1$*
 - 2.2 *Determine Newton's direction $\mathbf{d}^k = -(\nabla^2 f(\mathbf{x}^k))^{-1} \nabla f(\mathbf{x}^k)$*
 - 2.3 *Update $\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k + \mathbf{d}^k$*

- A second-order method: use of first-order gradient and second-order Hessian information to proceed from iteration to iteration
- More progress at each iteration than gradient descent
 - More work per iteration: $\mathcal{O}(n^3)$ operations
 - “Pure” Newton method relies on a step size of 1
 - Convergence criterion based on estimated improvement: $\lambda(\mathbf{x})^2/2 \leq \varepsilon$
 - Affine invariance: Newton's method independent of problem scaling
 - Newton's method for $f(\mathbf{x})$ and $\tilde{f}(\mathbf{y}) = f(\mathbf{T}\mathbf{y})$ yields $\mathbf{x}^{k+1} = \mathbf{T}\mathbf{y}^{k+1}$
 - Recall that this is not the case with gradient descent

Example: fitting a logistic regression model

$$\max f(\beta) = \sum_{i=1}^n \left\{ y_i \mathbf{x}_i^\top \beta - \log \left(1 + e^{\mathbf{x}_i^\top \beta} \right) \right\}$$

$$\nabla f(\beta) = \mathbf{X}^\top (\mathbf{y} - \mathbf{p}(\mathbf{X}, \beta)), \quad \text{with: } p_i(\mathbf{X}, \beta) = \frac{e^{\mathbf{x}_i^\top \beta}}{1 + e^{\mathbf{x}_i^\top \beta}}$$

$$\nabla^2 f(\beta) = -\mathbf{X}^\top \mathbf{W} \mathbf{X}, \quad \text{where: } \mathbf{W} = \text{diag} \left(\frac{e^{\mathbf{x}_i^\top \beta}}{1 + e^{\mathbf{x}_i^\top \beta}} \cdot \frac{1}{1 + e^{\mathbf{x}_i^\top \beta}} \right)$$

- A convex optimization problem: $\max f(\beta)$ with $\nabla^2 f(\beta) \preceq 0$
- Minimum at a stationary point $\mathbf{X}^\top (\mathbf{y} - \mathbf{p}(\mathbf{X}, \beta^*)) = 0$
- Applying Newton's algorithm to the logistic regression model:

$$\begin{aligned} \beta^{k+1} &= \beta^k + (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{p}(\mathbf{X}, \beta^k)) \\ &= (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{z}^k, \quad \text{with } \mathbf{z}^k = \mathbf{X} \beta^k + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p}(\mathbf{X}, \beta^k)) \end{aligned}$$

→ Interpretation as iterative re-weighted least squares:

$$\beta^{k+1} = \arg \min_{\beta} (\mathbf{z}^k - \mathbf{X} \beta)^\top \mathbf{W} (\mathbf{z}^k - \mathbf{X} \beta)$$

Local convergence

Local convergence of Newton's method

Definition (operator norm of a matrix)

$$\|M\| = \max\{\|Mx\| : \|x\| = 1\}$$

Theorem

$f(\cdot)$ twice continuously differentiable, and $\nabla f(x^*) = \mathbf{0}$. Assume that:

- $\|(\nabla^2 f(x^*))^{-1}\| \leq \frac{1}{m}$ for some $m > 0$
- $\nabla^2 f(x)$ is L -Lipschitz in the β -ball around x^* for some $\beta > 0$, $L > 0$:

$$\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq L\|x - y\|, \forall x, y \in \mathcal{B}(x^*, \beta)$$

Define $\delta = \min\left\{\beta, \frac{2m}{3L}\right\}$. The following holds:

1. If $\|x^k - x\| \leq \delta$, then $\|x^{k+1} - x\| \leq \delta$ for all $k = 0, 1, 2, \dots$
2. $\|x^{k+1} - x^*\| \leq \frac{3L}{2m}\|x^k - x^*\|^2$, $\forall k = 0, 1, 2, \dots$

Local convergence: interpretation and implications

Corollary

$$\|\mathbf{x}^k - \mathbf{x}^*\| \leq \frac{1}{C} (CX\|\mathbf{x}^0 - \mathbf{x}^*\|)^{2^k}, \text{ where } C = \frac{3L}{2m}$$

- Interpretation of the result:

- Once in a δ -neighborhood of \mathbf{x}^* , the algorithm stays there

- Quadratic convergence within the δ -neighborhood: $\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|^2} \leq \frac{3L}{2m}$

→ Solution within ε of the optimum after at most

$$\left\lceil \frac{\log\left(\frac{\log(C\varepsilon)}{\log(C\|\mathbf{x}^0 - \mathbf{x}^*\|)}\right)}{\log 2} \right\rceil \text{ iterations}$$

- Newton's method is attracted to local minima & maxima: $\nabla f(\mathbf{x}^*) = \mathbf{0}$
- β , m and L are hard to estimate, but not used in the algorithm
- The algorithm and local convergence do not require the convexity of f , only that $H(\mathbf{x}^*)$ is nonsingular and not badly behaved near \mathbf{x}^* .

Proof of the theorem (1/2)

- Notation: $g(\mathbf{x}) = \nabla f(\mathbf{x})$ and $H(\mathbf{x}) = \nabla^2 f(\mathbf{x})$

Lemma

$$g(\mathbf{x}^k) - g(\mathbf{x}^*) = \int_0^1 H(\mathbf{x}^* + t(\mathbf{x}^k - \mathbf{x}^*))(\mathbf{x}^k - \mathbf{x}^*)dt$$

- By definition of \mathbf{x}^{k+1} , we derive:

$$\begin{aligned}\mathbf{x}^{k+1} - \mathbf{x}^* &= \mathbf{x}^k - \mathbf{x}^* - (H(\mathbf{x}^k))^{-1}g(\mathbf{x}^k) \\ &= \mathbf{x}^k - \mathbf{x}^* - (H(\mathbf{x}^k))^{-1}(g(\mathbf{x}^k) - \underbrace{g(\mathbf{x}^*)}_{=0}) \\ &= (\mathbf{x}^k - \mathbf{x}^*) - (H(\mathbf{x}^k))^{-1} \int_0^1 H(\mathbf{x}^* + t(\mathbf{x}^k - \mathbf{x}^*))(\mathbf{x}^k - \mathbf{x}^*)dt \\ &= (H(\mathbf{x}^k))^{-1} \int_0^1 [H(\mathbf{x}^k) - H(\mathbf{x}^* + t(\mathbf{x}^k - \mathbf{x}^*))](\mathbf{x}^k - \mathbf{x}^*)dt\end{aligned}$$

Proof of the theorem (2/2)

Lemma

Under the conditions of the theorem, we have

$$\|H(\mathbf{x})^{-1}\| \leq \frac{1}{m - L\|\mathbf{x} - \mathbf{x}^*\|}, \quad \forall \mathbf{x} \in \mathcal{B}(\mathbf{x}^*, \beta)$$

$$\begin{aligned} \|\mathbf{x}^{k+1} - \mathbf{x}^*\| &\leq \|H(\mathbf{x}^k)^{-1}\| \int_0^1 \|H(\mathbf{x}^k) - H(\mathbf{x}^* + t(\mathbf{x}^k - \mathbf{x}^*))\| \|\mathbf{x}^k - \mathbf{x}^*\| dt \\ &\leq \frac{1}{m - L\|\mathbf{x} - \mathbf{x}^*\|} \cdot \int_0^1 L(1-t)\|\mathbf{x}^k - \mathbf{x}^*\| dt \cdot \|\mathbf{x}^k - \mathbf{x}^*\| \\ &= \frac{L}{2(m - L\|\mathbf{x} - \mathbf{x}^*\|)} \|(\mathbf{x}^k - \mathbf{x}^*)\|^2 \end{aligned}$$

- Since $L\|\mathbf{x}^k - \mathbf{x}^*\| \leq 2m/3$, we obtain:

$$1. \|\mathbf{x}^{k+1} - \mathbf{x}^*\| \leq \frac{2m/3}{2(m - 2m/3)} \|(\mathbf{x}^k - \mathbf{x}^*)\| = \|(\mathbf{x}^k - \mathbf{x}^*)\| \leq \delta$$

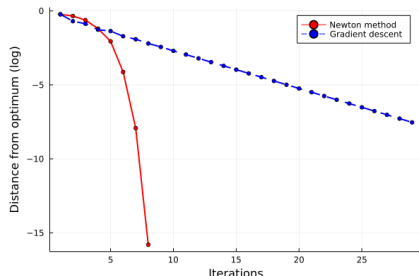
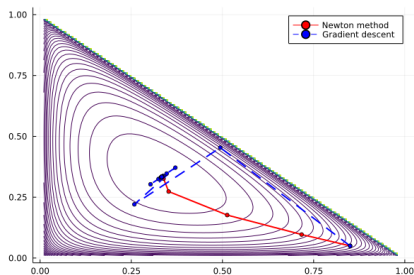
$$2. \|\mathbf{x}^{k+1} - \mathbf{x}^*\| \leq \frac{L}{2(m - 2m/3)} \|(\mathbf{x}^k - \mathbf{x}^*)\|^2 = \frac{3L}{2m} \|(\mathbf{x}^k - \mathbf{x}^*)\|^2$$

Example and illustration

$$f(x_1, x_2) = -\log(1 - x_1 - x_2) - \log x_1 - \log x_2$$

$$\nabla f(x_1, x_2) = \left(\frac{1}{1 - x_1 - x_2} - \frac{1}{x_1}; \frac{1}{1 - x_1 - x_2} - \frac{1}{x_2} \right)^{\top}$$

$$\nabla^2 f(x_1, x_2) = \begin{bmatrix} \left(\frac{1}{1 - x_1 - x_2} \right)^2 + \left(\frac{1}{x_1} \right)^2 & \left(\frac{1}{1 - x_1 - x_2} \right)^2 \\ \left(\frac{1}{1 - x_1 - x_2} \right)^2 & \left(\frac{1}{1 - x_1 - x_2} \right)^2 + \left(\frac{1}{x_2} \right)^2 \end{bmatrix}$$



Global convergence

Global convergence: issues and solutions

- Quadratic convergence in Newton's method is "local"
 - There is no guarantee that $f(\mathbf{x}^{k+1}) \leq f(\mathbf{x}^k)$ at each iteration
 - Newton's method is attracted to local minima & maxima: $\nabla f(\mathbf{x}^*) = \mathbf{0}$
 - What happens if we start "far" away from \mathbf{x}^* ?

→ Augment algorithm with line search: $\alpha^k = \arg \min_{\alpha} f(\mathbf{x}^k + \alpha \mathbf{d}^k)$

Proposition

If $\nabla^2 f(\mathbf{x}) \succeq \mathbf{0}$, then $\mathbf{d} = -\nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x}) \neq \mathbf{0}$ is a descent direction:
 $f(\mathbf{x} + \alpha \mathbf{d}) < f(\mathbf{x})$ for α sufficiently small.

Algorithm (Newton's method with exact line search)

1. Initialization: starting point $\mathbf{x}^0 \in \mathbb{R}^n$, and iteration counter $k = 0$
2. Repeat, until termination criterion is reached
 - 2.1 Update iteration counter: $k \leftarrow k + 1$
 - 2.2 Determine Newton's direction $\mathbf{d}^k = -(\nabla^2 f(\mathbf{x}^k))^{-1} \nabla f(\mathbf{x}^k)$
 - 2.3 Determine α^k via exact line search, and update $\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k + \alpha^k \mathbf{d}^k$

Global convergence: (strong) convexity to the rescue

Theorem

$f(\cdot)$ twice continuously differentiable. Assume that:

- f is M -smooth
- f is m -strongly convex
- $\nabla^2 f(\mathbf{x})$ is L -Lipschitz: $\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \forall \mathbf{x}, \mathbf{y}$

Newton's method with line search satisfies, for $0 < \eta \leq m^2/L$, $\gamma > 0$:

1. If $\|\nabla f(\mathbf{x}^k)\| \geq \eta$: $f(\mathbf{x}^{k+1}) - f(\mathbf{x}^k) \leq -\gamma$
2. If $\|\nabla f(\mathbf{x}^k)\| < \eta$: $\|\nabla f(\mathbf{x}^{k+1})\| < \eta$ & $\|\nabla f(\mathbf{x}^{k+1})\| \leq \frac{L}{2m^2} \|\nabla f(\mathbf{x}^k)\|^2$

→ Two phases in Newton's method:

1. Damped phase: progress of at least η per iteration
2. Local phase: quadratic convergence within a local neighborhood

→ Main objective: getting quickly into a good neighborhood

$$\frac{f(\mathbf{x}^0 - \mathbf{z}^*)}{\gamma} + \log_2 \log_2 \left(\frac{\varepsilon_0}{\varepsilon} \right) \text{ iterations, for some constant } \varepsilon_0$$

Damped phase: proof

- $\nabla f(\mathbf{x}^k)^\top \mathbf{d}^k = -\nabla f(\mathbf{x}^k)^\top (\nabla^2 f(\mathbf{x}^k))^{-1} \nabla f(\mathbf{x}^k) = -\lambda(\mathbf{x}^k)^2$
- Due to the M -smoothness of the function f :

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

$$\lambda(\mathbf{x}^k)^2 = \nabla f(\mathbf{x}^k)^\top (\nabla^2 f(\mathbf{x}^k))^{-1} \nabla f(\mathbf{x}^k) \geq \frac{1}{M} \|\nabla f(\mathbf{x}^k)\|^2$$

- Due to the strong convexity of the function f :

$$\lambda(\mathbf{x}^k)^2 = (\mathbf{d}^k)^\top \nabla^2 f(\mathbf{x}^k) \mathbf{d}^k \geq m \|\mathbf{d}^k\|^2$$

- By combining the two properties, we obtain:

$$f(\mathbf{x}^k + \alpha \mathbf{d}^k) \leq f(\mathbf{x}^k) - \alpha \lambda(\mathbf{x}^k)^2 + \frac{M\alpha^2}{2m} \lambda(\mathbf{x}^k)^2$$

- Exact line search: $\alpha^k \in \arg \min_{\alpha} f(\mathbf{x}^k + \alpha \mathbf{d}^k)$
- By minimizing over α on both sides, we obtain:

$$\begin{aligned} f(\mathbf{x}^k + \alpha^k \mathbf{d}^k) &\leq f(\mathbf{x}^k) - \frac{m}{2M} \lambda(\mathbf{x}^k)^2 = f(\mathbf{x}^k) - \frac{m}{2M^2} \|\nabla f(\mathbf{x}^k)\|^2 \\ \implies f(\mathbf{x}^{k+1}) - f(\mathbf{x}^k) &= f(\mathbf{x}^k + \alpha^k \mathbf{d}^k) - f(\mathbf{x}^k) \leq -\frac{m}{2M^2} \eta^2 := -\gamma \end{aligned}$$

Conclusion

Summary

Takeaway

Newton's method was originally developed to find roots of equations, and then extended to solve optimization problems.

Takeaway

Newton's method augments descent methods by leveraging second-order (Hessian) information.

Takeaway

*Convergence of Newton's method is very fast (quadratic) locally.
⇒ two steps: (i) finding a solution near a stationary point; and (ii) finding the stationary point, exploiting quadratic convergence.*

Takeaway

Numerous enhancements exist to address global convergence issues.