Second-Order Optimization

Optimization Techniques (ENGG*6140)

School of Engineering, University of Guelph, ON, Canada

Course Instructor: Benyamin Ghojogh Winter 2023 Newton's Method

<u>Newton-Raphson</u> root finding method $\frac{t_{a}h(x)}{t_{a}h(x)}$ $x^{2} + 3x = 0$

- We can find the root of a function $f: x \mapsto f(x)$ by solving the equation $f(x) \stackrel{\text{set}}{=} 0$.
- The root of function can be found iteratively where we get closer to the root over iterations.
- One of the iterative root-finding methods is the <u>Newton-Raphson method</u> [1]. In every iteration, it finds the next solution as:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \frac{f(\mathbf{x}^{(k)})}{\nabla f(\mathbf{x}^{(k)})},$$
(1)

where $\nabla f(\mathbf{x}^{(k)})$ is the derivative of function w.r.t. \mathbf{x} .

Univariate Newton's method

• Recall Eq. (1): We saw that Newton-Raphson method solves $f(x) \stackrel{\text{set}}{=} 0$ by:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \frac{f(\mathbf{x}^{(k)})}{\nabla f(\mathbf{x}^{(k)})}.$$

• In unconstrained optimization, we can find the <u>extremum (minimum or maximum)</u> of the function by setting its derivative to zero, i.e. with f(x)• Therefore, the root of $\nabla f(x) \stackrel{\text{set}}{=} 0$ can be found by Newton-Raphson method. We replace f(x) with $\nabla f(x)$ in Eq. (1): $x^{(k+1)} := x^{(k)} - \eta^{(k)} \nabla f(x^{(k)})$ (2)

where $\nabla^2 f(\mathbf{x}^{(k)})$ is the second derivative of function w.r.t. \mathbf{x} and we have included a step size at iteration k denoted by $\eta^{(k)} > 0$. This step size can be either fixed or adaptive.

Multivariate Newton's method

Recall Eq. (2):

• If x is multivariate, i.e. $x \in \mathbb{R}^d$, Eq. (2) is written as:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \eta^{(k)} \left(\nabla^2 f(\mathbf{x}^{(k)}) \right)^{-1} \nabla f(\mathbf{x}^{(k)}),$$
(3)

 $\frac{|k|}{\lambda} - \frac{|k|}{1} \left(\frac{\nabla^2 f(\lambda^{k})}{\nabla^2 f(\lambda^{k})} \right)$

where $\nabla f(\mathbf{x}^{(k)}) \in \mathbb{R}^d$ is the gradient of function w.r.t. \mathbf{x} and $\nabla^2 f(\mathbf{x}^{(k)}) \in \mathbb{R}^{d \times d}$ is the Hessian matrix w.r.t. \mathbf{x} .

 $\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \eta^{(k)} \frac{\nabla f(\mathbf{x}^{(k)})}{\nabla^2 f(\mathbf{x}^{(k)})}$

• Because of the second derivative or the Hessian, this optimization method is a second-order method. The name of this method is the **Newton's method**.

Newton's Method for Unconstrained Optimization

Newton's Method for Unconstrained Optimization

• Consider the following optimization problem:

$$\begin{array}{c} \underset{x}{\text{minimize}} \quad f(x). \end{array}$$
(4)

where f(.) is a convex function.

 Iterative optimization can be <u>first-order</u> or <u>second-order</u>. Iterative optimization updates solution iteratively:

$$\boldsymbol{x}^{(k+1)} := \boldsymbol{x}^{(k)} + \Delta \boldsymbol{x}, \qquad (5)$$

The update continues until Δx becomes very small which is the convergence of optimization.

• Recall that in the first-order optimization, the step of updating is $\Delta x := -\nabla f(x)$.

Newton's Method for Unconstrained Optimization

 Near the optimal point <u>x</u>*, gradient is very small so the second-order Taylor series expansion of function becomes:

$$f(\mathbf{x}) \approx f(\mathbf{x}^*) + \underbrace{\nabla f(\mathbf{x}^*)^{\top}}_{\approx 0} (\mathbf{x} - \mathbf{x}^*) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^{\top} \nabla^2 f(\mathbf{x}^*) (\mathbf{x} - \mathbf{x}^*)$$
$$\approx f(\mathbf{x}^*) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^{\top} \nabla^2 f(\mathbf{x}^*) (\mathbf{x} - \mathbf{x}^*).$$
(6)

This shows that the function is almost quadratic near the optimal point.

• Following this intuition, Newton's method uses Hessian $\nabla^2 f(\mathbf{x})$ in its updating step:

$$\Delta \mathbf{x} := -\nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x}).$$
(7)

• In the literature, this equation is sometimes restated to:

$$\nabla^2 f(\mathbf{x}) \Delta \mathbf{x} := -\nabla f(\mathbf{x}).$$
(8)

Newton's Method for Equality Constrained Optimization

Newton's method for equality constrained optimization

• The optimization problem may have equality constraints:

• After a step of update by
$$p = \Delta x$$
, this optimization becomes:
minimize $f(x)$ $f(x + p)$ $f(x + p)$ $f(x + p) = b$.
• The Lagrangian of this optimization problem is:
 $\mathcal{L} \neq f(x + p) + \frac{\nu^{\top}(A(x + p) - b)}{F(x + p) + \frac{\nu^{\top}(A(x + p) - b)}{F(x + p) + \frac{1}{2}p^{\top}\nabla^{2}f(x)p}$. (10)

• Substituting this into the Lagrangian gives:

$$\bigstar \quad \mathcal{L} = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top \mathbf{p} + \frac{1}{2} \mathbf{p}^\top \nabla^2 f(\mathbf{x}) \mathbf{p} + \mathbf{\nu}^\top (\mathbf{A}(\mathbf{x} + \mathbf{p}) - \mathbf{b}).$$

- 2

Newton's method for equality constrained optimization

We found:

$$\mathbf{\mathbf{\mathcal{L}}} = \mathbf{\mathbf{\mathcal{L}}} = \mathbf{\mathbf{\mathcal{L}}} \mathbf{\mathbf{\mathcal{L}}} + \mathbf{\mathbf{\nabla}} \mathbf{\mathbf{\mathcal{L}}} \mathbf{\mathbf{\mathcal{L}}}^{\top} \mathbf{\mathbf{p}} + \frac{1}{2} \mathbf{\mathbf{p}}^{\top} \mathbf{\mathbf{\nabla}}^{2} \mathbf{\mathbf{\mathcal{L}}} \mathbf{\mathbf{p}} + \mathbf{\mathbf{\nu}}^{\top} \mathbf{\mathbf{\mathcal{L}}} \mathbf{\mathbf{\mathcal{L}}} \mathbf{\mathbf{\mathcal{L}}} \mathbf{\mathbf{\mathbf{p}}} \mathbf{\mathbf{\mathcal{L}}} \mathbf{\mathbf{\mathcal{L}}} \mathbf{\mathbf{\mathcal{L}}} \mathbf{\mathbf{\mathbf{p}}} \mathbf{\mathbf{\mathcal{L}}} \mathbf{\mathcal{L}} \mathbf{\mathbf{\mathcal{L}}} \mathbf{\mathcal{L}}$$

• According to KKT conditions, the primal and dual residuals must be zero:

where we have $\nabla^3 f(\mathbf{x}) \approx 0$ because the third-order gradient is usually very small compared to the first and second gradients and (a) is because of the constraint $A\mathbf{x} - \mathbf{b} = \mathbf{0}$ in problem (9).

Eqs. (12) and (13) can be written as a system of equations:

$$+ \left(\begin{bmatrix} \nabla^2 f(\mathbf{x})^\top & \mathbf{A}^\top \\ \mathbf{A} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{p} \\ \nu \end{bmatrix} = \begin{bmatrix} -\nabla f(\mathbf{x}) \\ \mathbf{0} \end{bmatrix} \cdot \mathbf{\gamma} \mathbf{\gamma} \right)$$
(14)

Solving this system of equations gives the desired step p (i.e., Δx) for updating the solution at the iteration.

Starting with non-feasible initial point

- Newton's method can even start with a non-feasible point which does not satisfy all the constraints.
- If the initial point for optimization is not a feasible point, i.e., $(Ax b \neq 0)$ Eq. (13) becomes:

$$\nabla_{\nu} \mathcal{L} = \mathbf{A}(\mathbf{x} + \mathbf{p}) - \mathbf{b} \stackrel{\text{set}}{=} \mathbf{0} \implies \mathbf{A}\mathbf{p} = -(\mathbf{A}\mathbf{x} - \mathbf{b}). \qquad (15)$$

• Therefore, for the first iteration, we solve the following system rather than Eq. (14):

$$\underbrace{\left[\begin{array}{ccc} \nabla^{2} f(\mathbf{x})^{\top} & \mathbf{A}^{\top} \\ \mathbf{A} & \mathbf{0} \end{array}\right]}_{\mathbf{A}} \underbrace{\left[\begin{array}{ccc} \nabla f(\mathbf{x}) \\ \mathbf{A} & \mathbf{0} \end{array}\right]}_{\mathbf{A}} \underbrace{\left[\begin{array}[ccc} \nabla f(\mathbf{x}) \\ \mathbf{A} & \mathbf{0} \end{array}\right]}_{\mathbf{A}} \underbrace{\left[\begin{array}[ccc} \nabla f(\mathbf{x}) \\ \mathbf{A} & \mathbf{0} \end{array}\right]}_{\mathbf{A}} \underbrace{\left[\begin{array}[ccc} \nabla f(\mathbf{x}) \\ \mathbf{A} & \mathbf{A} \end{array}\right]}_{\mathbf{A}} \underbrace{\left[\begin{array}[ccc} \nabla f(\mathbf{x}) \\ \mathbf{A} & \mathbf{A} \end{array}\right]}_{\mathbf{A}} \underbrace{\left$$

and we use Eq. (16) for the rest of iterations because the next points will be in the feasibility set (because we force the solutions to satisfy Ax = b).

Interior-Point and Barrier Methods: Newton's Method for Inequality Constrained Optimization

The optimization problem may have inequality constraints:

- We can solve constrained optimization problems using **Barrier methods**, also known as interior-point methods [2, 3, 4, 5].
- Interior-point methods were first proposed in 1967 [6].
- The interior-point method is also referred to as the <u>Unconstrained Minimization</u> <u>Technique (UMT)</u> or <u>Sequential UMT (SUMT)</u> [7] because it converts the problem to an <u>unconstrained</u> problem and solves it iteratively.

• The barrier methods or the interior-point methods, convert inequality constrained problems to equality constrained or unconstrained problems. Ideally, we can do this conversion using the indicator function I(.) which is zero if its input condition is satisfied and is infinity otherwise:

$$\mathbb{I}(\boldsymbol{x}\in\mathcal{S}) = \begin{cases} 0 & \text{if } \boldsymbol{x}\in\mathcal{S} \\ \infty & \text{if } \boldsymbol{x}\notin\mathcal{S}. \end{cases}$$
(18)



• The indicator function is not differentiable because it is not smooth:

$$\mathbb{I}(y_i(\boldsymbol{x}) \leq 0) := \begin{cases} 0 & \text{if } y_i(\boldsymbol{x}) \leq 0 \\ \infty & \text{if } y_i(\boldsymbol{x}) > 0. \end{cases}$$



• One of the barrier functions is logarithm, named the **logarithmic barrier** or **log barrier** in short. It approximates the indicator function by:

$$\mathbb{I}(y_i(\mathbf{x}) \leq 0) \approx -\frac{1}{t} \log(-y_i(\mathbf{x})),$$
(21)

where t > 0 (usually a large number such as $t = 10^6$) and the approximation becomes more accurate by $t \to \infty$.



Y.W=1

• The problem had become:

This optimization problem is an equality constrained optimization problem which we already explained how to solve.

• Note that there exist many approximations for the barrier. One of mostly used methods is the logarithmic barrier.

• If the problem is convex, the iterative solutions of the interior-point method satisfy:

$$\{ \mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots \} \to \mathbf{x}^{*}, \\ \{ \boldsymbol{\nu}^{(0)}, \boldsymbol{\nu}^{(1)}, \boldsymbol{\nu}^{(2)}, \dots \} - \mathbf{\nu}^{*}, \\ \underline{f(\mathbf{x}^{(0)}) \ge f(\mathbf{x}^{(1)}) \ge f(\mathbf{x}^{(2)}) \ge \dots \ge f(\mathbf{x}^{*}), \\ \underline{g(\boldsymbol{\nu}^{(0)}) \le g(\boldsymbol{\nu}^{(1)}) \le \dots \le g(\boldsymbol{\nu}^{*}). \end{cases}$$
(23)

- If the optimization problem is a convex problem, the solution of interior-point method is the global solution; otherwise, the solution is local.
- The interior-point and barrier methods are used in many optimization toolboxes such as CVX [9].

Accuracy of the log barrier method

Theorem (On the sub-optimality of log-barrier method)

Let the optimum of problems (17) and (22) be denoted by f^* and f^*_r , respectively. We have:

$$f^* - \frac{m}{t} \leq f_r^* \leq f^*,$$
(24)

meaning that the optimum of problem (22) is no more than m_1/t from the optimum of problem (17).

Proof.

See our tutorial [10] for proof. Also explained in the next slide.

- The above theorem indicates that by $t \to \infty$, the log-barrier method is more accurate; i.e., the solution of problem (22) is more accurately close to the solution of problem (17).
- This is expected since the approximation in Eq. (21) gets more accurate by increasing t.
- Note that by increasing *t*, optimization gets more accurate but harder to solve and slower to converge.

Acknowledgement

- Some slides of this slide deck are inspired by the lectures of Prof. Stephen Boyd at the Stanford University.
- Our tutorial also has the materials of this slide deck: [10]

References

- J. Stoer and R. Bulirsch, *Introduction to numerical analysis*, vol. 12. Springer Science & Business Media, 2013.
- Y. Nesterov and A. Nemirovskii, Interior-point polynomial algorithms in convex programming.
 SIAM, 1994.
- [3] F. A. Potra and S. J. Wright, "Interior-point methods," *Journal of computational and applied mathematics*, vol. 124, no. 1-2, pp. 281–302, 2000.
- [4] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [5] M. Wright, "The interior-point revolution in optimization: history, recent developments, and lasting consequences," *Bulletin of the American mathematical society*, vol. 42, no. 1, pp. 39–56, 2005.
- [6] I. Dikin, "Iterative solution of problems of linear and quadratic programming," in *Doklady Akademii Nauk*, vol. 174, pp. 747–748, Russian Academy of Sciences, 1967.
- [7] A. V. Fiacco and G. P. McCormick, "The sequential unconstrained minimization technique (SUMT) without parameters," *Operations Research*, vol. 15, no. 5, pp. 820–827, 1967.
- [8] Y. Nesterov, Lectures on convex optimization, vol. 137. Springer, 2018.

References (cont.)

- [9] M. Grant, S. Boyd, and Y. Ye, "CVX: Matlab software for disciplined convex programming," 2009.
- [10] B. Ghojogh, A. Ghodsi, F. Karray, and M. Crowley, "KKT conditions, first-order and second-order optimization, and distributed optimization: Tutorial and survey," arXiv preprint arXiv:2110.01858, 2021.