

# Non-convex Optimization

Optimization Techniques (ENGG\*6140)

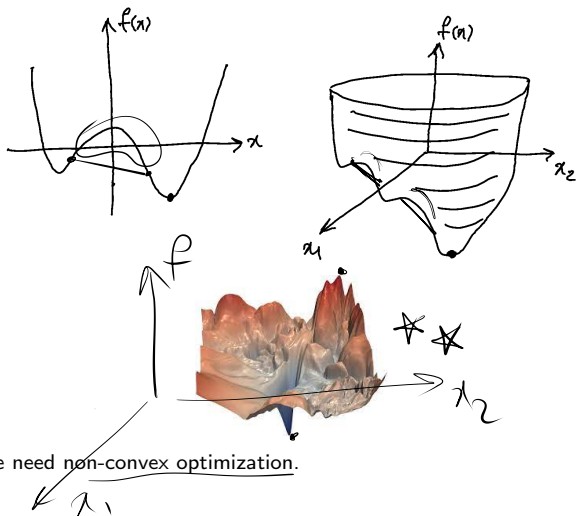
School of Engineering,  
University of Guelph, ON, Canada

Course Instructor: Benyamin Ghojogh  
Winter 2023

## **Non-convex Function and Optimization**

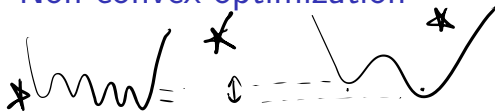
# Non-convex function

- We might have non-convex cost functions (recall the preliminaries) (credit of second image: [1]):

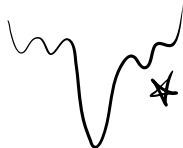


- In these cases, we need non-convex optimization.

# Non-convex optimization



- Consider the following optimization problem:



$$\begin{aligned} & \underset{x}{\text{minimize}} && f(x) \\ & \text{subject to} && y_i(x) \leq 0, \quad i \in \{1, \dots, m_1\}, \\ & && \underline{h_i(x) = 0}, \quad i \in \{1, \dots, m_2\}, \end{aligned} \quad (1)$$

where the functions  $f(\cdot)$ ,  $y_i(\cdot)$ , and  $h_i(\cdot)$  are not necessarily convex.

- The introduced optimization methods can also work for non-convex problems but they do not guarantee to find the global optimum.
- They can find local minimizers which depend on the random initial solution.
- For example, the optimization landscape of neural network is highly nonlinear and non-convex but backpropagation works very well for it.

layer 1      layer 2       $\star \quad 2 \rightarrow 3 \Rightarrow \text{kernel}$   
 $\star \quad 3 \rightarrow 2 \Rightarrow \text{dim reduction}$

# Non-convex optimization in neural networks

- The optimization landscape of neural network is highly nonlinear and non-convex but backpropagation works very well for it.

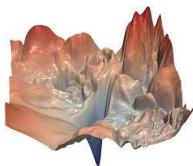
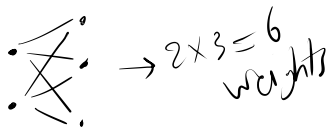
► **Question:** Assume the loss function of neural network is the mean squared error:

$V \rightarrow \text{w.r.t. } y \text{ or } l_i$

$C = \frac{1}{b} \sum_{i=1}^b \|y_i - l_i\|_2^2 \rightarrow \text{convex function}$

where  $b$  is the mini-batch size,  $y_i$  is the output of network, and  $l_i$  is the one-hot encoded label for the  $i$ -th data in the mini-batch.

- This loss is clearly quadratic and convex. Then, why is the neural network highly non-convex? (credit of image: [1])



$\min C$   
 $\star \{w_{ij}\}_{i=1}^6$

- $\star$  • The reason for this is explained in this way: every layer of neural network pulls data to the feature space such as in kernels [2]. In the high-dimensional feature space, all local minimizers are almost global minimizers because the local minimum values are almost equal in that space [3]. Also see [4, 5, 6] to understand why backpropagation optimization works well even in highly non-convex optimization.

# Non-convex optimization

- As was explained, the already introduced first-order and second-order optimization methods can work fairly well for non-convex problems by finding local minimizers depending on the initial solution.
- However, there exist some specific methods for non-convex programming, divided into two categories:
  - ▶ The local optimization methods are faster but do not guarantee to find the global minimizer.
  - ▶ The global optimization methods find the global minimizer but are usually slow to find the answer [7].
- Example for local optimization methods:
  - ▶ Sequential Convex Programming (SCP) [8] is an example for local optimization methods.
  - ▶ It is based on a sequence of convex approximations of the non-convex problem.
  - ▶ It is closely related to Sequential Quadratic Programming (SQP) [9] which is used for constrained nonlinear optimization.
- Example for global optimization methods:
  - ▶ Branch and bound method, first proposed in 1960 [10], is an example for the global optimization methods.
  - ▶ It divides the optimization landscape, i.e. the feasible set, into local parts by a binary tree and solves optimization in every part.
  - ▶ It checks whether the solution of a part is the global solution or not.
- ~~In this slide deck, we explain SCP which is a faster but local method.~~
- Note that another approach for highly non-convex optimization is metaheuristic (search-based) optimization which will be briefly introduced later.

## Convex Approximation

# Convex Approximation

- Recall Eq. (1):

$$\left[ \begin{array}{ll} \underset{x}{\text{minimize}} & f(x) \\ \text{subject to} & y_i(x) \leq 0, \quad i \in \{1, \dots, m_1\}, \\ & h_i(x) = 0, \quad i \in \{1, \dots, m_2\}. \end{array} \right. \quad \leftarrow \quad \star$$

- SCP iteratively solves a convex problem where, at every iteration, it approximates the non-convex problem (1) with a convex problem, based on the current solution, and restricts the variable to be in a so-called trust region [11].
- The trust region makes sure that the variable stays in a locally convex region of the optimization problem.
- At the iteration  $k$  of SCP, we solve the following convex problem:

$$\left[ \begin{array}{ll} \underset{x}{\text{minimize}} & \hat{f}(x) \\ \text{subject to} & \hat{y}_i(x) \leq 0, \quad i \in \{1, \dots, m_1\}, \\ & \hat{h}_i(x) = 0, \quad i \in \{1, \dots, m_2\}, \\ & \star \quad x \in \mathcal{T}^{(k)}, \quad \leftarrow \end{array} \right. \quad \star \quad (2)$$

where  $\hat{f}(\cdot)$ ,  $\hat{y}_i(\cdot)$ , and  $\hat{h}_i(\cdot)$ , are convex approximations of functions  $f(\cdot)$ ,  $y_i(\cdot)$ , and  $h_i(\cdot)$ , and  $\mathcal{T}^{(k)}$  is the trust region at iteration  $k$ .

- This approximated convex problem is also solved iteratively itself using one of the previously introduced methods such as the interior-point method.
- There exist several approaches for convex approximation of the functions. In the following, we introduce some of these approaches.





# Convex Approximation by Taylor Series Expansion

- The non-convex functions  $f(\cdot)$ ,  $y_i(\cdot)$ , and  $h_i(\cdot)$  can be approximated by affine functions (i.e., first-order Taylor series expansion) to become convex. For example, the function  $f(\cdot)$  is approximated as:

$$\hat{f}(\mathbf{x}) = f(\mathbf{x}^{(k)}) + \nabla f(\mathbf{x}^{(k)})^\top (\mathbf{x} - \mathbf{x}^{(k)}). \quad (3)$$

- The functions can also be approximated by quadratic functions (i.e., second-order Taylor series expansion) to become convex. For example, the function  $f(\cdot)$  is approximated as:


$$\hat{f}(\mathbf{x}) = \underbrace{f(\mathbf{x}^{(k)})}_{\text{constant}} + \underbrace{\nabla f(\mathbf{x}^{(k)})^\top (\mathbf{x} - \mathbf{x}^{(k)})}_{\text{linear term}} + \frac{1}{2}(\mathbf{x} - \mathbf{x}^{(k)})^\top \mathbf{P}(\mathbf{x} - \mathbf{x}^{(k)}), \quad (4)$$

where  $\mathbf{P} = \Pi_{\mathbb{S}_+^d}(\nabla^2 f(\mathbf{x}^{(k)}))$  is projection of Hessian onto the symmetric positive semi-definite cone. This projection is performed by setting the negative eigenvalues of Hessian to zero.


- The same approaches can be used for approximation of functions  $y_i(\cdot)$  and  $h_i(\cdot)$  using first- or second-order Taylor expansion.

# Convex Approximation by Particle Method

- We can approximate the functions  $f(\cdot)$ ,  $y_i(\cdot)$ , and  $h_i(\cdot)$  in the domain of trust region using regression. This approach is named the particle method [7].
- Let  $\{\mathbf{x}_i \in \mathcal{T}^{(k)}\}_{i=1}^m$  be  $m$  points which lie in the trust region. We can use least-squares quadratic regression to make the functions convex in the trust region:

★ 
$$\begin{aligned} & \underset{a \in \mathbb{R}, \mathbf{b} \in \mathbb{R}^d, \mathbf{P} \in \mathbb{S}_{++}^d}{\text{minimize}} \sum_{i=1}^m \left( \frac{1}{2} (\mathbf{x}_i - \mathbf{x}^{(k)})^\top \mathbf{P} (\mathbf{x}_i - \mathbf{x}^{(k)}) + \mathbf{b}^\top (\mathbf{x}_i - \mathbf{x}^{(k)}) + a - f(\mathbf{x}_i) \right)^2 \\ & \text{subject to } \mathbf{P} \succeq \mathbf{0}. \end{aligned} \quad (5)$$

★



- Then, the function  $f(\cdot)$  is replaced by its convex approximation:

$$\hat{f}(\mathbf{x}) = (1/2)(\mathbf{x}_i - \mathbf{x}^{(k)})^\top \mathbf{P} (\mathbf{x}_i - \mathbf{x}^{(k)}) + \mathbf{b}^\top (\mathbf{x}_i - \mathbf{x}^{(k)}) + a.$$

- The same approach can be used for approximation of functions  $y_i(\cdot)$  and  $h_i(\cdot)$ .



# Convex Approximation by Quasi-linearization



- Another approach for convex approximation of functions  $f(\cdot)$ ,  $y_i(\cdot)$ , and  $h_i(\cdot)$  is quasi-linearization.

- We should state the function  $f(\cdot)$  in the form  $f(\mathbf{x}) = \mathbf{A}(\mathbf{x})\mathbf{x} + c(\mathbf{x})$ .

- For example, we can use the second-order Taylor series expansion to do this:

$$f(\mathbf{x}) \approx \underbrace{\frac{1}{2}\mathbf{x}^\top \mathbf{P}\mathbf{x}}_{\text{quadratic term}} + \underbrace{\mathbf{b}^\top \mathbf{x} + a}_{\text{linear and constant terms}} = \left(\frac{1}{2}\mathbf{P}\mathbf{x} + \mathbf{b}\right)^\top \mathbf{x} + a$$

Handwritten notes:  $\frac{1}{2}\mathbf{x}^\top \mathbf{P}\mathbf{x}$  with an arrow pointing to  $\mathbf{P}$ , and  $\mathbf{b}^\top \mathbf{x}$  with an arrow pointing to  $\mathbf{b}$ .

so we use  $\mathbf{A}(\mathbf{x}) := \left(\frac{1}{2}\mathbf{P}\mathbf{x} + \mathbf{b}\right)^\top$  and  $c(\mathbf{x}) := a$  which depend on the Taylor expansion of  $f(\mathbf{x})$ .

- Hence, the convex approximation of function  $f(\cdot)$  can be:

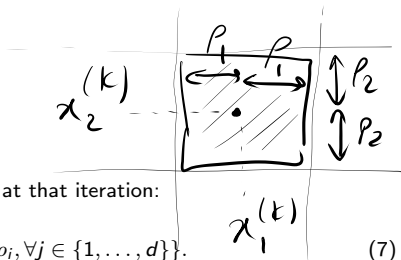
$a\mathbf{x} + \mathbf{b}$

$$\hat{f}(\mathbf{x}) = \mathbf{A}(\mathbf{x}^{(k)})\mathbf{x} + c(\mathbf{x}^{(k)}) = \left(\frac{1}{2}\mathbf{P}\mathbf{x}^{(k)} + \mathbf{b}\right)^\top \mathbf{x} + a. \quad (6)$$

- The same approach can be used for approximation of functions  $y_i(\cdot)$  and  $h_i(\cdot)$ .

**Trust Region**

# Formulation of Trust Region

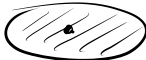


- The trust region can be a box around the point at that iteration:

$$\star \quad \mathcal{T}^{(k)} := \{ \mathbf{x} \mid |x_j - x_j^{(k)}| \leq \rho_i, \forall j \in \{1, \dots, d\} \}. \quad (7)$$

where  $x_j$  and  $x_j^{(k)}$  are the  $j$ -th element of  $\mathbf{x}$  and  $\mathbf{x}^{(k)}$ , respectively, and  $\rho_i$  is the bound of box for the  $j$ -th dimension.

- Another option for trust region is an ellipse around the point to have a quadratic trust region:

$$\star \quad \mathcal{T}^{(k)} := \{ \mathbf{x} \mid (\mathbf{x} - \mathbf{x}^{(k)})^\top \mathbf{P}^{-1} (\mathbf{x} - \mathbf{x}^{(k)}) \leq \rho \}, \quad (8)$$


where  $\mathbf{P} \in \mathbb{S}_{++}^d$  (is symmetric positive definite) and  $\rho > 0$  is the radius of ellipse.

# Updating Trust Region

- The trust region gets updated in every iteration of SCP. In the following, we explain how the trust region can be updated.

- Recall Eq. (1):

$$\begin{array}{l} \text{minimize}_x \quad f(x) \\ \text{subject to} \quad y_i(x) \leq 0, \quad i \in \{1, \dots, m_1\}, \\ \quad \quad \quad h_i(x) = 0, \quad i \in \{1, \dots, m_2\}. \end{array}$$

Handwritten notes:  $y_i'(\lambda) \equiv \max(y_i(\lambda), 0)$  with a star below it. An arrow points from the constraint  $y_i(x) \leq 0$  to the handwritten definition.

- First, we embed the constraints in the objective function of problem (1):

$$\text{minimize}_x \quad \phi(x) := f(x) + \lambda \left( \sum_{i=1}^{m_1} (\max(y_i(x), 0))^2 + \sum_{i=1}^{m_2} |h_i(x)|^2 \right), \quad (9)$$

Handwritten notes: A star is placed above the summation term. Arrows point from the handwritten definition of  $y_i'(\lambda)$  to the  $\max$  term in the equation.

where  $\lambda > 0$  is the regularization parameter. This is called the exact penalty method (1994) [12] because it penalizes violation from the constraints.

- For large enough regularization parameter (which gives importance to violation of constraints), the solution of problem (9) is exactly equal to the solution of problem (1). That is the reason for the term “exact” in the name “exact penalty method”.

# Updating Trust Region

- Recall Eq. (2):

$$\begin{aligned} & \text{minimize}_{\mathbf{x}} \quad \hat{f}(\mathbf{x}) \\ & \text{subject to} \quad \hat{y}_i(\mathbf{x}) \leq 0, \quad i \in \{1, \dots, m_1\}, \\ & \quad \quad \quad \hat{h}_i(\mathbf{x}) = 0, \quad i \in \{1, \dots, m_2\}, \\ & \quad \quad \quad \mathbf{x} \in \mathcal{T}^{(k)}, \end{aligned}$$

- We found Eq. (9):

$$\text{minimize}_{\mathbf{x}} \quad \phi(\mathbf{x}) := f(\mathbf{x}) + \lambda \left( \sum_{i=1}^{m_1} (\max(y_i(\mathbf{x}), 0))^2 + \sum_{i=1}^{m_2} |h_i(\mathbf{x})|^2 \right).$$

- Similar to Eq. (9), we define:

$$\hat{\phi}(\mathbf{x}) := \hat{f}(\mathbf{x}) + \lambda \left( \sum_{i=1}^{m_1} (\max(\hat{y}_i(\mathbf{x}), 0))^2 + \sum_{i=1}^{m_2} |\hat{h}_i(\mathbf{x})|^2 \right), \quad (10)$$

for the problem (2). At the iteration  $k$  of SCP, let  $\hat{\mathbf{x}}^{(k)}$  be the solution of the convex approximated problem (2) using any method such as the interior-point method.

- We calculate the predicted and exact decreases which are  $\hat{\delta} := \phi(\mathbf{x}^{(k)}) - \hat{\phi}(\hat{\mathbf{x}})$  and  $\delta := \phi(\mathbf{x}^{(k)}) - \phi(\hat{\mathbf{x}})$ , respectively.

# Updating Trust Region

★ We calculate the predicted and exact decreases which are  $\hat{\delta} := \phi(\mathbf{x}^{(k)}) - \hat{\phi}(\hat{\mathbf{x}})$  and  $\delta := \phi(\mathbf{x}^{(k)}) - \phi(\hat{\mathbf{x}})$ , respectively.

★ Two cases may happen:

- ▶ We have progress in optimization if  $\alpha\hat{\delta} \leq \delta$  where  $0 < \alpha < 1$  (e.g.,  $\alpha = 0.1$ ). In this case, we accept the approximate solution:

$$\star \quad \boxed{\mathbf{x}^{(k+1)} := \hat{\mathbf{x}}}$$

and we increase the size of trust region, for the next iteration of SCP, by:

$$\star \quad \boxed{\rho^{(k+1)} := \beta \rho^{(k)},}$$

where  $\beta \geq 1$  (e.g.,  $\beta = 1.1$ ).

- ▶ We do not have progress in optimization if  $\alpha\hat{\delta} > \delta$ . In this case, we reject the approximate solution:

$$\star \quad \boxed{\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)},}$$

and we decrease the size of trust region, for the next iteration of SCP, by:

$$\star \quad \boxed{\rho^{(k+1)} := \gamma \rho^{(k)},}$$

where  $0 < \gamma < 1$  (e.g.,  $\gamma = 0.5$ ).

- ★ In summary, the trust region is expanded if we find a good solution; otherwise, it is made smaller.



## **Branch and Bound Method**

# Branch and Bound Method

min:  $f_2^* \leq f_1^*$

max:  $f_2^* \geq f_1^*$

$P1$   
nonconvex  
 $f_1^*$

relax  $\rightarrow$   $P2$   
 $f_2^*$  convex relaxation

- Branch and bound method, proposed in 1960 [10], is a global optimization method, which is slow but find the global optimum.

- The idea of this method:

- ▶ divide the feasibility set region into two convex sets (e.g., rectangular (hyper-cubic) regions).
- ▶ in each region, find an upper bound and a lower bound.
- ▶ choose the better region among these two.
- ▶ do it iteratively until convergence.

- ★ • Approximation of lower bound of optimal value  $f^*$  in a region:

- ▶ it can be the solution of the convex relaxation of the optimization problem.

- ★ • Approximation of upper bound of optimal value  $f^*$  in a region:

- ▶ it can be the function value at any point in the region.

$l \rightarrow u$   
 $f^* \leq f(x) \forall x$

- ★ • The more we progress in the algorithm, the closer (tighter) the lower bound and the upper bound get to each other.

$u$   
 $l$   
 $l \leq f^* \leq u$   
 $u \downarrow$   
 $l \uparrow$   
 $u = l \rightarrow u = l = f^*$

# Branch and Bound Method

The algorithm of branch and bound method:

- 1 start with the feasibility set as the region.
- 2 calculate the lower bound  $l$  and upper bound  $u$  of the feasibility set. Set the overall lower bound and overall upper bound to the lower bound and upper bound, respectively:

$$l_{\text{overall}} = l, \quad u_{\text{overall}} = u$$



- 3 loop over iterations:
  - 1 split the region into left and right regions.
  - 2 calculate the lower bound and upper bound of the left and right regions,  $l_{\text{left}}, u_{\text{left}}, l_{\text{right}}, u_{\text{right}}$ .
  - 3 Update the overall lower bound if the largest lower bound among left and right regions is larger than the overall lower bound.

$$l_{\text{overall}} = \max(l_{\text{left}}, l_{\text{right}}, l_{\text{overall}}).$$

- 4 Update the overall upper bound if the smallest upper bound among left and right regions is smaller than the overall upper bound.

$$u_{\text{overall}} = \min(u_{\text{left}}, u_{\text{right}}, u_{\text{overall}}).$$

- 5 if  $u_{\text{overall}} - l_{\text{overall}} < \epsilon$ : terminate and, for minimization, return  $l_{\text{left}}$  and its corresponding  $x^*$  if  $l_{\text{left}} < l_{\text{right}}$ . But return  $l_{\text{right}}$  and its corresponding  $x^*$  if  $l_{\text{right}} < l_{\text{left}}$ .

# 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: [13]

# References

- [1] H. Li, Z. Xu, G. Taylor, C. Studer, and T. Goldstein, "Visualizing the loss landscape of neural nets," Advances in neural information processing systems, vol. 31, 2018.
- [2] B. Ghogh, A. Ghodsi, F. Karray, and M. Crowley, "Reproducing kernel Hilbert space, Mercer's theorem, eigenfunctions, Nyström method, and use of kernels in machine learning: Tutorial and survey," arXiv preprint arXiv:2106.08443, 2021.
- [3] S. Feizi, H. Javadi, J. Zhang, and D. Tse, "Porcupine neural networks:(almost) all local optima are global," arXiv preprint arXiv:1710.02196, 2017.
- 4 [4] M. Soltanolkotabi, A. Javanmard, and J. D. Lee, "Theoretical insights into the optimization landscape of over-parameterized shallow neural networks," IEEE Transactions on Information Theory, vol. 65, no. 2, pp. 742–769, 2018.
- 4 [5] Z. Allen-Zhu, Y. Li, and Y. Liang, "Learning and generalization in overparameterized neural networks, going beyond two layers," Advances in neural information processing systems, 2019.
- 4 [6] Z. Allen-Zhu, Y. Li, and Z. Song, "A convergence theory for deep learning via over-parameterization," in International Conference on Machine Learning, pp. 242–252, 2019.
- [7] J. Duchi, S. Boyd, and J. Mattingley, "Sequential convex programming," tech. rep., Notes for EE364b, Stanford University, 2018.

## References (cont.)

- [8] Q. T. Dinh and M. Diehl, “Local convergence of sequential convex programming for nonconvex optimization,” in *Recent Advances in Optimization and its Applications in Engineering*, pp. 93–102, Springer, 2010.
- [9] P. T. Boggs and J. W. Tolle, “Sequential quadratic programming,” *Acta numerica*, vol. 4, pp. 1–51, 1995.
- [10] A. H. Land and A. G. Doig, “An automatic method for solving discrete programming problems,” *Econometrica*, vol. 28, no. 3, pp. 497–520, 1960.
- [11] A. R. Conn, N. I. Gould, and P. L. Toint, *Trust region methods*. SIAM, 2000.
- [12] G. Di Pillo, “Exact penalty methods,” in *Algorithms for Continuous Optimization*, pp. 209–253, Springer, 1994.
- [13] 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.