

Lecture 07: Optimization

Introduction to Machine Learning [25737]

Sajjad Amini

Sharif University of Technology

Contents

- 1 Basic Definitions
- 2 First Order Methods
- 3 Second Order Methods
- 4 Stochastic Gradient Descent
- 5 Constrained Optimization

The material in the slides except cited are inspired from the following reference:

- Murphy, K. P. (2022). *Probabilistic machine learning: an introduction*. MIT press.

Section 1

Basic Definitions

Optimization Problem

Training as Optimization Problem

Assume function $\mathcal{L} : \Theta \rightarrow \mathbb{R}$. An optimization problem is the process of finding the value for vector $\theta \in \Theta$, denoted θ^* that minimizes L . We write this process as:

$$\theta^* \in \overbrace{\operatorname{argmin}_{\theta \in \Theta} \mathcal{L}(\theta)}^{\text{Set}}$$

where:

$\mathcal{L}(\theta)$	Loss function or cost function
$R(\theta) = -\mathcal{L}(\theta)$	Score function or reward function
$R(\theta), \mathcal{L}(\theta)$	Objective function
$\Theta \subseteq \mathbb{R}^D$	Parameter space
D	Number of Variables

Global vs Local Optimization

Global Minimum

The set for global minimum is:

$$\{\boldsymbol{\theta}^* : \forall \boldsymbol{\theta} \in \Theta, \mathcal{L}(\boldsymbol{\theta}^*) \leq \mathcal{L}(\boldsymbol{\theta})\}$$

The set for strict global minimum is:

$$\{\boldsymbol{\theta}^* : \forall \boldsymbol{\theta} \in \Theta, \mathcal{L}(\boldsymbol{\theta}^*) < \mathcal{L}(\boldsymbol{\theta})\}$$

Local Minimum

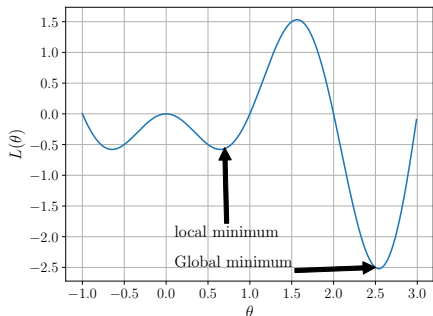
The set for local minimum is:

$$\{\boldsymbol{\theta}^* : \exists \delta > 0, \forall \boldsymbol{\theta} \in \Theta, \boldsymbol{\theta} \neq \boldsymbol{\theta}^*, \text{ if } \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| < \delta \text{ then } \mathcal{L}(\boldsymbol{\theta}^*) \leq \mathcal{L}(\boldsymbol{\theta})\}$$

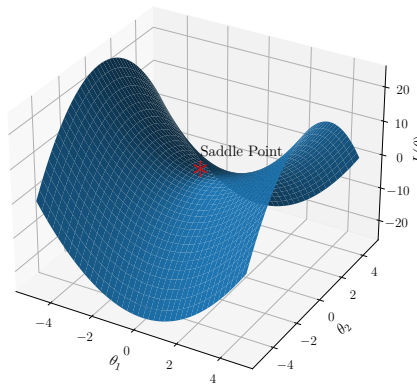
The set for strict local minimum is:

$$\{\boldsymbol{\theta}^* : \exists \delta > 0, \forall \boldsymbol{\theta} \in \Theta, \boldsymbol{\theta} \neq \boldsymbol{\theta}^*, \text{ if } \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| < \delta \text{ then } \mathcal{L}(\boldsymbol{\theta}^*) < \mathcal{L}(\boldsymbol{\theta})\}$$

Illustration of Local and Global Minimum



(a) Local minimum vs global minimum



(b) Saddle point

Local Minimum

Assume \mathcal{L} to be twice differentiable and $\Theta = \mathbb{R}^D$. Consider a point $\boldsymbol{\theta}^* \in \mathbb{R}^D$ and let $\mathbf{g}^* = \mathbf{g}(\boldsymbol{\theta})\big|_{\boldsymbol{\theta}^*}$ and $\mathbf{H}^* = \mathbf{H}(\boldsymbol{\theta})\big|_{\boldsymbol{\theta}^*}$ to be gradient vector and Hessian matrix at $\boldsymbol{\theta}^*$. Then:

- *Necessary condition:* If $\boldsymbol{\theta}^*$ is a local minimum, then we must have $\mathbf{g}^* = \mathbf{0}$ and $\mathbf{H}^* \succeq 0$.
- *Sufficient condition:* If $\mathbf{g}^* = \mathbf{0}$ and $\mathbf{H}^* \succ 0$, then $\boldsymbol{\theta}^*$ is a local minimum.

Constrained vs Unconstrained Optimization

Feasible Set

Feasible set is the subset of the parameter space that satisfies the constraints over the parameter vector as:

$$\mathcal{C} = \{\boldsymbol{\theta} : g_j(\boldsymbol{\theta}) \leq 0, j \in \mathcal{I} \text{ and } h_k(\boldsymbol{\theta}) = 0, k \in \epsilon\} \subseteq \mathbb{R}^D$$

where:

$g_j(\boldsymbol{\theta})$	Inequality constraints
$h_k(\boldsymbol{\theta}) = 0$	Equality constraints
\mathcal{I}	Index set for Inequality constraints
ϵ	Index set for Equality constraints

Constrained vs Unconstrained Optimization

$$\boldsymbol{\theta}^* \in \underset{\boldsymbol{\theta} \in \mathcal{C}}{\operatorname{argmin}} \mathcal{L}(\boldsymbol{\theta})$$

The above optimization problem is unconstrained if $\mathcal{C} \in \mathbb{R}^D$, otherwise it is constrained.

Section 2

First Order Methods

General Properties

First order methods are methods that:

- Leverage first order derivatives of the objective function
- Ignore curvature (higher order derivatives)

Procedure

- Specify starting point θ_0
- Perform update by:

$$\theta_{t+1} = \theta_t + \eta_t \mathbf{d}_t$$

where:

η_t

Step size or learning rate

\mathbf{d}_t

Descent direction

Gradient Direction

Using Taylor expansion we have:

$$\mathcal{L}(\boldsymbol{\theta} + \epsilon\boldsymbol{\lambda}) \simeq \mathcal{L}(\boldsymbol{\theta}) + \epsilon\mathbf{g}^T(\boldsymbol{\theta})\boldsymbol{\lambda} + \mathcal{O}(\epsilon^2)$$

Thus if we assume $\boldsymbol{\lambda} = -\mathbf{g}(\boldsymbol{\theta})$, then for a small enough ϵ , we have:

$$\mathcal{L}(\boldsymbol{\theta} + \epsilon\boldsymbol{\lambda}) - \mathcal{L}(\boldsymbol{\theta}) \simeq -\epsilon\|\mathbf{g}(\boldsymbol{\theta})\|^2 \leq 0$$

So $-\mathbf{g}(\boldsymbol{\theta})$ is a descent direction.

Gradient Descent

Gradient Descent (GD) method uses the following direction:

$$\mathbf{d}_t = -\mathbf{g}(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}_t}$$

Momentum

Momentum method uses the following direction:

$$\begin{aligned}\mathbf{m}_t &= \beta \mathbf{m}_{t-1} + \mathbf{g}_{t-1} \\ \boldsymbol{\theta}_t &= \boldsymbol{\theta}_{t-1} - \eta_t \mathbf{m}_t\end{aligned}$$

where \mathbf{m}_t is the momentum vector and $\beta < 1$

Momentum as Generalization of GD

For $\beta = 1$, momentum method degenerated to GD method.

Learning Rate Schedule

The sequence of step sizes $\{\eta_t\}$ is called the learning rate schedule.

Sample Schedules

- Constant: $\eta_t = \eta$
 - Too large values may fail convergence
 - Too small values lead to low convergence rate
- Armijo-Goldstein: Assume $m = \langle \mathbf{g}(\boldsymbol{\theta}_t), \mathbf{d}_t \rangle$ and select $\tau \in (0, 1)$, $c \in (0, 1)$ and η_0 , then:
 - $\gamma = -cm$ and $j = 0$
 - Until $\mathcal{L}(\boldsymbol{\theta}_t) - \mathcal{L}(\boldsymbol{\theta}_t + \eta_{t,j}\mathbf{d}_t) \geq \eta_{t,j}\gamma$, increment j and set $\eta_{t,j} = \tau\eta_{t,(j-1)}$
 - $\eta_{t,j}$ is the learning rate at iteration t .

Section 3

Second Order Methods

Descent Direction

Assume $\mathbf{H}_t \triangleq \nabla^2 \mathcal{L}(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}_t} \succ 0$. Then Second order approximation of $\mathcal{L}(\cdot)$ in $\boldsymbol{\theta} = \boldsymbol{\theta}_t$ is:

$$\mathcal{L}(\boldsymbol{\theta}) \simeq \mathcal{L}(\boldsymbol{\theta}_t) + \mathbf{g}_t^T (\boldsymbol{\theta} - \boldsymbol{\theta}_t) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_t)^T \mathbf{H}_t (\boldsymbol{\theta} - \boldsymbol{\theta}_t)$$

The minimizer for the above approximation is: $\boldsymbol{\theta}_t - \mathbf{H}_t^{-1} \mathbf{g}_t$. Thus:

$$\Rightarrow \mathbf{d}_t = -\mathbf{H}_t^{-1} \mathbf{g}_t$$

Algorithm 0: Optimization based on descent direction

Input : t_{\max} (Maximum iterations),
 $f_d(\cdot)$ (direction function),
 $f_l(\cdot)$ (learning rate function)

Initialization: $t = 0$, θ_0 , $flag_c = \text{True}$

begin

while $flag_c$ **do**

$\mathbf{d}_t = f_d(\theta_t)$

$\eta_t = f_l(\theta_t)$

$\theta_{t+1} = \theta_t - \eta_t \mathbf{d}_t$

$t \leftarrow t + 1$

if $\|\mathbf{g}_{t+1}\| \leq \delta$ **or** $t > t_{\max}$ **then**

$flag_c \leftarrow \text{False}$

end

end

end

Output : θ_t

Section 4

Stochastic Gradient Descent

Stochastic Gradient Descent

Loss Measurement Limitation

Previously we have seen that Gradient Descent (GD) method uses $\mathbf{d}_t = -\frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_t}$.

Now assume you only have access to a noisy version of loss function, denoted $\mathcal{L}(\boldsymbol{\theta}, \mathbf{z}_t)$, where $\mathbf{z}_t \sim q$ and we have:

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{q(\mathbf{z})}[\mathcal{L}(\boldsymbol{\theta}, \mathbf{z})]$$

Stochastic gradient descent is a solution to the aforementioned problem.

Stochastic Gradient Descent

The update rule for stochastic gradient descent is:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \nabla \mathcal{L}(\boldsymbol{\theta}_t, \mathbf{z}_t) = \boldsymbol{\theta}_t - \eta_t \mathbf{g}_t$$

The sequence $\{\boldsymbol{\theta}_t\}$ is guaranteed to converge to a stationary point provided:

- The step size η_t is decayed at a certain rate
- \mathbf{z} is independent of $\boldsymbol{\theta}$

Section 5

Constrained Optimization

Convex Set

Set \mathcal{S} is convex if, for any $\mathbf{x}, \mathbf{x}' \in \mathcal{S}$, we have:

$$\lambda \mathbf{x} + (1 - \lambda) \mathbf{x}' \in \mathcal{S}, \forall \lambda \in [0, 1]$$



(a) Convex sets



(b) Nonconvex sets

Convex Function

Function $f(\mathbf{x})$ is convex if it is defined on a convex set \mathcal{S} and if, for any $\mathbf{x}, \mathbf{y} \in \mathcal{S}$, and for any $0 \leq \lambda \leq 1$, we have:

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y})$$

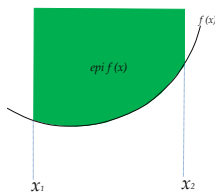


Figure: Convexity check based on epigraph

Hessian Matrix of Convex Function

A twice differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if and only if the Hessian $\nabla^2 f(\mathbf{x})$ is positive semi-definite for all $\mathbf{x} \in \mathbb{R}^n$ [1].

Hessian Matrix of Convex Function

Assume that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and differentiable. Then \mathbf{x}^* is a global minimizer of $f(\cdot)$, if and only if $\nabla f(\mathbf{x}^*) = \mathbf{0}$ [1].

Constrained Optimization

A constrained optimization is defined as:

$$\theta^* \in \underset{\theta \in \mathcal{C}}{\operatorname{argmin}} \mathcal{L}(\theta)$$

where $\mathcal{C} = \{\theta : g_j(\theta) \leq 0, j \in \mathcal{I} \text{ and } h_k(\theta) = 0, k \in \epsilon\} \subseteq \mathbb{R}^D$. The above optimization problem is unconstrained if $\mathcal{C} \in \mathbb{R}^D$, otherwise it is constrained.

Simple Case with One Equality Constraint

Assume we have $\theta^* \in \operatorname{argmin}_{h(\theta)=0} \mathcal{L}(\theta)$. Then:

- $\nabla h(\theta)$ is orthogonal to constraint surface because:

$$\begin{cases} h(\theta + \epsilon) \simeq h(\theta) + \epsilon^T \nabla h(\theta) \\ h(\theta) = h(\theta + \epsilon) \\ \epsilon \parallel \text{constraint surface} \end{cases} \Rightarrow \nabla h(\theta) \perp \text{constraint surface}$$

- If θ^* is optimizer then $\nabla \mathcal{L}(\theta^*) \perp \text{constraint surface}$

Altogether: $\nabla \mathcal{L}(\theta^*) = \lambda^* \nabla h(\theta^*)$, $\lambda^* \in \mathbb{R}$

Lagrange Multiplier

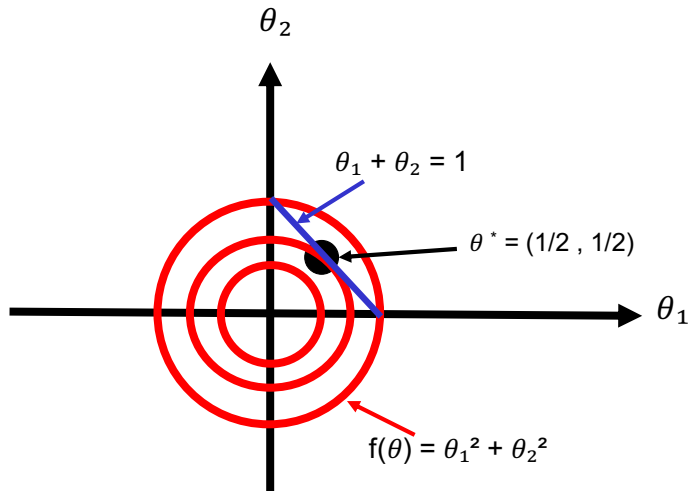


Figure: Solving problem $\theta^* \in \operatorname{argmin}_{\theta_1 + \theta_2 = 1} \theta_1^2 + \theta_2^2$

Lagrangian

Assume Lagrangian as:

$$L(\boldsymbol{\theta}, \lambda) \triangleq \mathcal{L}(\boldsymbol{\theta}) + \lambda h(\boldsymbol{\theta})$$

Then we have:

$$\nabla_{\boldsymbol{\theta}, \lambda} L(\boldsymbol{\theta}, \lambda) = 0 \Leftrightarrow \begin{cases} \lambda \nabla_{\boldsymbol{\theta}} h(\boldsymbol{\theta}) = -\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) \\ h(\boldsymbol{\theta}) = 0 \end{cases}$$

Thus the stationary points of Lagrangian satisfy constraints and lead to parallel gradient vectors.

M Equality constraints

For this case we simply find the stationary points of Lagrangian defined as:

$$L(\boldsymbol{\theta}, \boldsymbol{\lambda}) = \mathcal{L}(\boldsymbol{\theta}) + \sum_{j=1}^m \lambda_j h_j(\boldsymbol{\theta})$$

Constrained Optimization with M Equality Constraints

Assume the following optimization problem:

$$\theta^* \in \underset{\theta \in \mathcal{C}}{\operatorname{argmin}} \mathcal{L}(\theta)$$

$$\mathcal{C} = \{\theta : g_j(\theta) \leq 0, j \in \mathcal{I} \text{ and } h_k(\theta) = 0, k \in \epsilon\} \subseteq \mathbb{R}^D$$

- The necessary condition for θ^* is $L(\theta^*, \lambda^*) = \mathbf{0}$
- If $\mathcal{L}(\theta)$ is convex and equality constraints are Affine ($h_k(\theta) = \mathbf{a}_k \theta = 0$), then the optimization problem is convex and condition $L(\theta^*, \lambda^*) = \mathbf{0}$ is sufficient.

Constrained Optimization

Assume the following optimization problem:

$$\theta^* \in \operatorname{argmin}_{\theta \in \mathcal{C}} \mathcal{L}(\theta)$$

$$\mathcal{C} = \{\theta : g_j(\theta) \leq 0, j \in \mathcal{I} \text{ and } h_k(\theta) = 0, k \in \mathcal{E}\} \subseteq \mathbb{R}^D$$

We define the generalized Lagrangian as:

$$L(\theta, \mu, \lambda) = \mathcal{L}(\theta) + \sum_i \mu_i g_i(\theta) + \sum_j \lambda_j h_j(\theta)$$

Then the KKT (Karush–Kuhn–Tucker) conditions are:

- $\nabla \mathcal{L}(\theta) + \sum_i \mu_i \nabla g_i(\theta) + \sum_j \lambda_j \nabla h_j(\theta) = 0$ (Stationary point of Lagrangian)
- $\mathbf{g}(\theta) \leq 0, \mathbf{h}(\theta) = 0$ (Feasibility)
- $\mu \geq 0$ (Dual feasibility)
- $\mu \odot \mathbf{g} = \mathbf{0}$ (Complementary Slackness)

Constrained Optimization

Again assume the following optimization problem:

$$\theta^* \in \underset{\theta \in \mathcal{C}}{\operatorname{argmin}} \mathcal{L}(\theta)$$

$$\mathcal{C} = \{\theta : g_j(\theta) \leq 0, j \in \mathcal{I} \text{ and } h_k(\theta) = 0, k \in \epsilon\} \subseteq \mathbb{R}^D$$

Then KKT conditions are:

- Necessary for θ
- Sufficient for θ if above problem is convex ($\mathcal{L}(\theta)$ and $\{g_j(\theta)\}_{j \in \mathcal{I}}$ are convex functions and $\{h_k(\theta)\}_{k \in \epsilon}$ are Affine transforms).

KKT Conditions [2]

Consider the following convex optimization problem:

$$\begin{aligned} \min_{(x,y) \in \mathcal{S}} \quad & \frac{1}{x+y} \\ \text{subject to} \quad & \begin{cases} 2x + y^2 - 6 \leq 0 \\ 1 - x \leq 0 \\ 1 - y \leq 0 \end{cases} \end{aligned}$$

where $\mathcal{S} = \{(x, y) : x, y > 0\}$. Find the optimal point.

Solution: From zero Lagrangian gradient we have:

$$\begin{bmatrix} -\frac{1}{(x+y)^2} \\ -\frac{1}{(x+y)^2} \end{bmatrix} + \mu_1 \begin{bmatrix} 2 \\ 2y \end{bmatrix} + \mu_2 \begin{bmatrix} -1 \\ 0 \end{bmatrix} + \mu_3 \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

From complementary slackness equations we have:

$$\mu_1(2x + y^2 - 6) = \mu_2(1 - x) = \mu_3(1 - y) = 0$$

KKT Conditions [2] (Continue)

Assume $\mu_1 = 0$, then:

$$\begin{bmatrix} -\frac{1}{(x+y)^2} \\ \frac{1}{(x+y)^2} \end{bmatrix} + \mu_2 \begin{bmatrix} -1 \\ 0 \end{bmatrix} + \mu_3 \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow \mu_2 = \mu_3 = -\frac{1}{(x+y)^2} \leq 0$$

$$\Rightarrow \begin{cases} \text{Contradiction} \\ 2x + y^2 - 6 = 0 \end{cases}$$

No we assume $x = 1$ then:

$$x = 1 \Rightarrow y = \begin{cases} +2 \text{ (valid)} \\ -2 \text{ (invalid)} \end{cases} \Rightarrow \begin{cases} \mu_1 = \frac{1}{36} \\ \mu_2 = -\frac{1}{18} \\ \mu_3 = 0 \end{cases} \Rightarrow \begin{cases} \text{Contradiction} \\ x \neq 0 \end{cases}$$

No we assume $y = 1$ then:

$$y = 1 \Rightarrow x = 2.5 \Rightarrow \begin{cases} \mu_1 = \frac{2}{49} \\ \mu_2 = 0 \\ \mu_3 = 0 \end{cases} \Rightarrow \begin{cases} \boldsymbol{\theta}^* = (2.5, 1) \\ \boldsymbol{\mu} = (\frac{2}{49}, 0, 0) \end{cases}$$



Markus Grasmair,

“Basic properties of convex functions,”

Department of Mathematics, Norwegian University of Science and Technology, 2016.



“Chapter 5, lecture 6: Kkt theorem, gradient form,”

<https://faculty.math.illinois.edu/~mlavrov/docs/484-spring-2019/ch5lec6.pdf>,

Accessed: 2022-10-26.