

Appendix.02 Optimization Theory

Definition.A.2.1 Optimization Problem with constraint in general programming

Let $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ called an object or cost function,

$g_i(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ called an inequality constraint function, and

$h_j(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ called an equality constraint function be.

The following nonlinear programming is called a optimization problem.

$$\min_{\mathbf{x}} f(\mathbf{x}) \text{ s.t. } g_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, m, \quad h_j(\mathbf{x}) = 0, \quad j = 1, 2, \dots, k.$$

$g_i(\mathbf{x}) \leq 0$ is called an inequality constraint and $h_j(\mathbf{x}) = 0$ is called an equality constraint.

Definition.A.2.2 Lagrangian function

In a optimization problem,

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) + \sum_{j=1}^k \mu_j h_j(\mathbf{x}) \text{ where } \boldsymbol{\lambda} : m \times 1, \boldsymbol{\mu} : k \times 1$$

λ_i and μ_j are also called dual variables or weights. \mathbf{x} is called a primal variable.

Theorem.A.2.1 Karush-Kuhn-Tucker(KKT) conditions

If $\mathbf{x}^* \in \mathbb{R}^n$ is a local minimum,

then $\exists \boldsymbol{\lambda}^* \in \mathbb{R}^m$ and $\boldsymbol{\mu}^* \in \mathbb{R}^k$ s.t

(i) Stationarity :

$$\nabla f(\mathbf{x}) + \sum_i \lambda_i \nabla g_i(\mathbf{x}) + \sum_j \mu_j \nabla h_j(\mathbf{x}) = \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{0}$$

(ii) Primal feasibility :

$$g_i(\mathbf{x}^*) \leq 0, \quad \forall i, \quad h_j(\mathbf{x}^*) = 0, \quad \forall j$$

(iii) Dual feasibility :

$$\lambda_i^* \geq 0, \quad i = 1, \dots, m$$

(iv) Complementary slackness :

$$\lambda_i^* g_i(\mathbf{x}^*) = 0, \quad i = 1, \dots, m$$

This \mathbf{x}^* is also called a local solution of a optimization problem. KKT Conditions are necessary(not sufficient) for global optimality.

$$\mathbf{x}^*(\text{Satisfying the KKT Conditions}) \leftarrow \text{Global optimal}$$

Definition.A.2.3 Convex function

A function $f(\mathbf{x})$ is said to be convex if

$$\forall t \in [0, 1], \quad f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \leq tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2), \quad \forall \mathbf{x}_1, \mathbf{x}_2 \in D$$

Definition.A.2.4 Convex set

A set S is said to be convex if

$$\mathbf{x}, \mathbf{y} \in S \rightarrow t\mathbf{x} + (1-t)\mathbf{y} \in S, \quad t \in [0, 1]$$

Definition.A.2.5 Convex optimization

Let $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ and $\mathbf{g}(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be convex functions and $\mathbf{h}(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^k$ be an affine function.

An optimization problem is called convex if

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad s.t. \quad \mathbf{g}_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m, \quad \mathbf{h}_j(\mathbf{x}) = 0, \quad j = 1, \dots, k.$$

Theorem.A.2.2 KKT Conditions in convex problems

For convex problem, KKT conditions becomes necessary and also sufficient for global optimality.

(i) Stationarity :

$$\nabla f(\mathbf{x}^*) + \nabla (\boldsymbol{\lambda}^*{}^T \mathbf{h}(\mathbf{x}^*)) + \nabla (\boldsymbol{\mu}^*{}^T \mathbf{g}(\mathbf{x}^*)) = \mathbf{0}$$

(ii) Primal feasibility :

$$\mathbf{g}(\mathbf{x}^*) \leq \mathbf{0}, \quad \mathbf{h}(\mathbf{x}^*) = \mathbf{0}$$

(iii) Dual feasibility :

$$\boldsymbol{\lambda}^* \geq \mathbf{0}$$

(vi) Complementary slackness :

$$\boldsymbol{\lambda}^*{}^T \mathbf{g}(\mathbf{x}^*) = 0$$

Definition.A.2.6 Lagrange dual function

$$\begin{aligned} \mathcal{D}(\boldsymbol{\lambda}, \boldsymbol{\mu}) &= \min_{\mathbf{x} \in \mathcal{X}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \\ &= \min_{\mathbf{x} \in \mathcal{X}} \{f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{h}(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x})\} \end{aligned}$$

where $\mathcal{X} = \{\mathbf{x} : \mathbf{g}(\mathbf{x}) \leq \mathbf{0}, \quad \mathbf{h}(\mathbf{x}) = \mathbf{0}\}$

Definition.A.2.7 Lagrange dual problem

$$\begin{aligned} \max_{\boldsymbol{\lambda} \geq \mathbf{0}, \boldsymbol{\mu}} \mathcal{D}(\boldsymbol{\lambda}, \boldsymbol{\mu}) &= \max_{\boldsymbol{\lambda} \geq \mathbf{0}} \min_{\boldsymbol{\mu}} \min_{\mathbf{x} \in \mathcal{X}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \\ &= \max_{\boldsymbol{\lambda} \geq \mathbf{0}} \min_{\mathbf{x} \in \mathcal{X}} \{f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{h}(\mathbf{x})\} \end{aligned}$$

Theorem.A.2.3 Optimization-theoretic solution of distance between a vector and a hyperplane

Let a vector \mathbf{x}_0 and a hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$ be.

$$\min_{\mathbf{x}} \|\mathbf{x}_0 - \mathbf{x}\|^2 \text{ s.t. } \mathbf{w}^T \mathbf{x} + b = 0$$

$$\therefore \|\mathbf{x}_0 - \mathbf{x}^*\|^2 = \frac{(\mathbf{w}^T \mathbf{x}_0 + b)^2}{\|\mathbf{w}\|^2}$$

Proof)

Let the function $\mathcal{L}(\mathbf{x}, \lambda, \mu) = \|\mathbf{x}_0 - \mathbf{x}\|^2 + \mu(\mathbf{w}^T \mathbf{x} + b)$

(i)

$$\nabla_{\mathbf{x}} \mathcal{L} = 2(\mathbf{x}_0 - \mathbf{x}^*) + \mu \mathbf{w} = \mathbf{0}$$

$$2\mathbf{w}^T(\mathbf{x}_0 - \mathbf{x}^*) = -\mu \|\mathbf{w}\|^2 \quad (\because \text{inner product with } (\mathbf{x}_0 - \mathbf{x}^*))$$

$$2\|\mathbf{x}_0 - \mathbf{x}^*\|^2 = \frac{2(\mathbf{w}^T(\mathbf{x}_0 - \mathbf{x}^*))^2}{\|\mathbf{w}\|^2}$$

$$\therefore \|\mathbf{x}_0 - \mathbf{x}^*\|^2 = \frac{\{\mathbf{w}^T(\mathbf{x}_0 - \mathbf{x}^*)\}^2}{\|\mathbf{w}\|^2} \quad \blacksquare$$

Theorem.A.2.4 Gradient descent method

Gradient descent is an iterative method to find a stationary point of an unconstraint optimization problem :

$$\min_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$

$$L(\boldsymbol{\theta} + \eta \mathbf{d}) \approx L(\boldsymbol{\theta}) + \eta \nabla_{\boldsymbol{\theta}}^T L(\boldsymbol{\theta}) \mathbf{d} \quad \text{where } \eta > 0, \|\mathbf{d}\| = 1$$

$$L(\boldsymbol{\theta} + \eta \mathbf{d}) - L(\boldsymbol{\theta}) \approx \eta \nabla_{\boldsymbol{\theta}}^T L(\boldsymbol{\theta}) \mathbf{d} = \eta \cos(\phi) \|\nabla_{\boldsymbol{\theta}}^T L(\boldsymbol{\theta})\|$$

Find the directional vector \mathbf{d} that minimizes $L(\boldsymbol{\theta} + \eta \mathbf{d}) - L(\boldsymbol{\theta}) \leq 0$

$$\cos(\phi) = -1 \rightarrow \mathbf{d} = -\frac{\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})}{\|\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})\|}$$

$$\therefore \boldsymbol{\theta} + \eta \mathbf{d} = \boldsymbol{\theta} - \eta \frac{\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})}{\|\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})\|} = \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$

Theorem.A.2.4 Types of gradient descent method

(i) Standard (or steepest) Gradient Descent

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla \mathbb{E}[J(\mathbf{w})]$$

- Practically infeasible
- Thus, we need distribution about data \mathbf{x} (Contradiction)
- So, We can use sample mean

(ii) Stochastic Gradient Descent

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla J_i(\mathbf{w})$$

- Simple to implement
- Effective for large-scale problem
- Much less memory
- Unstable : zigzagging
- Purpose : We just consider one of data. Just one.
- It can be convergent. But there is little unstable.

(iii) Batch gradient Descent

$$\mathbf{w} \leftarrow \eta \nabla \sum_{i=1}^N J_i(\mathbf{w})$$

- Accurate estimation of gradients
- Parallelization of learning
- Large memory
- Big time-complexity can be problem in this method.(So slow)
- But, there isn't problem in convergence.
- Purpose : We consider all of data!

(vi) Mini-Batch Gradient Descent

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla \sum_{i \in \mathcal{F}} J_i(\mathbf{w}), \quad 1 \leq |\mathcal{F}| \leq N$$

- Most generalized version
- Effective to deal with large
- Amount of training data
- Purpose : We just consider several datas.

Reference :

Deep Learning - Yosha Benjio

Neural networks and learning machines - Simon Haykin

<https://wikipedia.org> (<https://wikipedia.org>)

<https://ko.wikipedia.org/wiki/%EC%9C%84%ED%82%A4%EB%B0%B1%EA%B3%BC%EB%AC%B8%EB%B2%95>

<https://ko.wikipedia.org/wiki/%EC%9C%84%ED%82%A4%EB%B0%B1%EA%B3%BC%EB%AC%B8%EB%B2%95>