



Method of Lagrange

Cheng Soon Ong
Marc Peter Deisenroth

December 2020



Motivation

- ▶ In machine learning, we use gradients to train predictors
- ▶ For functions $f(\boldsymbol{x})$ we can directly obtain its gradient $\nabla f(\boldsymbol{x})$
- ▶ We may wish to calculate a gradient for functions with constraints
- ▶ Assume that we can calculate gradient for corresponding unconstrained problem

Key idea

How to take a gradient with respect to a constrained optimization problem.

History

- ▶ Joseph-Louis Lagrange (b. 1736) worked on calculus of variations and founded the idea
- ▶ Nonlinear programming was coined by Harold W. Kuhn and Albert W. Tucker in 1951.
- ▶ Turns out William Karush already wrote about similar conditions in his master's thesis in 1939, and was reintroduced into the literature by Akira Takayama in 1974.
- ▶ Classic backpropagation paper, by Yann LeCun in 1988, observed the relation between backpropagation and Lagrange multipliers.

Constrained optimization

Primal optimization problem

Given $f : \mathbb{R}^D \rightarrow \mathbb{R}$ and $g_i : \mathbb{R}^D \rightarrow \mathbb{R}$ for $i = 1, \dots, m$,

$$\begin{array}{ll} \min_{\mathbf{x}} & f(\mathbf{x}) \\ \text{subject to} & g_i(\mathbf{x}) \leq 0 \quad \text{for all } i = 1, \dots, m \end{array}$$

Lagrange multipliers as relaxation

We can convert the primal optimization problem to an unconstrained problem.

$$J(\mathbf{x}) = f(\mathbf{x}) + \sum_{i=1}^m \mathbf{1}(g_i(\mathbf{x})),$$

where $\mathbf{1}(z)$ is an infinite step function

$$\mathbf{1}(z) = \begin{cases} 0 & \text{if } z \leq 0 \\ \infty & \text{otherwise} \end{cases}.$$

Key idea

Approximate indicator function by linear function.

Lagrangian

We associate the primal optimization problem with a **Lagrangian**, by introducing Lagrange multipliers λ_i for each constraint g_i .

Lagrangian

$$\begin{aligned}\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) &= f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) \\ &= f(\mathbf{x}) + \boldsymbol{\lambda}^\top \mathbf{g}(\mathbf{x})\end{aligned}$$

Why Lagrangian?

- ▶ $\min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ is an unconstrained optimization problem for a given value of $\boldsymbol{\lambda}$
- ▶ Turns out: If solving $\min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ is easy, then the overall problem is easy

Dual optimization problem

Key idea

Convert problem in $\mathbf{x} \in \mathbb{R}^D$, to another problem in $\boldsymbol{\lambda} \in \mathbb{R}^m$.

- ▶ Primal and dual problems are related
- ▶ The objective values at optimality related via gradients
- ▶ the dual problem (maximization over $\boldsymbol{\lambda}$) is a maximum over a set of affine functions, and hence is a concave function, even though $f(\cdot)$ and $g_i(\cdot)$ may be non-convex

Primal dual pair

The problem

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{subject to} \quad & g_i(\mathbf{x}) \leq 0 \quad \text{for all } i = 1, \dots, m \end{aligned}$$

is known as the **primal problem**, corresponding to the primal variables \mathbf{x} . The associated **Lagrangian dual problem** is given by

$$\begin{aligned} \max_{\boldsymbol{\lambda} \in \mathbb{R}^m} \quad & \mathfrak{D}(\boldsymbol{\lambda}) \\ \text{subject to} \quad & \boldsymbol{\lambda} \geq \mathbf{0}, \end{aligned}$$

where $\boldsymbol{\lambda}$ are the dual variables and $\mathfrak{D}(\boldsymbol{\lambda}) = \min_{\mathbf{x} \in \mathbb{R}^d} \mathfrak{L}(\mathbf{x}, \boldsymbol{\lambda})$.

Minimax inequality

Why Lagrangian duality is useful? Because **primal values are greater than dual values**.

Key idea

For any function with two arguments, the maximin is less than the minimax.

$$\max_{\mathbf{y}} \min_{\mathbf{x}} \varphi(\mathbf{x}, \mathbf{y}) \leq \min_{\mathbf{x}} \max_{\mathbf{y}} \varphi(\mathbf{x}, \mathbf{y}).$$

Weak duality

Key idea

Primal values $f(\mathbf{x})$ are always greater than (or equal to) dual values $\mathfrak{D}(\boldsymbol{\lambda})$

We can see this by applying the minimax inequality

$$\min_{\mathbf{x} \in \mathbb{R}^d} \underbrace{\max_{\boldsymbol{\lambda} \geq 0} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})}_{f(\mathbf{x})} \geq \max_{\boldsymbol{\lambda} \geq 0} \underbrace{\min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})}_{\mathfrak{D}(\boldsymbol{\lambda})} .$$

Motivating the Karush Kuhn Tucker conditions

Lagrangian

$$\mathcal{L}(\boldsymbol{x}, \boldsymbol{\lambda}) = f(\boldsymbol{x}) + \sum_{i=1}^m \lambda_i g_i(\boldsymbol{x})$$

Stationarity

Solve for optimal value by taking gradient w.r.t. \boldsymbol{x} and set to zero

Lagrange multipliers are orthogonal to constraints

- ▶ Recall that we have one Lagrange multiplier per constraint
- ▶ If $\lambda_i = 0$, then $g_i(\boldsymbol{x})$ can be any value (has slack)
- ▶ If $g_i(\boldsymbol{x}) = 0$ (active), then λ_i has slack

Retain feasibility

Karush Kuhn Tucker conditions

stationarity

$$\mathbf{0} \in \nabla f(\mathbf{x}) + \sum_{i=1}^m \lambda_i \nabla g_i(\mathbf{x}). \quad (1)$$

complementary slackness

$$\lambda_i g_i(\mathbf{x}) = 0 \quad \text{for all } i = 1, \dots, m. \quad (2)$$

primal feasibility

$$g_i(\mathbf{x}) \leq 0 \quad \text{for all } i = 1, \dots, m. \quad (3)$$

dual feasibility

$$\lambda_i \geq 0 \quad \text{for all } i = 1, \dots, m. \quad (4)$$

Constrained optimization

Key idea

Objective gradient is a conic combination of active constraint gradients

By stationarity condition

$$\nabla f(\mathbf{x}) = - \sum_{i=1}^m \lambda_i \nabla g_i(\mathbf{x})$$

for a given value of λ .

KKT implies strong duality

Key idea

KKT conditions are sufficient for zero duality gap

Proof sketch:

- ▶ By weak duality, $f(\mathbf{x}) \geq \mathfrak{D}(\boldsymbol{\lambda})$. Let \mathbf{x}^* and $\boldsymbol{\lambda}^*$ denote optimal values
- ▶ By stationarity conditions

$$\mathfrak{D}(\boldsymbol{\lambda}^*) = f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* g_i(\mathbf{x}^*),$$

- ▶ By complementary slackness

$$f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* g_i(\mathbf{x}^*) = f(\mathbf{x}^*).$$

- ▶ Hence duality gap is zero

KKT implies strong duality

Key idea

Challenge is to express geometric constraints as algebraic equations

Constraint qualification additionally needed for KKT to necessarily imply zero duality gap.

- ▶ Mangasarian Fromovitz
- ▶ Linear independence
- ▶ Fritz John
- ▶ Slater's condition

Summary

Key idea

How to take a gradient with respect to a constrained optimization problem.

- ▶ Objective gradient is a conic combination of active constraint gradients

$$\nabla f(\mathbf{x}) = - \sum_{i=1}^m \lambda_i \nabla g_i(\mathbf{x})$$

- ▶ Check that one of Lagrange multipliers or primal constraints is zero
- ▶ The KKT conditions are sufficient for a zero duality gap