

Statistical Optimization: Lecture 12

Penalty and Augmented Lagrangian Methods

Zijian Guo

Zhejiang University
Center for Data Science

April 21, 2026

Objective

Goal. We want to solve constrained optimization problems.

Main idea. We will compare three viewpoints:

- dual updates for enforcing feasibility,
- penalty terms for discouraging constraint violation,
- the augmented Lagrangian method, which combines both.

Route of this lecture.

- Dual-ascent intuition
- Quadratic penalty method
- Augmented Lagrangian method

Problem setup

$$\begin{aligned} \min_{\theta} \quad & f_0(\theta) \\ \text{s.t.} \quad & h_i(\theta) = 0, \quad i \in \mathcal{E}. \end{aligned} \tag{1}$$

This is the **equality-constrained problem**.

$$\begin{aligned} \min_{\theta} \quad & f_0(\theta) \\ \text{s.t.} \quad & h_i(\theta) = 0, \quad i \in \mathcal{E}, \\ & f_i(\theta) \leq 0, \quad i \in \mathcal{I}. \end{aligned} \tag{2}$$

This is the **general constrained problem**.

Outline

Dual-ascent intuition

Quadratic penalty method

Augmented Lagrangian method

Ordinary Lagrangian and dual-ascent idea

For the equality-constrained problem (1), define the ordinary Lagrangian

$$L(\theta, \nu) := f_0(\theta) + \sum_{i \in \mathcal{E}} \nu_i h_i(\theta). \quad (3)$$

The associated dual function is

$$g(\nu) := \inf_{\theta} L(\theta, \nu), \quad \max_{\nu} g(\nu). \quad (4)$$

If

$$\theta(\nu) \in \arg \min_{\theta} L(\theta, \nu),$$

then, under suitable conditions, Danskin's theorem gives

$$\nabla g(\nu) = h(\theta(\nu)). \quad (5)$$

Therefore, gradient ascent on the dual function suggests the update

$$\nu^{k+1} = \nu^k + \alpha_k h(\theta_k), \quad \theta_k \approx \arg \min_{\theta} L(\theta, \nu^k), \quad \alpha_k > 0. \quad (6)$$

Why plain dual ascent is not enough

Without strong convexity, the minimizer of $L(\theta, \nu)$ with respect to θ may fail to be unique.

Then the dual function may become less regular.

As a result:

- convergence may be very slow,
- the iterates may oscillate,
- with poor step-size choices, the method may even fail to converge.

This motivates adding a **quadratic penalty term** to stabilize the primal step.

Outline

Dual-ascent intuition

Quadratic penalty method

Augmented Lagrangian method

Equality-constrained: quadratic penalty method

For the equality-constrained problem (1), define the quadratic penalty function

$$Q(\theta; \mu) := f_0(\theta) + \frac{\mu}{2} \sum_{i \in \mathcal{E}} h_i(\theta)^2, \quad \mu > 0. \quad (7)$$

Idea. As μ increases, violations of the equality constraints are penalized more severely.

General constrained: quadratic penalty method

For the general constrained problem (2), define

$$Q(\theta; \mu) := f_0(\theta) + \frac{\mu}{2} \sum_{i \in \mathcal{E}} h_i(\theta)^2 + \frac{\mu}{2} \sum_{i \in \mathcal{I}} ([f_i(\theta)]_+)^2, \quad \mu > 0, \quad (8)$$

where

$$[y]_+ := \max(y, 0).$$

The inequality term penalizes only the violated constraints.

Remarks on quadratic penalty subproblems

- The quadratic penalty subproblem is generally an unconstrained optimization problem, but it need not be convex.
- Even if f_0 , h_i , and f_i are smooth, Q may be less smooth because of the inequality penalty term.
- For the general constrained problem, the function is typically still differentiable, but the Hessian may fail to exist at points where

$$f_i(\theta) = 0,$$

because the term $([f_i(\theta)]_+)^2$ is usually not twice differentiable there.

Pseudo-code (Quadratic penalty method)

```
Input:  $\mu_0 > 0$ , tolerances  $\tau_k \downarrow 0$ , initial point  $\theta_0^s$ 
for  $k = 0, 1, 2, \dots$  do
    approximately minimize  $Q(\theta; \mu_k)$  from  $\theta_k^s$ 
    until  $\|\nabla_{\theta} Q(\theta; \mu_k)\| \leq \tau_k$ 
    set  $\theta_k \leftarrow$  obtained approximate minimizer
    if  $\|h(\theta_k)\| \leq \varepsilon_{\text{feas}}$  then
        return  $\theta_k$ 
    end if
    choose  $\mu_{k+1} > \mu_k$ 
    choose new starting point  $\theta_{k+1}^s$ 
end for
```

In convergence analysis, we restrict our attention to the equality-constrained problem (1).

Theorem (Ideal convergence behavior)

Suppose that each θ_k is the exact global minimizer of (7) and that

$$\mu_k \uparrow \infty.$$

Then every limit point θ^ of the sequence $\{\theta_k\}$ is a global minimizer of the equality-constrained problem*

$$\min_{\theta} f_0(\theta) \quad \text{s.t.} \quad h(\theta) = 0.$$

Theorem (Algorithm convergence behavior)

Suppose that the tolerances and penalty parameters satisfy

$$\tau_k \rightarrow 0, \quad \mu_k \uparrow \infty.$$

Let θ^* be a limit point of the sequence $\{\theta_k\}$.

- If θ^* is infeasible, then it is a stationary point of $\|h(\theta)\|^2$.
- If θ^* is feasible and the constraint gradients $\nabla h_i(\theta^*)$ are linearly independent, then θ^* is a KKT point.

Moreover, along any infinite subsequence \mathcal{K} with $\theta_k \rightarrow \theta^*$ ($k \in \mathcal{K}$), we have

$$\mu_k h_i(\theta_k) \rightarrow \nu_i^*, \quad i \in \mathcal{E},$$

where ν^* is the multiplier vector that satisfies the KKT conditions for the equality-constrained problem.

Intuition for the limit $\mu_k h_i(\theta_k) \rightarrow \nu^*$

At the approximate minimizer θ_k of the quadratic penalty subproblem, we have

$$\nabla Q(\theta_k; \mu_k) \approx 0.$$

This gives

$$\nabla f_0(\theta_k) + \sum_{i \in \mathcal{E}} \mu_k h_i(\theta_k) \nabla h_i(\theta_k) \approx 0.$$

Compare this with the KKT stationarity condition at a feasible limit point θ^* :

$$\nabla f_0(\theta^*) + \sum_{i \in \mathcal{E}} \nu_i^* \nabla h_i(\theta^*) = 0.$$

So the quantity

$$\mu_k h_i(\theta_k)$$

plays the role of an *approximate multiplier*. Hence, along a subsequence with $\theta_k \rightarrow \theta^*$, it is natural to expect

$$\mu_k h_i(\theta_k) \rightarrow \nu_i^*.$$

Main drawback of quadratic penalty

Interpretation of the convergence result.

- The sequence may be attracted to infeasible stationary points.
- Even at feasible limit points, the method only guarantees convergence to KKT points.
- The quantities $\mu_k h_i(\theta_k)$ behave like multiplier estimates.

Main drawback.

- As μ_k increases, the Hessian matrix becomes increasingly **ill-conditioned**.
- The quadratic penalty function is usually **not exact**: even for a fixed $\mu > 0$, its minimizer does not necessarily coincide with the constrained solution.

Outline

Dual-ascent intuition

Quadratic penalty method

Augmented Lagrangian method

From dual ascent and quadratic penalty to ALM

The dual-ascent update (6) is natural, but it behaves best when the θ -subproblem is sufficiently well conditioned.

The quadratic penalty function (7) stabilizes the primal step, and at an approximate minimizer θ_k of the penalty subproblem we have

$$\nabla f_0(\theta_k) + \sum_{i \in \mathcal{E}} \mu_k h_i(\theta_k) \nabla h_i(\theta_k) \approx 0.$$

So the quantity $\mu_k h_i(\theta_k)$ already behaves like an **approximate multiplier**.

This suggests keeping an explicit multiplier estimate ν^k , while still using the quadratic term to stabilize the primal step. This leads to

$$L_A(\theta, \nu; \mu) := f_0(\theta) + \sum_{i \in \mathcal{E}} \nu_i h_i(\theta) + \frac{\mu}{2} \sum_{i \in \mathcal{E}} h_i(\theta)^2. \quad (9)$$

So ALM can be viewed both as a **stabilized dual-ascent method** and as a **quadratic penalty method with multiplier correction**.

Multiplier update

At iteration k , fix $\mu_k > 0$ and the current multiplier ν^k , and approximately minimize

$$L_A(\theta, \nu^k; \mu_k) \quad \text{with respect to } \theta.$$

If θ_k is an approximate minimizer, then

$$0 \approx \nabla_{\theta} L_A(\theta_k, \nu^k; \mu_k) = \nabla f_0(\theta_k) + \sum_{i \in \mathcal{E}} (\nu_i^k + \mu_k h_i(\theta_k)) \nabla h_i(\theta_k). \quad (10)$$

Comparing (10) with the KKT stationarity condition suggests defining

$$\nu_i^{k+1} := \nu_i^k + \mu_k h_i(\theta_k), \quad i \in \mathcal{E}. \quad (11)$$

This motivates the multiplier update

$$\nu^{k+1} = \nu^k + \mu_k h(\theta_k). \quad (12)$$

Main point. If ν^k is already close to ν^* , then we can often obtain good feasibility with a moderate μ_k , avoiding the severe ill-conditioning of pure quadratic penalty.

Pseudo-code (Augmented Lagrangian method)

```
Input:  $\mu_0 > 0$ , tolerance  $\tau_0 > 0$ , initial point  $\theta_0^s$ , initial multiplier  $\nu^0$ 
for  $k = 0, 1, 2, \dots$  do
    approximately minimize  $L_A(\theta, \nu^k; \mu_k)$  from  $\theta_k^s$ 
    until  $\|\nabla_{\theta} L_A(\theta, \nu^k; \mu_k)\| \leq \tau_k$ 
    set  $\theta_k \leftarrow$  obtained approximate minimizer
    if  $\|h(\theta_k)\| \leq \varepsilon_{\text{feas}}$  then
        return  $\theta_k$ 
    end if
    update multipliers:  $\nu^{k+1} = \nu^k + \mu_k h(\theta_k)$ 
    choose new penalty parameter  $\mu_{k+1} \geq \mu_k$ 
    set  $\theta_{k+1}^s \leftarrow \theta_k$ 
    choose new tolerance  $\tau_{k+1}$ 
end for
```

Summary of augmented Lagrangian method

- Unlike the quadratic penalty method, ALM improves the iterates not only by increasing μ , but also by updating the multiplier estimate ν^k .
- Therefore, good feasibility and accuracy can often be obtained with a moderate penalty parameter μ .
- A smaller μ helps avoid the severe ill-conditioning that appears in the quadratic penalty Hessian, so the subproblems are usually easier to solve.
- In this sense, ALM keeps the penalty idea, while reducing one of its main numerical drawbacks.
- Conceptually, ALM can be viewed as a stabilized dual-ascent method: the multiplier is still updated by the current constraint violation, while the primal step uses the augmented Lagrangian.