

## Constrained Optimization

Winter 2023

Clayton Scott

Constrained optimization problems abound in machine learning. In these notes we introduce the basics of constrained optimization, with an emphasis on Lagrange multiplier theory. This will give us enough machinery to derive the support vector machine and understand several other machine learning algorithms that are formulated as constrained optimization problems.

A constrained optimization problem has the form

$$\begin{aligned} \min_{\mathbf{u}} \quad & f(\mathbf{u}) \\ \text{s.t.} \quad & g_i(\mathbf{u}) \leq 0, \quad i = 1, \dots, r \\ & h_j(\mathbf{u}) = 0, \quad j = 1, \dots, s \end{aligned}$$

where  $\mathbf{u} \in \mathbb{R}^p$  and  $f, g_1, \dots, g_r, h_1, \dots, h_s : \mathbb{R}^p \rightarrow \mathbb{R}$ . As an example, the optimal soft-margin hyperplane

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{n} \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad \forall i \\ & \xi_i \geq 0, \quad \forall i \end{aligned} \tag{1}$$

can be expressed in this form where

$$\mathbf{u} = [\mathbf{w}^T \ b \ \xi_1 \ \dots \ \xi_n]^T,$$

$p = d + 1 + n$ ,  $r = 2n$ , and  $s = 0$ .

If  $\mathbf{u}$  satisfies all of the constraints, it is said to be *feasible*. The set of all feasible points is called the *feasible set*, and is assumed to be nonempty. Also assume  $f$  is defined on the feasible set..

## 1 The Lagrangian

The *Lagrangian* is the function

$$L(\mathbf{u}, \boldsymbol{\lambda}, \boldsymbol{\nu}) := f(\mathbf{u}) + \sum_{i=1}^r \lambda_i g_i(\mathbf{u}) + \sum_{j=1}^s \nu_j h_j(\mathbf{u})$$

and  $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_r]^T$  and  $\boldsymbol{\nu} = [\nu_1, \dots, \nu_s]^T$  are called *Lagrange multiplier* or *dual variables*.

The *Lagrange dual function* is

$$L_D(\boldsymbol{\lambda}, \boldsymbol{\nu}) := \min_{\mathbf{u}} L(\mathbf{u}, \boldsymbol{\lambda}, \boldsymbol{\nu})$$

**Note 1.**  $L_D$  is concave, being the point wise minimum of a family of affine functions. See Figure 1.

The *dual* optimization problem is

$$\max_{\boldsymbol{\lambda}, \boldsymbol{\nu}: \lambda_i \geq 0} L_D(\boldsymbol{\lambda}, \boldsymbol{\nu})$$

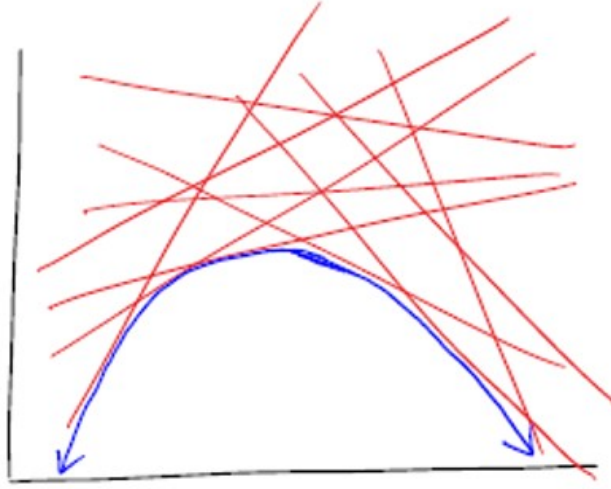


Figure 1: Illustration of  $L_D(\mathbf{u}, \boldsymbol{\nu})$

Similarly, the *primal function* is

$$L_p(\mathbf{u}) := \max_{\boldsymbol{\lambda}, \boldsymbol{\nu}: \lambda_i \geq 0} L(\mathbf{u}, \boldsymbol{\lambda}, \boldsymbol{\nu})$$

and the *primal optimization problem* is

$$\min_{\mathbf{u}} L_p(\mathbf{u})$$

Notice that

$$L_p(\mathbf{u}) = \begin{cases} f(\mathbf{u}) & \text{if } \mathbf{u} \text{ is feasible} \\ \infty & \text{otherwise} \end{cases}$$

and therefore, the primal problem and the original problem have the same solution(s) and optimal objective values, yet the primal problem is unconstrained.

## 2 Duality

Denote the optimal objective function values of the primal problem and dual problem by

$$\begin{aligned} p^* &= \min_{\mathbf{u}} L_p(\mathbf{u}) = \min_{\mathbf{u}} \max_{\boldsymbol{\lambda}, \boldsymbol{\nu}: \lambda_i \geq 0} L(\mathbf{u}, \boldsymbol{\lambda}, \boldsymbol{\nu}) \\ d^* &= \max_{\boldsymbol{\lambda}, \boldsymbol{\nu}: \lambda_i \geq 0} L_D(\mathbf{u}) = \max_{\boldsymbol{\lambda}, \boldsymbol{\nu}: \lambda_i \geq 0} \min_{\mathbf{u}} L(\mathbf{u}, \boldsymbol{\lambda}, \boldsymbol{\nu}). \end{aligned}$$

These values are related.

### 2.1 Weak Duality

The following property is always true, and is referred to as *weak duality*.

**Proposition 1.**  $d^* \leq p^*$

*Proof.* Let  $\tilde{\mathbf{u}}$  be feasible. Then for any  $\boldsymbol{\lambda}, \boldsymbol{\nu}$  with  $\lambda_i \geq 0$

$$L(\tilde{\mathbf{u}}, \boldsymbol{\lambda}, \boldsymbol{\nu}) = f(\tilde{\mathbf{u}}) + \sum_{i=1}^r \lambda_i g_i(\tilde{\mathbf{u}}) + \sum_{j=1}^s \nu_j h_j(\tilde{\mathbf{u}}) \leq f(\tilde{\mathbf{u}})$$

Hence

$$L_D(\boldsymbol{\lambda}, \boldsymbol{\nu}) = \min_{\mathbf{u}} L(\mathbf{u}, \boldsymbol{\lambda}, \boldsymbol{\nu}) \leq f(\tilde{\mathbf{u}}).$$

This is true for any feasible  $\tilde{\mathbf{u}}$ , so

$$L_D(\boldsymbol{\lambda}, \boldsymbol{\nu}) \leq \min_{\tilde{\mathbf{u}} \text{ feasible}} f(\tilde{\mathbf{u}}) = p^*.$$

Taking the max order  $\boldsymbol{\lambda}, \boldsymbol{\nu} : \lambda_i \geq 0$ , we have

$$d^* = \max_{\boldsymbol{\lambda}, \boldsymbol{\nu} : \lambda_i \geq 0} L_D(\boldsymbol{\lambda}, \boldsymbol{\nu}) \leq p^*.$$

□

The difference  $p^* - d^*$  is called the *duality gap*.

## 2.2 Strong Duality

If  $p^* = d^*$ , we say *strong duality* holds. The original constrained optimization problem is said to be *convex* if  $f$  and  $g_1, \dots, g_r$  are convex functions and  $h_1, \dots, h_s$  are affine. We state the following without proof.

**Theorem 1.** *If the original problem is convex and a constraint qualification holds, then  $p^* = d^*$ .*

Two examples of constraint qualifications are

- All  $g_i$  are affine
- (strict feasibility)  $\exists \mathbf{u}$  s.t.  $h_j(\mathbf{u}) = 0 \ \forall j$  and  $g_i(\mathbf{u}) < 0 \ \forall i$ .

## 3 KKT Conditions

In this section, assume  $f, g_1, \dots, g_r, h_1, \dots, h_s$  are differentiable. For unconstrained optimization, we know  $\nabla f(\mathbf{u}^*) = 0$  is necessary for  $\mathbf{u}^*$  to be a global minimizer, and sufficient if  $f$  is additionally convex. The following two results generalize these properties to constrained optimization.

**Theorem 2.** *(Necessity) If  $p^* = d^*$ ,  $\mathbf{u}^*$  is primal optimal, and  $(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$  is dual optimal, then the Karush-Kuhn-Tucker (KKT) conditions hold:*

- (1)  $\nabla_{\mathbf{u}} f(\mathbf{u}^*) + \sum_{i=1}^r \lambda_i^* \nabla_{\mathbf{u}} g_i(\mathbf{u}^*) + \sum_{j=1}^s \nu_j^* \nabla_{\mathbf{u}} h_j(\mathbf{u}^*) = 0$
- (2)  $g_i(\mathbf{u}^*) \leq 0 \ \forall i$
- (3)  $h_j(\mathbf{u}^*) = 0 \ \forall j$
- (4)  $\lambda_i^* \geq 0 \ \forall i$
- (5)  $\lambda_i^* g_i(\mathbf{u}^*) = 0 \ \forall i$  (complimentary slackness).

*Proof.* (2) - (3) hold since  $\mathbf{u}^*$  is feasible. (4) holds by definition of the dual problem. To prove (5) and (1):

$$\begin{aligned}
f(\mathbf{u}^*) &= L_D(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*) \quad [\text{by strong duality}] \\
&= \min_{\tilde{\mathbf{u}}} \left( f(\tilde{\mathbf{u}}) + \sum_{i=1}^r \lambda_i^* g_i(\tilde{\mathbf{u}}) + \sum_{j=1}^s \nu_j^* h_j(\tilde{\mathbf{u}}) \right) \\
&\leq f(\mathbf{u}^*) + \sum_{i=1}^r \lambda_i^* g_i(\mathbf{u}^*) + \sum_{j=1}^s \nu_j^* h_j(\mathbf{u}^*) \\
&\leq f(\mathbf{u}^*) \quad [\text{by (2) - (4)}]
\end{aligned}$$

and therefore the two inequalities are equalities. Equality of the last two lines implies  $\lambda_i^* g_i(\mathbf{u}^*) = 0 \forall i$ . Equality of the 2nd and 3rd lines implies  $\mathbf{u}^*$  is a minimizer of  $L(\mathbf{u}, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$  with respect to  $\mathbf{u}$ . Therefore,

$$\nabla_{\mathbf{u}} L(\mathbf{u}^*, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*) = \mathbf{0},$$

which is (1).  $\square$

If  $(\tilde{\mathbf{u}}, \tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\nu}})$  satisfies the KKT conditions, we say that  $(\tilde{\mathbf{u}}, \tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\nu}})$  is a *feasible KKT point*. The preceding result shows that if  $\mathbf{u}^*$  is primal optimal, and  $(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$  is dual optimal, then  $(\mathbf{u}^*, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$  is a feasible KKT point.

**Theorem 3.** (*Sufficiency*) *If the original problem is convex and  $\tilde{\mathbf{u}}, \tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\nu}}$  satisfy the KKT conditions, then  $\tilde{\mathbf{u}}$  is primal optimal,  $(\tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\nu}})$  is dual optimal, and strong duality holds.*

*Proof.* By (2) and (3),  $\tilde{\mathbf{u}}$  is feasible. By (4),  $L(\mathbf{u}, \tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\nu}})$  is convex in  $\mathbf{u}$ . By (1),  $\tilde{\mathbf{u}}$  is a minimizer of  $L(\mathbf{u}, \tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\nu}})$ . Then

$$\begin{aligned}
L_D(\tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\nu}}) &= L(\tilde{\mathbf{u}}, \tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\nu}}) \\
&= f(\tilde{\mathbf{u}}) + \sum_{i=1}^r \tilde{\lambda}_i g_i(\tilde{\mathbf{u}}) + \sum_{j=1}^s \tilde{\nu}_j h_j(\tilde{\mathbf{u}}) \\
&= f(\tilde{\mathbf{u}}). \quad [\text{by (5) and (3)}]
\end{aligned}$$

Therefore  $p^* = d^*$  and the result follows.  $\square$

In conclusion, if a constrained optimization problem is differentiable, convex, and a constraint qualification holds, then the KKT conditions are necessary and sufficient for primal/dual optimality (with zero duality gap). The KKT conditions can be used to simplify and/or solve such problems.

## 4 Kernel Ridge Regression Revisited

When we studied kernel ridge regression (KRR), we used the matrix inversion lemma to kernelize ridge regression (RR). You may have wondered why this is valid when the feature space associated to the kernel is infinite-dimensional, as in the case of the Gaussian kernel. We will now use the above theory to rederive KRR while avoiding this issue.

Let  $(\mathbf{x}_i, y_i)$ ,  $i = 1, \dots, n$  be training data for regression. For simplicity, we let's focus on KRR w/o offset (the w/ offset case is left as an exercise). We may view ridge regression as the solution of

$$\begin{aligned}
\min_{\mathbf{w}, \xi} \quad & \lambda \|\mathbf{w}\|^2 + \sum_{i=1}^n \xi_i^2 \\
\text{s.t.} \quad & y_i - \mathbf{w}^T \mathbf{x}_i = \xi_i, \quad i = 1, \dots, n.
\end{aligned}$$

Denote the dual variables by  $\alpha_i \in \mathbb{R}$ . Notice that the above problem is convex and differentiable. Therefore, we can apply the KKT sufficiency theorem. Our goal is thus to find  $(\mathbf{w}^*, \boldsymbol{\xi}^*, \boldsymbol{\alpha}^*)$  satisfying the KKT conditions, and by the theorem, we will know that  $(\mathbf{w}^*, \boldsymbol{\xi}^*)$  solves the original/primal problem, and that  $\boldsymbol{\alpha}^*$  solves the dual. Since this problem has no inequality constraints, we are only concerned with KKT conditions (1) and (3).

The Lagrangian is

$$L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}) = \lambda \|\mathbf{w}\|^2 + \sum_{i=1}^n \xi_i^2 + \sum_{i=1}^n \alpha_i (y_i - \mathbf{w}^T \mathbf{x}_i - \xi_i).$$

According to the first KKT condition, we require

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{w}} &= 2\lambda \mathbf{w} - \sum_i \alpha_i \mathbf{x}_i = \mathbf{0} \\ \frac{\partial L}{\partial \xi_i} &= 2\xi_i - \alpha_i = 0 \quad \forall i. \end{aligned}$$

This leads to

$$\mathbf{w} = \frac{1}{2\lambda} \sum_i \alpha_i \mathbf{x}_i \tag{2}$$

$$\xi_i = \frac{1}{2} \alpha_i \quad \forall i. \tag{3}$$

According to the third KKT condition, we need

$$y_i - \mathbf{w}^T \mathbf{x}_i = \xi_i \quad \forall i.$$

Plugging (2) and (3) into this equation, and reformulating as a vector equation, we have

$$\mathbf{y} = \frac{1}{2\lambda} (\mathbf{K} + \lambda \mathbf{I}) \boldsymbol{\alpha}$$

where  $\mathbf{K}$  is the gram matrix and  $\mathbf{y} = [y_1 \cdots y_n]^T$ . Therefore, the KKT conditions are satisfied by

$$\begin{aligned} \boldsymbol{\alpha}^* &= 2\lambda (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y} \\ \mathbf{w}^* &= \frac{1}{2\lambda} \sum_i \alpha_i^* \mathbf{x}_i \\ \xi_i^* &= \frac{1}{2} \alpha_i^*. \end{aligned}$$

We may express

$$\begin{aligned} \mathbf{w}^* &= \frac{1}{2\lambda} \mathbf{X}^T \boldsymbol{\alpha}^* \\ &= \mathbf{X}^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y} \end{aligned}$$

where

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}.$$

The KRR (w/o offset) predictor is then

$$\begin{aligned} \hat{f}(\mathbf{x}) &= (\mathbf{w}^*)^T \mathbf{x} \\ &= \mathbf{y}^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{k}(\mathbf{x}), \end{aligned}$$

where

$$\mathbf{k}(\mathbf{x}) = \begin{bmatrix} \langle \mathbf{x}_1, \mathbf{x} \rangle \\ \vdots \\ \langle \mathbf{x}_n, \mathbf{x} \rangle \end{bmatrix}.$$

This agrees with our earlier derivation but did not require the matrix inversion lemma.

## Exercises

1. (★★) Prove the following property. If  $\mathbf{u}^*$  is primal optimal,  $(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$  is dual optimal, and strong duality holds, then  $(\mathbf{u}^*, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$  is a saddle point of  $L$ , i.e.,

$$L(\mathbf{u}^*, \boldsymbol{\lambda}, \boldsymbol{\nu}) \leq L(\mathbf{u}^*, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*) \leq L(\mathbf{u}, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$$

for all  $\mathbf{u} \in \mathbb{R}^p$ ,  $\boldsymbol{\lambda} \in \mathbb{R}^m$  with  $\lambda_i \geq 0$ , and  $\boldsymbol{\nu} \in \mathbb{R}^n$ . See Figure 2.

2. (★★) Repeat the KRR example for KRR with offset.
3. (☆☆) Use the definition of concavity to prove that for any constrained optimization problem, the dual function is concave.

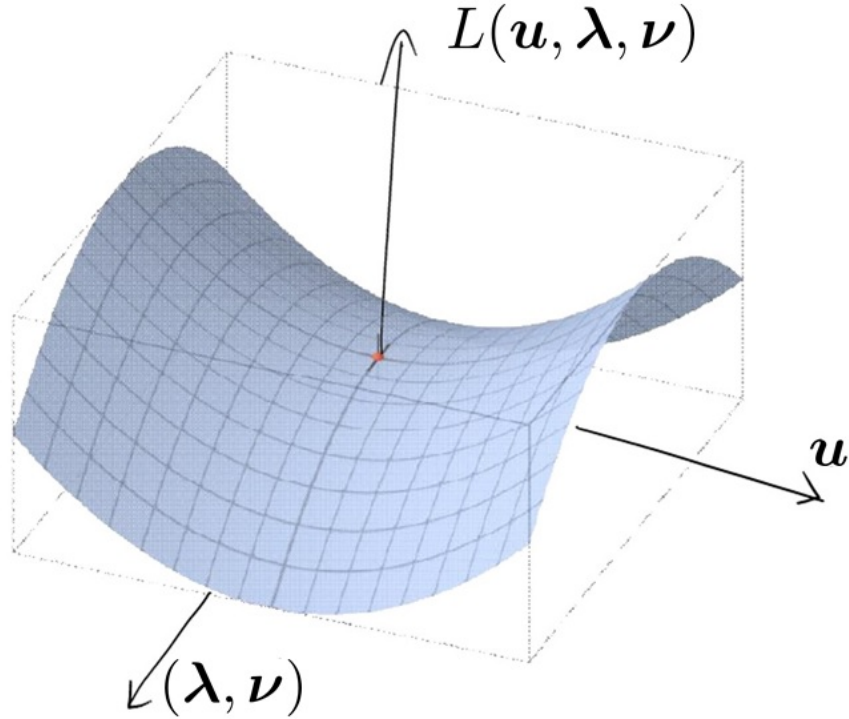


Figure 2: Illustration of a saddle point.