# **Proximal Methods**

October 9, 2024

#### **1** Moreau decomposition

In this section, we will explore some applications of duality in settings related to proximal gradient methods. First, recall the definition of a proximal operator:

$$\operatorname{prox}_f(v) = \arg\min_x \left(\frac{1}{2} \|x - v\|_2^2 + f(x)\right).$$

A useful fact for manipulating and extending proximal operators is known as **Moreau decomposition**. It states that the following relationship always holds:

$$v = \operatorname{prox}_f(v) + \operatorname{prox}_{f^*}(v),$$

where

$$f^*(y) = \max_x \left( y^\top x - f(x) \right).$$

Moreau's decomposition is "the main relationship between proximal operators and duality" and follows from the properties of sub-gradients and conjugate functions.

Notice that this is a generalization of orthogonal decomposition. Let L be a subspace

of a vector space U. For any  $v \in U$ , we have

$$v = \Pi_L(v) + \Pi_{L^\perp}(v).$$

To illustrate the usefulness of this decomposition, we review a simple example. If f(x) = ||x||, then  $f^*(y) = I_B(y)$ , where  $B = \{z : ||z||_* \le 1\}$  is a unit ball according to the dual norm. By Moreau decomposition,

$$v = \operatorname{prox}_{f}(v) + \operatorname{prox}_{f^{*}}(v)$$
$$= \operatorname{prox}_{\parallel \cdot \parallel}(v) + \operatorname{prox}_{I_{B}}(v),$$

where

$$prox_{I_B}(v) = \arg\min_{x} \left(\frac{1}{2} \|x - v\|_2^2 + I_B(x)\right)$$
  
=  $\arg\min_{x} \frac{1}{2} \|x - v\|_2^2$  s.t.  $x \in B$   
=  $\Pi_B(v)$ .

It follows that

$$\operatorname{prox}_{\|\cdot\|}(v) = v - \operatorname{prox}_{I_B}(v) = v - \Pi_B(v).$$

### 2 Extending the Moreau Decomposition

Starting from the identity

$$\operatorname{prox}_f(v) = v - \operatorname{prox}_{f^*}(v).$$

we want to derive a similar identity when we replace f by  $\lambda f$  for some  $\lambda > 0$ . We want to show that

$$\operatorname{prox}_{\lambda f}(v) = v - \operatorname{prox}_{(\lambda f)^*}(v) = v - \lambda \operatorname{prox}_{f^*/\lambda}(v/\lambda)$$

First, we find the convex conjugate of  $\lambda f$ :

$$(\lambda f)^*(v) = \max_y \left( v^\top y - \lambda f(y) \right)$$
$$= \max_y \lambda \left( \frac{v}{\lambda} y^\top - f(y) \right)$$
$$= \lambda \max_y \left( \frac{v}{\lambda} y^\top - f(y) \right)$$
$$= \lambda f^* \left( \frac{v}{\lambda} \right).$$

Then, we get

$$\operatorname{prox}_{(\lambda f)^*}(v) = \arg\min_{y} \left[ (\lambda f)^*(y) + \frac{1}{2} \|y - v\|_2^2 \right]$$
$$= \arg\min_{y} \left[ \lambda f^*\left(\frac{y}{\lambda}\right) + \frac{1}{2} \|y - v\|_2^2 \right]$$
$$= \arg\min_{y} \left[ f^*\left(\frac{y}{\lambda}\right) + \frac{1}{2\lambda} \|y - v\|_2^2 \right].$$

Now, we write  $y = \lambda z$  to get

$$\operatorname{prox}_{(\lambda f)^{*}}(v) = \arg\min_{\lambda z} \left[ f^{*}(z) + \frac{1}{2\lambda} \|\lambda z - v\|_{2}^{2} \right]$$
$$= \lambda \arg\min_{z} \left[ f^{*}(z) + \frac{\lambda}{2} \|z - \frac{v}{\lambda}\|_{2}^{2} \right]$$
$$= \lambda \operatorname{prox}_{f^{*}/\lambda} \left( \frac{v}{\lambda} \right).$$

Finally, we have the identity

$$\operatorname{prox}_{\lambda f}(v) = v - \operatorname{prox}_{(\lambda f)^*}(v) = v - \lambda \operatorname{prox}_{f^*/\lambda}(v/\lambda).$$

If  $f = \| \cdot \|$  is a general norm on  $\mathbb{R}^n$ , then

$$f^*(v) = I_B(v) = \begin{cases} 0 & \text{if } \|v\|_* \le 1, \\ \infty & \text{otherwise.} \end{cases}$$

where  $B = \{x : \|x\|_* \le 1\}$  is the unit-ball in  $(\mathbb{R}^n, \|\cdot\|_*)$ . Observe that

$$f^*/\lambda = I_B/\lambda = I_B.$$

Then by Moreau decomposition, we get:

$$\operatorname{prox}_{\lambda \parallel \cdot \parallel}(v) = v - \lambda \Pi_B\left(\frac{v}{\lambda}\right).$$

### **3** From Proximal to Projection

**Euclidean norm penalty.** Here,  $f = f^* = \| \cdot \|_2$ . We project v onto the Euclidean unit ball B as follows:

$$\Pi_B(v) = \begin{cases} v/\|v\|_2 & \text{if } \|v\|_2 > 1\\ 0 & \text{if } \|v\|_2 \le 1. \end{cases}$$

We get:

$$\operatorname{prox}_{\lambda \|\cdot\|_{2}}(v) = v - \lambda \Pi_{B} \left(\frac{v}{\lambda}\right)$$
$$= \begin{cases} (1 - \lambda/\|v\|_{2}) v & \text{if } \|v\|_{2} \ge \lambda \\ 0 & \text{if } \|v\|_{2} < \lambda \end{cases}$$
$$= (1 - \lambda/\|v\|_{2})_{+} v,$$

where

$$(z)_{+} = \begin{cases} z & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \,. \end{cases}$$

**Group lasso penalty.** This is how you compute proximal for each group in **group lasso**. For  $x \in \mathbb{R}^p$ ,

$$f(x) = \sum_{g=1}^{G} w_g ||x_g||_2$$

where  $\{1,...,p\}$  is partitioned into G groups. We get

$$\operatorname{prox}_{\lambda f}(v) = \arg \min_{x} \frac{1}{2t} \|v - x\|_{2}^{2} + \lambda f(x)$$
$$= \arg \min_{x} \frac{1}{2t} \|v - x\|_{2}^{2} + \lambda \sum_{g=1}^{G} w_{g} \|x_{g}\|_{2}.$$

So, for  $g \in \{1, ..., G\}$ ,

$$[\operatorname{prox}_{\lambda f}(v)]_g = \arg\min_{x_g} \frac{1}{2t} \|v_g - x_g\|_2^2 + \lambda w_g \|x_g\|_2$$
$$= \operatorname{prox}_{\lambda w_g \|x_g\|_2}(v_g)$$
$$= \left(1 - \frac{t\lambda w_g}{\|v_g\|_2}\right)_+ v_g.$$

**Sparse group lasso penalty.** For  $x \in \mathbb{R}^p$ , the sparse group lasso penalty is

$$f(x) = \sum_{g=1}^{G} w_g \left[ (1 - \alpha) \| x_g \|_2 + \alpha \| x_g \|_1 \right]$$

where  $\{1, ..., p\}$  is partitioned into G groups. This is how you compute proximal for each group in sparse **group lasso**. We get

$$\operatorname{prox}_{\lambda f}(v) = \arg\min_{x} \frac{1}{2t} \|v - x\|_{2}^{2} + \lambda f(x)$$

So, for  $g \in \{1, ..., G\}$ , define  $\tau = t\lambda w_g$ 

$$[\operatorname{prox}_{\lambda f}(v)]_{g} = \arg\min_{x_{g}} \frac{1}{2t} \|v_{g} - x_{g}\|_{2}^{2} + \lambda w_{g} \left[ (1 - \alpha) \|x_{g}\|_{2} + \alpha \|x_{g}\|_{1} \right]$$
  
$$= \operatorname{prox}_{\lambda f}(v_{g})$$
  
$$= \left( 1 - \frac{(1 - \alpha)\tau}{\|S_{\alpha \tau}(v_{g})\|_{2}} \right)_{+} S_{\alpha \tau}(v_{g}),$$

where

$$S_{\alpha\tau}(v_g) = \operatorname{sgn}(v_g)(|v_g| - \alpha\tau)_+.$$

 $l^1$  and  $l^\infty$  norms penalty. When  $f = \|\cdot\|_1$ , then  $f^* = I_B$ ,  $B = \{x : \|x\|_\infty \le 1\}$ . We project onto the  $\infty$ -norm unit ball B as follows:

$$(\Pi_B(v))_i = \begin{cases} 1 & : v_i > 1\\ v_1 & : |v_i| \le 1\\ -1 & : v_i < -1. \end{cases}$$

We get an alternative way of getting the proximal operator of lasso

$$\operatorname{prox}_{\lambda f}(v) = \operatorname{prox}_{\lambda \parallel \cdot \parallel_1}(v) = v - \lambda \prod_B \left(\frac{v}{\lambda}\right).$$

So

$$\left[\operatorname{prox}_{\lambda f}(v)\right]_{i} = \begin{cases} v_{i} - \lambda & : v_{i} > \lambda \\ 0 & : |v_{i}| \leq \lambda \\ v_{i} + \lambda & : v_{i} < \lambda. \end{cases}$$

When  $f = \| \cdot \|_{\infty}$ , then  $f^* = I_B$ ,  $B = \{x : \|x\|_1 \le 1\}$ . See paper for how to project on B.

**Hierarchical grouped norms.** Assume the variables  $X_1, ..., X_p$  have a hierarchical structure. The variables are selected according to the following rule, for  $i \in \{1, ..., p\}$ :

if 
$$\beta_i \neq 0$$
, then  $\beta_j \neq 0$  for all  $\beta_j \in \text{ancestors}(\beta_i)$ .

We define the following penalty:

$$\Omega(\beta) = \sum_{g \in G} w_g \| (\beta_g, \operatorname{descendents}(\beta_g)) \|_2,$$

where G is the set of all nodes. The proximal operator for this penalty is:

$$\operatorname{prox}_{\lambda\Omega}(v) = \arg\min_{u\in\mathbb{R}^p} \frac{1}{2} \|v-u\|_2^2 + \lambda\Omega(u)$$

Dual of the proximal problem. Let  $v \in \mathbb{R}^p$ . Consider

$$\max_{\xi \in \mathbb{R}^{p \times |G|}} -\frac{1}{2} \left( \| (v - \sum_{g \in G} \xi^g) \|_2^2 - \| v \|_2^2 \right)$$

such that for all  $g \in G$ ,  $\|\xi^g\|_* \le \lambda w_g$  and  $\xi_j^g = 0$  if  $j \notin g$ .

## **4** Applications

#### 4.1 Multitask sparse learning

**Data:** K data sources  $\{\mathbf{y}^{(k)}, \mathbf{X}^{(k)}\}_{k=1}^{K}$ : k-th data has  $n_k$  observations

- **Response:**  $\mathbf{y}^{(k)} = (y_1^{(k)}, \dots, y_{n_k}^{(k)})^{\mathsf{T}}$
- Predictors:  $\mathbf{X}^{(k)} = (\mathbf{x}_1^{(k)}, \dots, \mathbf{x}_{n_k}^{(k)})^\intercal$

- 
$$\mathbf{x}_{i}^{(k)} = (x_{i1}^{(k)}, \dots, x_{ip}^{(k)})^{\mathsf{T}} \in \mathbb{R}^{p}$$

Model: For continuous data, assume

$$E(y_i^{(k)}|\mathbf{x}_i^{(k)}) = \mathbf{x}_i^{(k)\top} \boldsymbol{\beta}^{(k)},$$

and for binary outcome use logistic regression setting

logit 
$$\left[P\left(y_i^{(k)}=1|x_i^{(k)}\right)\right] = \mathbf{x}_i^{(k)\top}\boldsymbol{\beta}^{(k)},$$

where  $y_i^{(k)} = \{-1, 1\}, i = 1, \dots, n_k$ . Here

$$\boldsymbol{\beta}^{(k)} = (\beta_1^{(k)}, \cdots, \beta_p^{(k)})^\top$$

is the coefficient vector for task k. Here  $\beta_j^{(k)}$  is the *j*-th element of  $\beta^{(k)}$ , for j = 1, ..., p. And the vector

$$\boldsymbol{\beta}_j = (\beta_j^{(1)}, \cdots, \beta_j^{(K)})^{\top},$$

contains the *j*-th elements of task 1 to task K. The whole coefficient can be written as a  $p \times K$  matrix

$$\boldsymbol{\beta} = (\boldsymbol{\beta}_1^{\top}, \cdots, \boldsymbol{\beta}_p^{\top})^{\top} \in \mathbb{R}^{p \times K}$$

To estimate  $\beta$ , for continuous outcome we minimize an aggregated least squares loss function

$$\ell(\boldsymbol{\beta}) = \frac{1}{n} \sum_{k=1}^{K} \left[ \mathbf{y}^{(k)} - \mathbf{X}^{(k)} \boldsymbol{\beta}^{(k)} \right]^{\top} \left[ \mathbf{y}^{(k)} - \mathbf{X}^{(k)} \boldsymbol{\beta}^{(k)} \right],$$
(1)

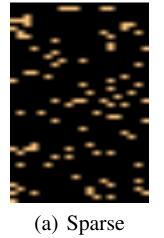
for binary outcome, we use

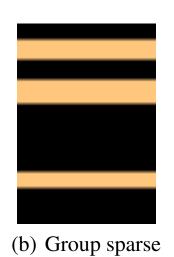
$$\ell(\boldsymbol{\beta}) = \frac{1}{n} \sum_{k=1}^{K} \sum_{i=1}^{n_k} \left[ 1 + \exp\left(-y_i^{(k)} \mathbf{x}_i^{(k)\top} \boldsymbol{\beta}^{(k)}\right) \right]$$
(2)

To consider structural sparsity

- Common relevant covariates across data sources
- Source-specific relevant covariates

$$\begin{bmatrix} \boldsymbol{\beta}_{1}^{\mathsf{T}} \\ \boldsymbol{\beta}_{2}^{\mathsf{T}} \\ \vdots \\ \boldsymbol{\beta}_{p}^{\mathsf{T}} \end{bmatrix} = \begin{bmatrix} \boldsymbol{\beta}_{1}^{(1)} & \boldsymbol{\beta}_{1}^{(2)} & \cdots & \boldsymbol{\beta}_{1}^{(K)} \\ \beta_{2}^{(1)} & \beta_{2}^{(2)} & \cdots & \boldsymbol{\beta}_{1}^{(K)} \\ \beta_{2}^{(1)} & \beta_{2}^{(2)} & \cdots & \boldsymbol{\beta}_{1}^{(K)} \\ \vdots & \vdots & \vdots \\ \boldsymbol{\beta}_{p}^{(1)} & \boldsymbol{\beta}_{p}^{(2)} & \cdots & \boldsymbol{\beta}_{p}^{(K)} \end{bmatrix}$$







(c) Group sparse *plus* sparse

We use composite  $L_1/L_2$  penalty

$$P_{\alpha,\lambda}(\boldsymbol{\beta}) = \lambda \sum_{j=1}^{p} v_j \left[ (1-\alpha) ||\boldsymbol{\beta}_j||_2 + \alpha ||\boldsymbol{\beta}_j||_1 \right]$$

Lasso: when  $\alpha = 1$ 

- Group lasso: when  $\alpha = 0$
- Sparse group lasso: when  $0<\alpha<1$