# MATH 680 Computation Intensive Statistics

November 27, 2019

**Proximal Methods for Penalization**

# 1 Proximal methods

## 1.1 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:

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

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

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

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 $B = \{z \colon \|z\|_* \leq 1\}$ is a unit ball according to the dual norm. By Moreau decomposition,

$$v = \text{prox}_f(v) + \text{prox}_{f^*}(v)$$
$$= \text{prox}_{\|\cdot\|}(v) + \text{prox}_{I_B}(v),$$

where

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

It follows that

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

## 1.2 Extending the Moreau Decomposition

Starting from the identity

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

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

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

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

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

Then, we get

$$
\begin{aligned}
\mathrm{prox}_{(\lambda f)^*}(v) &= \arg\min_y \left[ (\lambda f)^*(y) + \frac{1}{2} \|y - v\|_2^2 \right] \\
&= \arg\min_y \left[ \lambda f^* \left( \frac{y}{\lambda} \right) + \frac{1}{2} \|y - v\|_2^2 \right] \\
&= \arg\min_y \left[ f^* \left( \frac{y}{\lambda} \right) + \frac{1}{2\lambda} \|y - v\|_2^2 \right].
\end{aligned}
$$

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

$$
\begin{aligned}
\mathrm{prox}_{(\lambda f)^*}(v) &= \arg\min_{\lambda z} \left[ f^*(z) + \frac{1}{2\lambda} \|\lambda z - v\|_2^2 \right] \\
&= \lambda \arg\min_z \left[ f^*(z) + \frac{\lambda}{2} \left\| z - \frac{v}{\lambda} \right\|_2^2 \right] \\
&= \lambda \mathrm{prox}_{f^*/\lambda} \left( \frac{v}{\lambda} \right).
\end{aligned}
$$

Finally, we have the identity

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

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

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

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

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

Then by Moreau decomposition, we get:

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

## 1.3   From Proximal to Projection

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

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

We get:

$$\begin{aligned} \text{prox}_{\lambda\|\cdot\|_2}(v) &= v - \lambda\Pi_B\left(\frac{v}{\lambda}\right) \\ &= \begin{cases} (1 - \lambda/\|v\|_2)\, v & \text{if } \|v\|_2 \geq \lambda \\ 0 & \text{if } \|v\|_2 < \lambda \end{cases} \\ &= (1 - \lambda/\|v\|_2)_+\, v, \end{aligned}$$

where

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

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

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

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

$$\text{prox}_{\lambda f}(v) = \arg\min_x \frac{1}{2}\|v - x\|_2^2 + \lambda f(x)$$

$$= \arg\min_x \frac{1}{2}\|v - x\|_2^2 + \lambda \sum_{g=1}^{G} \|x_g\|_2.$$

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

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

$$= \text{prox}_{\lambda\|x_g\|_2}(v_g)$$

$$= \left[\left(1 - \frac{\lambda}{\|v_g\|_2}\right)_+ v_g\right]_g.$$

$l^1$ **and** $l^\infty$ **norms.** 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

$$\text{prox}_{\lambda f}(v) = \text{prox}_{\lambda\|\cdot\|_1}(v) = v - \lambda \Pi_B\left(\frac{v}{\lambda}\right).$$

So

$$\left[\text{prox}_{\lambda f}(v)\right]_i = \begin{cases} v_i - \lambda & : v_i > \lambda \\ 0 & : |v_i| \le \lambda \\ v_i + \lambda & : v_i < \lambda. \end{cases}$$

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 $i \in \{1, ..., p\}$:

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

We define the following penalty:

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

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

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

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

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

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