

On Decomposing the Proximal Map

Authored by:

Yao-Liang Yu

Abstract

The proximal map is the key step in gradient-type algorithms, which have become prevalent in large-scale high-dimensional problems. For simple functions this proximal map is available in closed-form while for more complicated functions it can become highly nontrivial. Motivated by the need of combining regularizers to simultaneously induce different types of structures, this paper initiates a systematic investigation of when the proximal map of a sum of functions decomposes into the composition of the proximal maps of the individual summands. We not only unify a few known results scattered in the literature but also discover several new decompositions obtained almost effortlessly from our theory.

1 Paper Body

Regularization has become an indispensable part of modern machine learning algorithms. For example, the ‘2 -regularizer for kernel methods [1] and the ‘1 -regularizer for sparse methods [2] have led to immense successes in various fields. As real data become more and more complex, different types of regularizers, usually nonsmooth functions, have been designed. In many applications, it is thus desirable to combine regularizers, usually taking their sum, to promote different structures simultaneously. Since many interesting regularizers are nonsmooth, they are harder to optimize numerically, especially in large-scale high-dimensional settings. Thanks to recent advances [3?5], gradient-type algorithms have been generalized to take nonsmooth regularizers explicitly into account. And due to their cheap per-iteration cost (usually linear-time), these algorithms have become prevalent in many fields recently. The key step of such gradient-type algorithms is to compute the proximal map (of the nonsmooth regularizer), which is available in closed-form for some specific regularizers. However, the proximal map becomes highly nontrivial when we start to combine regularizers. The main goal of this paper is to systematically investigate when the proximal map of a sum of functions decomposes into the composition of the proximal maps of the individual functions, which we simply term prox-decomposition. Our motivation comes from a few known decomposition results scattered in the literature [6?8], all in the form of our interest. The

study of such proxdecompositions is not only of mathematical interest, but also the backbone of popular gradient-type algorithms [3?5]. More importantly, a precise understanding of this decomposition will shed light on how we should combine regularizers, taking computational efforts explicitly into account. After setting the context in Section 2, we motivate the decomposition rule with some justifications, as well as some cautionary results. Based on a sufficient condition presented in Section 3.1, we study how ?invariance? of the subdifferential of one function would lead to nontrivial proxdecompositions. Specifically, we prove in Section 3.3 that when the subdifferential of one function is scaling invariant, then the prox-decomposition always holds if and only if another function is radial?which is, quite unexpectedly, exactly the same condition proven recently for the validity of the representer theorem in the context of kernel methods [9, 10]. The generalization to cone invariance is considered in Section 3.4, and enables us to recover most known prox-decompositions, as well as some new ones falling out quite naturally. 1

Our notations are mostly standard. We use $\mathbb{I}_C(x)$ for the indicator function that takes 0 if $x \in C$ and ∞ otherwise, and $1_C(x)$ for the indicator that takes 1 if $x \in C$ and 0 otherwise. The symbol Id stands for the identity map and the extended real line $\mathbb{R} \cup \{\infty\}$ is denoted as $\bar{\mathbb{R}}$. paper we denote $\partial f(x)$ as the subdifferential of the function f at point x .

2

Preliminary

Let our domain be some (real) Hilbert space $(H, \langle \cdot, \cdot \rangle, \|\cdot\|)$, with the induced Hilbertian norm $\|\cdot\|$. If needed, we will assume some fixed orthonormal basis $\{e_i\}_{i \in I}$ is chosen for H , so that for $x \in H$ we are able to refer to its ?coordinates? $x_i = \langle x, e_i \rangle$. we define its Moreau envelop as [11] For any closed convex proper function $f : H \rightarrow \bar{\mathbb{R}}$, $\forall y \in H$, $M_f(y) = \min_{x \in H} \frac{1}{2} \|x - y\|^2 + f(x)$,

(1)

and the related proximal map $P_f(y) = \arg\min_{x \in H} \frac{1}{2} \|x - y\|^2 + f(x)$.

(2)

$x \in H$

Due to the strong convexity of $\frac{1}{2} \|\cdot\|^2$ and the closedness and convexity of f , $P_f(y)$ always exists and is unique. Note that $M_f : H \rightarrow \bar{\mathbb{R}}$ while $P_f : H \rightarrow H$. When $f = \mathbb{I}_C$ is the indicator of some closed convex set C , the proximal map reduces to the usual projection. Perhaps the most interesting property of M_f , known as Moreau's identity, is the following decomposition [11] $M_f(y) + M_f^*(y) = \frac{1}{2} \|y\|^2$,

(3)

?

where $f^*(z) = \sup_{x \in H} \langle z, x \rangle - f(x)$ is the Fenchel conjugate of f . It can be shown that M_f is Frechet differentiable, hence taking derivative w.r.t. y in both sides of (3) yields $P_f(y) + P_f^*(y) = y$.

3

(4)

Main Results

Our main goal is to investigate and understand the equality (we always assume $f + g \in \Gamma_0(H)$) ?

$$\begin{aligned} &? \\ &Pf + g = Pf \circ Pg = Pg \circ Pf, \\ &(5) \end{aligned}$$

where $f, g \in \Gamma_0(H)$, the set of all closed convex proper functions on H , and $f \circ g$ denotes the mapping composition. We present first some cautionary results. ?1 ?1 Note that $Pf = (Id + \frac{1}{2}f)^{-1}$, hence under minor technical assumptions $Pf + g = (P^{-1} \frac{1}{2} Id + \frac{1}{2} Pf + P^{-1} \frac{1}{2} g)^{-1}$. However, computationally this formula is of little use. On the other hand, it is possible to develop forward-backward splitting procedures¹ to numerically compute $Pf + g$, using only Pf and Pg as subroutines [12]. Our focus is on the exact closed-form formula (5). Interestingly, under some ?shrinkage? assumption, the prox-decomposition (5), even if not necessarily hold, can still be used in subgradient algorithms [13].

Our first result is encouraging: Proposition 1. If $H = \mathbb{R}$, then for any $f, g \in \Gamma_0(\mathbb{R})$, there exists $h \in \Gamma_0(\mathbb{R})$ such that $Ph = Pf \circ Pg$. Proof: In fact, Moreau [11, Corollary 10.c] proved that $P : H \rightarrow H$ is a proximal map iff it is nonexpansive and it is the subdifferential of some convex function in $\Gamma_0(H)$. Although the latter condition in general is not easy to verify, it reduces to monotonic increasing when $H = \mathbb{R}$ (note that P must be continuous). Since both Pf and Pg are increasing and nonexpansive, it follows easily that so is $Pf \circ Pg$, hence the existence of $h \in \Gamma_0(\mathbb{R})$ so that $Ph = Pf \circ Pg$. In a general Hilbert space H , we again easily conclude that the composition $Pf \circ Pg$ is always a nonexpansion, which means that it is ?close? to be a proximal map. This justifies the composition $Pf \circ Pg$ as a candidate for the decomposition of $Pf + g$. However, we note that Proposition 1 indeed can fail already in \mathbb{R}^2 : 1 In some sense, this procedure is to compute $Pf + g = \lim_{t \rightarrow \infty} (Pf \circ Pg)^t$, modulo some intermediate steps. Essentially, our goal is to establish the one-step convergence of that iterative procedure.

2

Example 1. Let $H = \mathbb{R}^2$. Let $f = \{x_1 = x_2\}$ and $g = \{x_2 = 0\}$. Clearly both f and g are in $\Gamma_0(\mathbb{R}^2)$. The $2x_1 + x_2$ proximal maps in this case are simply projections: $Pf(x) = (\frac{x_1 + x_2}{2}, \frac{x_1 + x_2}{2})$ and $Pg(x) = (x_1, 0)$. 2 Therefore $Pf(Pg(x)) = (\frac{x_1}{2}, \frac{x_1}{2})$. We easily verify that the inequality $\|Pf(Pg(x)) - Pf(Pg(y))\| \leq \|Pf(Pg(x)) - Pf(Pg(y))\|$, $x \neq y$ is not always true, contradiction if $Pf \circ Pg$ was a proximal map [11, Eq. (5.3)]. Even worse, when Proposition 1 does hold, in general we can not expect the decomposition (5) to be true without additional assumptions. 1 Example 2. Let $H = \mathbb{R}$ and $q(x) = \frac{1}{2}x^2$. It is easily seen that $Pq(x) = 1 + x$. Therefore $1 \circ Pq \circ Pq = 4Id \neq 3Id = Pq + q$. We will give an explanation for this failure of composition shortly. Nevertheless, as we will see, the equality in (5) does hold in many scenarios, and an interesting theory can be suitably developed. 3.1

A Sufficient Condition

We start with a sufficient condition that yields (5). This result, although easy to obtain, will play a key role in our subsequent development. Using the first order optimality condition and the definition of the proximal map (2), we have $Pf + g(y) \in y + \partial(f + g)(Pf + g(y)) \cap \partial Pg(y) \in y + \partial g(Pg(y)) \cap \partial Pf$

$$(Pg(y)) \in Pg(y) + \partial f(Pg(y)) \quad (6)$$

$$(7) \quad (8)$$

Adding the last two equations we obtain $Pf(Pg(y)) \in y + \partial g(Pg(y)) + \partial f(Pg(Pg(y))) \quad (9)$

Comparing (6) and (9) gives us Theorem 1. A sufficient condition for $Pf + g = Pf \circ Pg$ is $x \in H, \partial g(Pf(x)) \subseteq \partial g(x)$.

$$(10)$$

Proof: Let $x = Pg(y)$. Then by (9) and the subdifferential rule $\partial(f + g) \subseteq \partial f + \partial g$ we verify that $Pf(Pg(y))$ satisfies (6), hence follows $Pf + g = Pf \circ Pg$ since the proximal map is single-valued. We note that a special form of our sufficient condition has appeared in the proof of [8, Theorem 1], whose main result also follows immediately from our Theorem 4 below. Let us fix f , and define $Kf = \{g \in \mathcal{C} : f + g \in \mathcal{C}, (f, g) \text{ satisfy (10)}\}$. Immediately we have Proposition 2. For any $f \in \mathcal{C}$, Kf is a cone. Moreover, if $g_1 \in Kf, g_2 \in Kf$, $f + g_1 + g_2 \in \mathcal{C}$ and $\partial(g_1 + g_2) = \partial g_1 + \partial g_2$, then $g_1 + g_2 \in Kf$ too. The condition $\partial(g_1 + g_2) = \partial g_1 + \partial g_2$ in Proposition 2 is purely technical; it is satisfied when, say g_1 is continuous at a single, arbitrary point in $\text{dom } g_1 \cap \text{dom } g_2$. For comparison purpose, we note that it is not clear how $Pf + g + h = Pf \circ Pg + h = Pf \circ Ph$. This is the main motivation to consider the sufficient condition (10). In particular Definition 1. We call $f \in \mathcal{C}$ self-prox-decomposable (s.p.d.) if $f \in Kf$ for all $\lambda \geq 0$. For any s.p.d. f , since Kf is a cone, $\lambda f \in Kf$ for all $\lambda \geq 0$. Consequently, $P(\lambda + \mu)f = P\lambda f + P\mu f = P\lambda f \circ P\mu f$. Remark 1. A weaker definition for s.p.d. is to require $f \in Kf$, from which we conclude that $\lambda f \in Kf$ for all $\lambda \geq 0$, in particular $P(m+n)f = Pnf \circ Pmf = Pmf \circ Pnf$ for all natural numbers m and n . The two definitions coincide for positive homogeneous functions. We have not been able to construct a function that satisfies this weaker definition but not the stronger one in Definition 1. Example 3. We easily verify that all affine functions $f(x) = \langle a, x \rangle + b$ are s.p.d., in fact, they are the only differentiable functions that are s.p.d., which explains why Example 2 must fail. Another trivial class of s.p.d. functions are projectors to closed convex sets. Also, univariate gauges are s.p.d., due to Theorem 4 below. Some multivariate s.p.d. functions are given in Remark 5 below. 2

A gauge is a positively homogeneous convex function that vanishes at the origin.

3

The next example shows that (10) is not necessary. Example 4. Fix $z \in H$, $f = \delta_{\{z\}}$, and $g \in \mathcal{C}$ with full domain. Clearly for any $x \in H$, $Pf + g(x) = z + g(x)$. However, since x is arbitrary, $\partial g(Pf(x)) = \partial g(z) \subseteq \partial g(x)$ if g is not linear. On the other hand, if f, g are differentiable, then we actually have equality in (10), which is clearly necessary in this case. Since convex functions are almost everywhere differentiable (in the interior of their domain), we expect the sufficient condition (10) to be necessary ‘almost everywhere’ too. Thus we see that the key for the decomposition (5) to hold is to let the proximal map of f and the subdifferential of g ‘interact well’ in the sense of (10). Interestingly,

both are fully equivalent to the function itself. Proposition 3 ([11, ?8]). Let $f, g : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$. $f = g + c$ for some $c \in \mathbb{R}$ iff $Pf = Pg$. Proof: The first implication is clear. The second follows from the optimality condition $Pf = (\text{Id} + \gamma \partial f)^{-1}$. Lastly, $Pf = Pg$ implies that $Mf = Mg + c$ for some $c \in \mathbb{R}$ (by integration). Conjugating we get $f = g + c$ for some $c \in \mathbb{R}$. Therefore some properties of the proximal map will transfer to some properties of the function f itself, and vice versa. The next result is easy to obtain, and appeared essentially in [14]. Proposition 4. Let $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ and $x \in H$ be arbitrary, then i). Pf is odd iff f is even; ii). $Pf(Ux) = U Pf(x)$ for all unitary U iff $f(Ux) = f(x)$ for all unitary U ; iii). $Pf(Qx) = Q Pf(x)$ for all permutation Q (under some fixed basis) iff f is permutation invariant, that is $f(Qx) = f(x)$ for all permutation Q . In the following, we will put some invariance assumptions on the subdifferential of g and accordingly find the right family of f whose proximal map respects that invariance. This way we will meet (10) by construction therefore effortlessly have the decomposition (5). 3.2

No Invariance

To begin with, consider first the trivial case where no invariance on the subdifferential of g is assumed. This is equivalent as requiring (10) to hold for all $g : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$. Not surprisingly, we end up with a trivial choice of f . Theorem 2. Fix $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$. $Pf + g = Pf \circ Pg$ for all $g : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ if and only if $f = \dim(H) \cdot \mathbb{I} + c$ or $f = \mathbb{I}_{\{w\}} + c$ for some $c \in \mathbb{R}$ and $w \in H$; $\dim(H) = 1$ and $f = \mathbb{I}_C + c$ for some closed and convex set C and $c \in \mathbb{R}$. Proof: Straightforward calculations, see [15] for details. We first prove that f is constant on its domain even when g is restricted to indicators. Indeed, let $x \in \text{dom } f$ and take $g = \mathbb{I}_{\{x\}}$. Then $x = Pf + g(x) = Pf[Pg(x)] = Pf(x)$, meaning that $x \in \text{argmin } f$. Since $x \in \text{dom } f$ is arbitrary, f is constant on its domain. The case $\dim(H) = 1$ is complete. We consider the other case where $\dim(H) \geq 2$ and $\text{dom } f$ contains at least two points. If $\text{dom } f \neq H$, there exists $z \notin \text{dom } f$ such that $Pf(z) = y$ for some $y \in \text{dom } f$, and closed convex set $C \subset \text{dom } f$ with $y \in C$ and $z \notin C$. Let $g = \mathbb{I}_C$ we obtain $Pf + g(z) \in C \subset \text{dom } f$ while $Pf(Pg(z)) = Pf(z) = y \notin C$, contradiction. Observe that the decomposition (5) is not symmetric in f and g , also reflected in the next result: Theorem 3. Fix $g : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$. $Pf + g = Pf \circ Pg$ for all $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ iff g is a continuous affine function. Proof: If $g = h + a$, then $Pg(x) = x + a$. Easy calculation reveals that $Pf + g(x) = Pf(x + a) = Pf[Pg(x)]$. The converse is true even when f is restricted to continuous linear functions. Indeed, let $a \in H$ be arbitrary and consider $f = h + a$. Then $Pf + g(x) = Pg(x + a) = Pf(Pg(x)) = Pg(x) + a$. Letting $a = x$ yields $Pg(x) = x + Pg(0) = Ph + Pg(0)$. Therefore by Proposition 3 we know g is equal to a continuous affine function. 4

Naturally, the next step is to put invariance assumptions on the subdifferential of g , effectively restricting the function class of g . As a trade off, the function class of f , that satisfies (10), becomes larger so that nontrivial results will arise. 3.3

Scaling Invariance

The first invariance property we consider is scaling-invariance. What kind of convex functions have their subdifferential invariant to (positive) scaling?

Assuming $0 \in \text{dom } g$ and by simple integration $\int_0^t \int_0^t h(g(sx), x) ds = t \int_0^1 [g(x) - g(0)] g(sx) ds = g(tx) - g(0) = 0$

where the last equality follows from the scaling invariance of the subdifferential of g . Therefore, up to some additive constant, g is positive homogeneous (p.h.). On the other hand, if $g \in \mathcal{P}_0$ is p.h. (automatically $0 \in \text{dom } g$), then from definition we verify that g is scaling-invariant. Therefore, under the scaling-invariance assumption, the right function class for g is the set of all p.h. functions in \mathcal{P}_0 , up to some additive constant. Consequently, the right function class for f is to have the proximal map $Pf(x) = \inf_{y \in [0, 1]} x$ for some $\gamma \in [0, 1]$ that may depend on x as well³. The next theorem completely characterizes such functions. Theorem 4. Let $f \in \mathcal{P}_0$. Consider the statements i). $f = h(k \cdot k)$ for some increasing function $h : \mathbb{R}_+ \rightarrow \mathbb{R}$; ii). $x \leq y \Rightarrow f(x + y) \leq f(y)$; iii). $Pf(u) = \inf_{y \in [0, 1]} u$ for some $\gamma \in [0, 1]$ (that may itself depend on u); iv). $0 \in \text{dom } f$ and $Pf + \gamma = Pf \cdot P\gamma$ for all p.h. (up to some additive constant) function $\gamma \in \mathcal{P}_0$. Then we have i) \Leftrightarrow ii) \Leftrightarrow iii) \Leftrightarrow iv). Moreover, when $\dim(H) \geq 2$, ii) \Leftrightarrow i) as well, in which case $Pf(u) = Ph(ku)/ku \cdot u$ (where we interpret $0/0 = 0$). Remark 2. When $\dim(H) = 1$, ii) is equivalent as requiring f to attain its minimum at 0, in which case the implication ii) \Rightarrow iv), under the redundant condition that f is differentiable, was proved by Combettes and Pesquet [14, Proposition 3.6]. The implication ii) \Rightarrow iii) also generalizes [14, Corollary 2.5], where only the case $\dim(H) = 1$ and f differentiable is considered. Note that there exists non-even f that satisfies Theorem 4 when $\dim(H) = 1$. Such is impossible for $\dim(H) \geq 2$, in which case any f that satisfies Theorem 4 must also enjoy all properties listed in Proposition 4. Proof: i) \Rightarrow ii): $x \leq y \Rightarrow kx + yk \leq kyk$. ii) \Rightarrow iii): Indeed, by definition $Mf(u) = \min_{k \in \mathbb{R}_+} kx \cdot uk^2 + f(x) = \min_{k \in \mathbb{R}_+} k^2u + f(u) = \min_{k \in [0, 1]} k^2u + f(u) = \min_{k \in [0, 1]} k^2u + f(u) = \min_{k \in [0, 1]} k^2u + f(u) = \min_{k \in [0, 1]} k^2u + f(u)$

where the third equality is due to ii), and the nonnegative constraint in the last equality can be seen as follows: For any $\gamma \geq 0$, by increasing it to 0 we can only decrease both terms; similar argument for $\gamma \leq 1$. Therefore there exists $\gamma \in [0, 1]$ such that γu minimizes the Moreau envelop Mf hence we have $Pf(u) = \gamma u$ due to uniqueness. iii) \Rightarrow iv): Note first that iii) implies $0 \in \text{dom } f$, therefore $0 \in \text{dom } f$. Since the subdifferential of f is scaling-invariant, iii) implies the sufficient condition (10) hence iv). iv) \Rightarrow iii): Fix y and construct the gauge function

0 , if $z = \gamma y$ for some $\gamma \in [0, 1]$, otherwise $P\gamma(y) = y$, hence $Pf(P\gamma(y)) = Pf(y) = Pf + \gamma(y)$ by iv). On the other hand, $Mf + \gamma(y) = \min_{k \in \mathbb{R}_+} kx \cdot yk^2 + f(x) + \gamma(x) = \min_{k \in [0, 1]} k^2y \cdot yk^2 + f(y) = x$

3

Note that $\gamma \in [0, 1]$ is necessary since any proximal map is nonexpansive.

5

(11)

Take $y = 0$ we obtain $Pf + \gamma(0) = 0$. Thus $Pf(0) = 0$, i.e. $0 \in \text{dom } f$, from which we deduce that $Pf(y) = Pf + \gamma(y) = \gamma y$ for some $\gamma \in [0, 1]$, since $f(\gamma y)$ in (11) is increasing on $[1, \gamma]$. iii) \Rightarrow ii): First note that iii) implies that $Pf(0) = 0$

$= 0$ hence $0 \in \partial f(0)$, in particular, $0 \in \text{dom } f$. If $\dim(H) = 1$ we are done, so we assume $\dim(H) \geq 2$ in the rest of the proof. In this case, it is known, cf. [9, Theorem 1] or [10, Theorem 3], that ii) \Leftrightarrow i) (even without assuming f convex). All we left is to prove iii) \Leftrightarrow ii) or equivalently i), for the case $\dim(H) \geq 2$. We first prove the case when $\text{dom } f = H$. By iii), $Pf(x) = \gamma x$ for some $\gamma \in [0, 1]$ (which may depend on x as well). Using the first order optimality condition for the proximal map we have $0 \in \gamma x + \partial f(\gamma x)$, that is $(1 - \gamma)y \in \partial f(y)$ for each $y \in \text{ran}(Pf) = H$ due to our assumption $\text{dom } f = H$. Now for any $x \in y$, by the definition of the subdifferential,

$$f(x + y) \leq f(y) + \langle h, x \rangle, \quad \langle \gamma y, x \rangle \leq f(y) + \langle x, (1 - \gamma)y \rangle = f(y). \quad (12)$$

For the case when $\text{dom } f \neq H$, we consider the proximal average [16] $g = A(f, q) = [(1 - \gamma)f + \gamma q]^\gamma$,

where $q = \frac{1}{2} \| \cdot \|^2$. Importantly, since q is defined on the whole space, the proximal average g has full domain too [16, Corollary 4.7]. Moreover, $Pg(x) = (1 - \gamma)Pf(x) + \gamma x = (\frac{1 - \gamma}{2} + \frac{\gamma}{2})x$. Therefore by our previous argument, g satisfies ii) hence also i). It is easy to check that i) is preserved under taking the Fenchel conjugation (note that the convexity of f implies that of h). Since we have shown that g satisfies i), it follows from (12) that f satisfies i) hence also ii). As mentioned, when $\dim(H) \geq 2$, the implication ii) \Rightarrow i) was shown in [9, Theorem 1]. The formula $Pf(u) = Ph(ku)/ku$ for $f = h(\| \cdot \|)$ follows from straightforward calculation. We now discuss some applications of Theorem 4. When $\dim(H) \geq 2$, iii) in Theorem 4 automatically implies that the scalar constant γ depends on x only through its norm. This fact, although not entirely obvious, does have a clear geometric picture: Corollary 1. Let $\dim(H) \geq 2$, $C \subset H$ be a closed convex set that contains the origin. Then the projection onto C is simply a shrinkage towards the origin iff C is a ball (of the norm $\| \cdot \|$). Proof: Let $f = \text{dist}(\cdot, C)$ and apply Theorem 4. Example 5. As usual, denote $q = \frac{1}{2} \| \cdot \|^2$. In many applications, in addition to the regularizer γ (usually a gauge), one adds the ‘ ℓ_2 regularizer’ γq either for stability or grouping effect or strong convexity. This incurs no computational cost in the sense of computing the proximal map: We easily compute that $P(\gamma q) = \gamma \text{Id}$. By Theorem 4, for any gauge γ , $P(\gamma + \gamma q) = \gamma + \gamma P\gamma$, whence it is also clear that adding an extra ‘ ℓ_2 regularizer’ tends to double ‘shrink’ the solution. In particular, let $H = \mathbb{R}^d$ and take $\gamma = \| \cdot \|_1$ (the sum of absolute values) we recover the proximal map for the elastic-net regularizer [17]. Example 6. The Berhu regularizer $h(x) = \frac{1}{2} \| x \|_1 + \frac{1}{2} \| x \|_2^2$

$$= \frac{1}{2} \| x \|_1 + \frac{1}{2} \| x \|_2^2, \quad \text{being the reverse of Huber's function, is proposed in [18] as a bridge between}$$

the lasso (‘ ℓ_1 regularization’) and ridge regression (‘ ℓ_2 regularization’). Let $f(x) = h(x) = \frac{1}{2} \| x \|_1 + \frac{1}{2} \| x \|_2^2$. Clearly, f satisfies ii) of Theorem 4 (but not differentiable), hence $Ph = Pf = P\gamma$, whereas simple calculation verifies that $\gamma Pf(x) = \text{sign}(x) = \min\{-\|x\|_1, 1 + \frac{1}{2}(-\|x\|_1 + 1)\}$, and of course $P\gamma(x) = \text{sign}(x) = \max\{-\|x\|_1, 1, 0\}$. Note that this regularizer is not s.p.d. Corollary 2. Let

$\dim(H) \geq 2$, then the p.h. function $f \geq 0$ satisfies any item of Theorem 4 iff it is a positive multiple of the norm $k \geq k$. Proof: [10, Theorem 4] showed that under positive homogeneity, i) implies that f is a positive multiple of the norm. Therefore (positive multiples of) the Hilbertian norm is the only p.h. convex function f that satisfies $Pf + \gamma = Pf \circ P\gamma$ for all gauge γ . In particular, this means that the norm $k \geq k$ is s.p.d. Moreover, we easily recover the following result that is perhaps not so obvious at first glance: 6

Corollary 3 (Jenatton et al. [7]). Fix the orthonormal basis $\{e_i\}_{i \in I}$ of H . Let $G \geq 2I$ be a collection of tree-structured groups, that is, either $g \geq g_0$ or $g_0 \geq g$ or $g \geq g_0 = \gamma$ for all $g, g_0 \in G$. Then $PP_m = Pk \circ kg_1 \circ \gamma \circ \gamma \circ Pk \circ kg_m$, $i=1, \dots, k \circ kg_i$ where we arrange the groups so that $g_i \circ g_j = \gamma \circ i \circ j$, and the notation $k \circ kg_i$ denotes the Hilbertian norm that is restricted to the coordinates indexed by the group g_i . Proof: Let $f = k \circ kg_1$ and $\gamma = i=2, \dots, k \circ kg_i$. Clearly they are both p.h. (and convex). By the tree-structured assumption we can partition $\gamma = \gamma_1 + \gamma_2$, where $g_i \circ g_1$ for all g_i appearing in γ_1 while $g_j \circ g_1 = \gamma$ for all g_j appearing in γ_2 . Restricting to the subspace spanned by the variables in g_1 we can treat f as the Hilbertian norm. Apply Theorem 4 we obtain $Pf + \gamma_1 = Pf \circ P\gamma_1$. On the other hand, due to the non-overlapping property, nothing will be affected by adding γ_2 , thus $PP_m = Pk \circ kg_1 \circ P\gamma_2 \circ P\gamma_1 \circ Pk \circ kg_i$, $i=1, \dots, k \circ kg_i$. We can clearly iterate the argument to unravel the proximal map as claimed. For notational clarity, we have chosen not to incorporate weights in the sum of group seminorms: Such can be absorbed into the seminorm and the corollary clearly remains intact. Our proof also reveals the fundamental reason why Corollary 3 is true: The ℓ_2 norm admits the decomposition (5) for any gauge g ! This fact, to the best of our knowledge, has not been recognized previously. 3.4

Cone Invariance

In the previous subsection, we restricted the subdifferential of g to be constant along each ray. We now generalize this to cones. Specifically, consider the gauge function $\gamma(x) = \max_{j \in J} \langle a_j, x \rangle$, where J is a finite index set and each $a_j \in H$. Such polyhedral gauge functions have become extremely important due to the work of Chandrasekaran et al. [19]. Define the polyhedral cones $K_j = \{x \in H : \langle a_j, x \rangle = \gamma(x)\}$.

(13)

where J is a finite index set and each $a_j \in H$. Such polyhedral gauge functions have become extremely important due to the work of Chandrasekaran et al. [19]. Define the polyhedral cones $K_j = \{x \in H : \langle a_j, x \rangle = \gamma(x)\}$.

(14)

Assume $K_j \neq \emptyset$ for each j (otherwise delete j from J). Since $\gamma(x) = \max_{j \in J} \langle a_j, x \rangle$, the sufficient condition (10) becomes $\gamma_j \in J$, $Pf(K_j) \subseteq K_j$.

(15)

In other words, each cone K_j is γ -fixed under the proximal map of f . Although it would be very interesting to completely characterize f under (15), we show that in its current form, (15) already implies many known results, with some new generalizations falling out naturally. Corollary 4. Denote E a collection of pairs (m, n) , and define the total variational norm $\|x\|_{\text{TV}} = \sum_{(m,n) \in E} |w_{m,n} x_m - x_n|$, where $w_{m,n} \geq 0$. Then for any permutation invariant function f , $Pf + k\|\cdot\|_{\text{TV}} = Pf \circ Pk\|\cdot\|_{\text{TV}}$. Proof: Pick an arbitrary pair $(m, n) \in E$ and let $\gamma = |x_m - x_n|$. Clearly $J = \{1, 2\}$, $K_1 = \{x_m \geq x_n\}$ and $K_2 = \{x_m \leq x_n\}$. Since f is permutation invariant, its proximal map Pf

(x) maintains the order of x , hence we establish (15). Finally apply Proposition 2 and Theorem 1. Remark 3. The special case where $E = \{(1, 2), (2, 3), \dots\}$ is a chain, $w_{m,n} \geq 1$ and f is the ‘1 norm, appeared first in [6] and is generally known as the fused lasso. The case where f is the ‘p norm appeared in [20]. We call the permutation invariant function f symmetric if $f(x) = f(x)$, where $|x|$ denotes the componentwise absolute value. The proof for the next corollary is almost the same as that of Corollary 4, except that we also use the fact $\text{sign}([Pf(x)]_m) = \text{sign}(x_m)$ for symmetric functions. P Corollary 5. As in Corollary 4, define the norm $\|x\|_{k \circ k \circ k} = \{m, n\} \in E, w_{m,n} \max\{|x_m|, |x_n|\}$. Then for any symmetric function f , $Pf + k \circ k \circ k = Pf + Pk \circ k \circ k$. 4

All we need is the weaker condition: For all $\{m, n\} \in E$, $x_m \geq x_n \Rightarrow [Pf(x)]_m \geq [Pf(x)]_n$.

7

Remark 4. This norm $\|x\|_{k \circ k \circ k}$ is proposed in [21] for feature grouping. Surprisingly, Corollary 5 appears to be new. The proximal map $Pk \circ k \circ k$ is derived P in [22], which turns out to be another decomposition result. Indeed, for $i \geq 2$, define $\gamma_i(x) = \max\{|x_i|, |x_{i-1}|\}$. Thus $\|x\|_{k \circ k \circ k} = \max_i \gamma_i(x)$.

Importantly, we observe that γ_i is symmetric on the first $i-1$ coordinates. We claim that $Pk \circ k \circ k = P\gamma_1 + \dots + P\gamma_{i-1}$. The proof is by recursion: Write $\|x\|_{k \circ k \circ k} = f + g$, where $f = \gamma_1$. Note that the subdifferential of g depends only on the ordering and sign of the first $i-1$ coordinates while the proximal map of f preserves the ordering and sign of the first $i-1$ coordinates (due to symmetry). If we pre-sort x , the individual proximal maps $P\gamma_i(x)$ become easy to compute sequentially and we recover the algorithm in [22] with some bookkeeping. Corollary 6. As in Corollary 3, let $G \subseteq 2I$ be a collection of tree-structured groups, then $PP_m = Pk \circ k \circ k_1, k \circ k \circ k_2, \dots, k \circ k \circ k_m$, $i=1, \dots, k \circ k \circ k_i, k \circ k \circ k$ where we arrange the groups so that $g_i \leq g_j \Rightarrow i \leq j$, and $\|x\|_{k \circ k \circ k_i, k \circ k \circ k_j} = \max\{|x_{g_i}|, |x_{g_j}|\}$ is the sum of the k (absolute-value) largest elements in the group g_i , i.e., Ky-Fan’s k -norm. Proof: Similar as in the proof of Corollary 3, we need only prove that $Pk \circ k \circ k_1, k \circ k \circ k_2, k \circ k \circ k = Pk \circ k \circ k_1, k \circ k \circ k_2, k \circ k \circ k$, where w.l.o.g. we assume g_1 contains all variables while $g_2 \subseteq g_1$. Therefore $\|x\|_{k \circ k \circ k_1, k \circ k \circ k}$ can be treated as symmetric and the rest follows the proof of Corollary 5. Note that the case $k \in \{1, |I|\}$ was proved in [7] and Corollary 6 can be seen as an interpolation. Interestingly, there is another interpolated result whose proof should be apparent now. Corollary 7. Corollary 6 remains true if we replace Ky-Fan’s k -norm with $\|x\|_{k \circ k \circ k, k} = \max\{|x_{i_1}|, \dots, |x_{i_k}|\}$. (16) $1 \leq i_1 \leq i_2 \leq \dots \leq i_k \leq |I|$

Therefore we can employ the norm $\|x\|_{k \circ k \circ k, 2}$ for feature grouping in a hierarchical manner. Clearly we can also combine Corollary 6 and Corollary 7. Corollary 8. For any symmetric f , $Pf + k \circ k \circ k, k = Pf + Pk \circ k \circ k, k$. Similarly for Ky-Fan’s k -norm. Remark 5. The above corollary implies that Ky-Fan’s k -norm and the norm $\|x\|_{k \circ k \circ k, k}$ defined in (16) are both s.p.d. (see Definition 1). The special case for the ‘p norm where $p \in \{1, 2, \dots\}$ was proved in [23, Proposition 11], with a substantially more complicated argument. As pointed out in [23], s.p.d. regularizers allow us to perform lazy updates in gradient-type

algorithms. We remark that we have not exhausted the possibility to have the decomposition (5). It is our hope to stimulate further work in understanding the prox-decomposition (5). Added after acceptance: We have managed to extend the results in this subsection to the Lov sz extension of submodular set functions. Details will be given elsewhere.

4

Conclusion

The main goal of this paper is to understand when the proximal map of the sum of functions decomposes into the composition of the proximal maps of the individual functions. Using a simple sufficient condition we are able to completely characterize the decomposition when certain scaling invariance is exhibited. The generalization to cone invariance is also considered and we recover many known decomposition results, with some new ones obtained almost effortlessly. In the future we plan to generalize some of the results here to nonconvex functions.

Acknowledgement The author thanks Bob Williamson and Xinhua Zhang from NICTA/Canberra for their hospitality during the author's visit when part of this work was performed; Warren Hare, Yves Lucet, and Heinz Bauschke from UBC/Okanagan for some discussions around Theorem 4; and the reviewers for their valuable comments. 8

2 References

- [1] Bernhard Scholkopf and Alexander J. Smola. Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond. MIT Press, 2001.
- [2] Peter B hlmann and Sara van de Geer. Statistics for High-Dimensional Data. Springer, 2011.
- [3] Patrick L. Combettes and Val rie R. Wajs. Signal recovery by proximal forward-backward splitting. Multiscale Modeling and Simulation, 4(4):1168–1200, 2005.
- [4] Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. SIAM Journal on Imaging Sciences, 2(1):183–202, 2009.
- [5] Yurii Nesterov. Gradient methods for minimizing composite functions. Mathematical Programming, Series B, 140:125–161, 2013.
- [6] Jerome Friedman, Trevor Hastie, Holger H fing, and Robert Tibshirani. Pathwise coordinate optimization. The Annals of Applied Statistics, 1(2):302–332, 2007.
- [7] Rodolphe Jenatton, Julien Mairal, Guillaume Obozinski, and Francis Bach. Proximal methods for hierarchical sparse coding. Journal of Machine Learning Research, 12:2297–2334, 2011.
- [8] Jiayu Zhou, Jun Liu, Vaibhav A. Narayan, and Jieping Ye. Modeling disease progression via fused sparse group lasso. In Conference on Knowledge Discovery and Data Mining, 2012.
- [9] Francesco Dinuzzo and Bernhard Sch lkopf. The representer theorem for Hilbert spaces: a necessary and sufficient condition. In NIPS, 2012.
- [10] Yao-Liang Yu, Hao Cheng, Dale Schuurmans, and Csaba Szepesv ri. Characterizing the representer theorem. In ICML, 2013.
- [11] Jean J. Moreau. Proximit  et dualit  dans un espace Hilbertien. Bulletin de la Soci t  Math matique de France, 93:273–299, 1965.
- [12] Patrick L. Combettes,

Vinh Dũng, and Ba'ng C'ng V'u. Proximity for sums of composite functions. *Journal of Mathematical Analysis and Applications*, 380(2):680?688, 2011. [13] André F. T. Martins, Noah A. Smith, Eric P. Xing, Pedro M. Q. Aguiar, and Mário A. T. Figueiredo. Online learning of structured predictors with multiple kernels. In *Conference on Artificial Intelligence and Statistics*, 2011. [14] Patrick L. Combettes and Jean-Christophe Pesquet. Proximal thresholding algorithm for minimization over orthonormal bases. *SIAM Journal on Optimization*, 18(4):1351?1376, 2007. [15] Yaoliang Yu. Fast Gradient Algorithms for Structured Sparsity. PhD thesis, University of Alberta, 2013. [16] Heinz H. Bauschke, Rafal Goebel, Yves Lucet, and Xianfu Wang. The proximal average: Basic theory. *SIAM Journal on Optimization*, 19(2):766?785, 2008. [17] Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society B*, 67:301?320, 2005. [18] Art B. Owen. A robust hybrid of lasso and ridge regression. In *Prediction and Discovery*, pages 59?72. AMS, 2007. [19] V. Chandrasekaran, B. Recht, P. A. Parrilo, and A. S. Willsky. The convex geometry of linear inverse problems. *Foundations of Computational Mathematics*, 12(6):805?849, 2012. [20] Xinhua Zhang, Yaoliang Yu, and Dale Schuurmans. Polar operators for structured sparse estimation. In *NIPS*, 2013. [21] Howard Bondell and Brian Reich. Simultaneous regression shrinkage, variable selection, and supervised clustering of predictors with oscar. *Biometrics*, 64(1):115?123, 2008. [22] Leon Wenliang Zhong and James T. Kwok. Efficient sparse modeling with automatic feature grouping. In *ICML*, 2011. [23] John Duchi and Yoram Singer. Efficient online and batch learning using forward backward splitting. *Journal of Machine Learning Research*, 10:2899?2934, 2009.