

## Notes on Learned Proximal Networks

- Let  $\underline{x} \in \mathbb{R}^n$ ,  $A: \mathbb{R}^m \rightarrow \mathbb{R}^n$ , and consider the inverse problem

$$\underline{y} = A(\underline{x}) + \underline{v}, \quad (1)$$

where  $\underline{v}$  is some noise/nuisance term. Our goal is to recover a solution  $\underline{\hat{x}} \in \mathbb{R}^n$  to this problem that approximates the real signal  $\underline{x}$  or even recover it.

- This is an **ill-posed inverse problem**. The inverse problem may have an infinite number of solutions all approximating the signal ... To overcome this, we need to regularize the problem.

- One approach: Use a **prior** to regularize the problem. This is a function  $R: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  that:

- Promote certain desirable properties in  $\underline{x}$  or
- Encapsulate the knowledge we have about  $\underline{x}$  or
- Promote a solution likely under the distribution of  $\underline{x}$ ,

or any combinations of i) - iii).



Using a prior function  $R$  does **NOT** mean necessarily that  $\underline{x}$  is sampled from a distribution  $p \sim e^{-RG(\underline{x})}$  or is likely a sample from  $e^{-RG(\underline{x})}$  or that samples of  $\sim e^{-RG(\cdot)}$  are similar to  $\underline{x}$ . There are counterexamples to this. (E.g.,  $R = \|\cdot\|_1$ .)

- Assume now that the noise  $\nu$  is Gaussian. Then an appropriate **data fidelity term** is the quadratic norm  $\frac{1}{2} \|\cdot\|_2^2$ .

- We can recover a solution with the data fidelity term and the prior  $R$  by minimizing their weighted sum:

$$\hat{\underline{x}}_R(y, t) \in \operatorname{argmin}_{\underline{x} \in \operatorname{dom} R} \frac{1}{2t} \|y - A(\underline{x})\|_2^2 + R(\underline{x}), \quad (2)$$

where  $t > 0$  is a hyperparameter.

(Note: We minimize over  $\operatorname{dom} R$  to allow for the prior  $R$  to take the extended value  $+\infty$  on some subset of  $\mathbb{R}^n$ . This allows for, e.g., priors defined on bounded domains.)

- Key to our work is the **proximal operator** of prior  $R$ , which we will denote by  $\operatorname{prox}_R$ :

$$\operatorname{prox}_R(y) \in \operatorname{argmin}_{\underline{x} \in \operatorname{dom} R} \frac{1}{2} \|y - \underline{x}\|_2^2 + R(\underline{x}). \quad (3)$$

When  $R$  is proper (i.e. not identically  $+\infty$  and nowhere equal to  $-\infty$ ), lower semicontinuous and convex, then  $\operatorname{prox}_R$  is single-valued.

When  $R$  is non-convex, then we do not have uniqueness of solutions in (3).

3/

- Following Gibonval and Nikolova (2020), we define the proximal operator of  $R$  as a selection over the solutions of (3).

- More precisely, a function  $f: \text{dom } R \rightarrow \mathbb{R}^n$  is a proximity operator of a prior  $R: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  if

$$f(y) \in \underset{x \in \text{dom } R}{\operatorname{argmin}} \frac{1}{2} \|y - x\|_2^2 + R(x) \quad (4)$$

for each  $y \in \text{dom } R$ .

- For short, we will write (4) as  $f(y) \in \text{prox}_R(y)$ .

- Theorem 1 + Corollary 1 from Gibonval and Nikolova (2020) yield the following:

Let  $\text{dom } R$  be non-empty and open. Then  $f$  is a (continuous) proximal operator of  $R$  if and only if there exists a convex, differentiable function  $\psi$  on  $\text{dom } R$  such that

$$f(y) = \nabla_y \psi(y) \text{ for each } y \in \text{dom } R.$$

Moreover, we have the identity

$$\psi(y) + \left( R(\nabla_y \psi(y)) + \frac{1}{2} \|\nabla_y \psi(y)\|_2^2 \right) = \langle y, \nabla_y \psi(y) \rangle. \quad (5)$$

L

# 11

## Neural Networks for parametrizing gradients of convex functions

- We seek a neural network (NN)  $\psi_\Theta : \mathbb{R}^n \rightarrow \mathbb{R}$  parametrized by some parameters  $\Theta$  such that  $\psi_\Theta$  is convex w.r.t. its input and such that

$$\phi_\Theta = \nabla \psi_\Theta$$

is a "good" approximation to the (continuous) proximal operator  $f$  of the prior  $R$ .

- Candidates: NNs whose
  - i) Nonlinear activation functions are convex and non-decreasing.
  - ii) Network weights are non-negative.

One way of achieving this is as follows:

- Pick some  $g : \mathbb{R} \rightarrow \mathbb{R}$  that is convex, non-decreasing and  $C^2$ .
  - Define  $g : \mathbb{R}^n \rightarrow \mathbb{R}^n$  via  $\underline{g}(z) = (g(z_1), g(z_2), \dots, g(z_n))$ .
  - Let  $\Theta = \{\underline{w}, \underline{b}, \{\underline{W}_k\}_{k=1}^K, \{\underline{H}_k, \underline{b}_k\}_{k=1}^K\}$  be a set of learnable parameters. Here,
- $\underline{w} \in \mathbb{R}^n, \underline{b} \in \mathbb{R}, \underline{W}_k \in \mathbb{R}^{n \times n}$  for  $k=1, \dots, K$ ,
- $(\underline{H}_k, \underline{b}_k) \in \mathbb{R}^{n \times n} \times \mathbb{R}^n$  for  $k=1, \dots, K$ .

- Let  $\boxed{\psi_\Theta(y) = \langle \underline{w}, z_k \rangle + \underline{b}}$ , where

$$z_1 = g(\underline{H}_1 y + \underline{b}_1) \text{ and } z_k = g(\underline{W}_k z_{k-1} + \underline{H}_k y + \underline{b}_k) \text{ for } k \in \{2, \dots, K\}.$$

5/

- Note 1: The conditions i) and ii) ensure that the function  $y \mapsto \gamma_0(y)$  is convex. This follows from the conditions needed to build a convex function from other convex functions, including:

- a) Positive combinations of convex functions is convex (all weights are non-negative).
- b) Pre-composition of a convex function with an affine mapping is convex (domain of the convex function and the image of the mapping must have non-zero intersection).
- c) Post-composition with a non-decreasing convex function is convex.

See Hiriart-Urruty and Lemaréchal (1996), II, Section 2.1.

- Note 2: Other convexity-preserving operations (e.g., supremum of convex functions, conjugations, dilations / perspectives) could be used, in principle.

Moreover, different/multiple convex, non-increasing  $C^2$  functions could be used in the layers instead of just one (i.e.,  $(g_1, g_2, \dots, g_k)$  instead of  $g$ .)

v ~~~~~

61

Connections between the convex function  $\Psi$  and the original minimization problem.

Let  $S: \mathbb{R}^n \rightarrow \mathbb{R}$  denote the value of the minimization problem underlying  $f \in \text{prox}_{\mathbb{R}}$ :

$$\begin{aligned} S(y) &:= \min_{x \in \text{dom } R} \frac{1}{2} \|x - y\|_2^2 + R(x). \\ &\equiv \frac{1}{2} \|f(y) - y\|_2^2 + R(f(y)). \end{aligned} \quad (6)$$

- What is the connection between the convex function  $\Psi$  from  $f = \nabla_y \Psi$  and the function  $S$ ?

- Expand the quadratic term in (6) and rearrange:

$$R(f(y)) + \frac{1}{2} \|f(y)\|_2^2 + \frac{1}{2} \|y\|_2^2 - S(y) = \langle y, f(y) \rangle$$

Comparing this with identity (5), we find

$$\Psi(y) = \frac{1}{2} \|y\|_2^2 - S(y)$$

$\Rightarrow$  This connects the convex function  $\Psi$  to the minimal value function  $S$ . In particular, if we are given samples/observations  $\{y_i, S(y_i)\}_{i=1}^n$ , then we effectively know the true value of  $\Psi$  at the points  $\{y_i\}_{i=1}^n$ .

7/

- Suppose we receive samples  $\{(y_i, S(y_i))\}_{i=1}^L$   
(e.g., observations / measurements from sensor).

We can set  $\psi(y_i) = \frac{1}{2} \|y_i\|_2^2 - S(y_i)$  and train a neural network by minimizing the loss function

$$\min_{\Theta} \frac{1}{2} \sum_{i=1}^L (\psi_\Theta(y_i) - \psi(y_i))^2$$

- After training, we obtain the function  $\hat{\psi}_\Theta(y)$  with trained parameters  $\Theta$ . In particular, we have

$$R(\nabla_y \hat{\psi}_\Theta(y_i)) = S(y_i) - \frac{1}{2} \|\nabla_y \hat{\psi}_\Theta(y_i) - y_i\|_2^2$$

- What if we want to evaluate the prior  $R$  at  $x$ ?

Write

$$R(x) = R(\nabla_y \hat{\psi}_\Theta(y^*)) \text{ such that } \nabla_y \hat{\psi}_\Theta(y^*) = x.$$

Fang, Buchanan and Sulam (2021) suggests to solve this as follows: Set  $\gamma > 0$  and minimize

$$\min_{y \in \mathbb{R}^n} \left\{ \hat{\psi}_\Theta(y) - \langle y, x \rangle + \frac{\gamma}{2} \|y\|_2^2 \right\}$$

This is  $\gamma$ -strongly convex with a global, unique minimum  $\hat{y}(\gamma)$  characterized by the first-order conditions

$$\nabla_y \hat{\psi}_\Theta(y^*(\gamma)) + \gamma \hat{y}'(\gamma) = x.$$

Then write  $R(x) \approx R(\nabla_y \hat{\psi}_\Theta(y^*(\gamma)) + \gamma \hat{y}'(\gamma))$ .

81

Note 1: This possibly can be done sequentially with a monotonically decreasing sequence  $\{x_j\}_{j=1}^N$ , using  $y^*(x_j)$  as a warm start for the subsequent minimum  $y^*(x_{j+1})$ .

Note 2: GPL thinks there is a rigorous basis to this procedure; he will try to find and cite the correct result to help with this.

TBC.