

Learning-based Image Reconstruction via Parallel Proximal Algorithm

Emrah Bostan*, Ulugbek S. Kamilov[†] and Laura Waller*

January 30, 2018

Abstract

In the past decade, sparsity-driven regularization has led to advancement of image reconstruction algorithms. Traditionally, such regularizers rely on analytical models of sparsity (e.g. total variation (TV)). However, more recent methods are increasingly centered around data-driven arguments inspired by deep learning. In this letter, we propose to generalize TV regularization by replacing the ℓ_1 -penalty with an alternative prior that is trainable. Specifically, our method learns the prior via extending the recently proposed fast parallel proximal algorithm (FPPA) to incorporate data-adaptive proximal operators. The proposed framework does not require additional inner iterations for evaluating the proximal mappings of the corresponding learned prior. Moreover, our formalism ensures that the training and reconstruction processes share the same algorithmic structure, making the end-to-end implementation intuitive. As an example, we demonstrate our algorithm on the problem of deconvolution in a fluorescence microscope.

1 Introduction

The problem of reconstructing an image from its noisy linear observations is fundamental in signal processing. Formulating the reconstruction as a linear inverse problem

$$\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{e}, \quad (1)$$

the unknown image $\mathbf{x} \in \mathbb{R}^N$ is computed from measurements $\mathbf{y} \in \mathbb{R}^M$. Here, the matrix $\mathbf{H} \in \mathbb{R}^{M \times N}$ models the response of the acquisition device, while $\mathbf{e} \in \mathbb{R}^M$ represents the measurement noise. In practice, the reconstruction often relies on the regularized least-squares approach:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_{\ell_2}^2 + \tau \mathcal{R}(\mathbf{x}) \right\}, \quad (2)$$

where \mathcal{R} is a regularization functional that promotes some desired properties in the solution and $\tau > 0$ controls the strength of the regularization.

In most reconstruction schemes, an analytical prior model is used. One of the most popular regularizers for images is total variation (TV) [1], defined as $\mathcal{R}_{\text{TV}}(\mathbf{x}) \triangleq \|\mathbf{D}\mathbf{x}\|_{\ell_1}$, where \mathbf{D} is the discrete gradient operator. The TV functional is a sparsity-promoting prior (via the ℓ_1 -norm) on the image gradient. Used in compressed sensing [2, 3], TV regularization has been central to inverse problems and successfully applied to a wide range of imaging applications [4–7].

*E. Bostan (email: bostan@berkeley.edu) and L. Waller (email: waller@berkeley.edu) are with the Department of Electrical Engineering & Computer Sciences, University of California, Berkeley, CA 94720, USA.

E. Bostan’s research is supported by the Swiss National Science Foundation (SNSF) under grant P2ELP2 172278.

[†]U. S. Kamilov (email: kamilov@wustl.edu) is with Computational Imaging Group (CIG), Washington University in St. Louis, St. Louis, MO 63130, USA.

Two commonly used methods for performing TV regularized image reconstructions are the (fast) iterative shrinkage/thresholding algorithm ((F)ISTA) [8] and alternating direction method of multipliers (ADMM) [9]. These algorithms reduce the complex optimization problem to a sequence of simpler operations applied to the iterates. Both methods require evaluating the proximal mapping of the TV regularizer at each iteration [10]. This amounts to solving a *denoising* problem that does not depend on \mathbf{H} and imposes piecewise-smoothness on the reconstruction [11].

From a fundamental standpoint, the modular structure of FISTA and ADMM algorithms separates the prior model (specified by the proximal) from the underlying physical model \mathbf{H} . To develop more effective regularizers than TV, researchers have thus modified the proximal operators based on practical grounds (notably, the subsequent mean-squared-error (MSE) performance) rather than analyticity. One class of algorithms called “plug-and-play” (PnP) [12–16] replaces the proximal step with powerful denoising techniques such as BM3D [17]. More recently, motivated by the success of neural networks [18] in image analysis applications [19], learning-based methods have also been proposed for designing regularization strategies. One popular approach is to *unfold* a specific iterative reconstruction algorithm that is derived for a TV-like regularization and consider a parametrized proximal step instead of a fixed one. Through the learning of parametrization coefficients in a data-driven fashion, such algorithms have adapted the regularizer to the underlying properties (deterministic and/or stochastic) of the data [20–24].

The efficiency of designing trainable regularizers is primarily determined by the algorithm that is chosen to be unfolded. The major challenge is that many proximal operators, such as that of TV, do not admit closed form solutions and require additional iterative solvers for computation [8, 25]. This complication might limit the learning process to differentiable models [22]. Alternatively, ISTA-based schemes can be used without such confinements for learning proximals that are simpler [23]. Using variable-splitting [9], ADMM-based learning has addressed these proximal-related problems. However, the final reconstruction algorithm obtained by this formulation is efficient only for a restricted class of forward models due to the inherent properties of ADMM [21, 26]. Moreover, since variable-splitting introduces auxiliary variables, such methods also require more memory, which becomes a bottle-neck for large-scale imaging problems [27].

In this letter, we propose a new learning-based image reconstruction method called the *trainable* parallel proximal algorithm (TPPA). Our algorithm extends the recently proposed fast parallel proximal algorithm (FPPA) [28] to its data-adaptive variant. At its core, FPPA uses a simple wavelet-domain soft-thresholding to compute the proximal of TV, eliminating the need for an additional iterative solver. Building upon this aspect, our framework: **1**) efficiently learns a TV-type regularization by replacing the soft-thresholding function by a parametric representation that is then learned for a given data-class, **2**) is general and does not put any restrictions on the forward model \mathbf{H} . We also show that the training and reconstruction processes share the same algorithmic structure, making TPPA’s end-to-end implementation very convenient. We apply the proposed method to the problem of deconvolution in fluorescence microscopy. Our results show that the learned regularization improves the deconvolution accuracy compared to TV and PnP models.

2 Mathematical Background

Our formalism starts with discussing the fundamentals of the FPPA method. This is then followed by the derivation of our method, which is the data-driven variant of FPPA.

2.1 FPPA for TV regularization

First, we provide some background on TV regularization via FPPA. The method uses wavelets to define (and generalize) the TV regularizer. To see this, we first define a transform $\mathbf{W} : \mathbb{R}^N \rightarrow \mathbb{R}^{N \times 4}$ that consists of the gradient operator $\mathbf{D} = (\mathbf{D}_x, \mathbf{D}_y)$, as well as an averaging operator $\mathbf{A} = (\mathbf{A}_x, \mathbf{A}_y)$. The averaging operator \mathbf{A} computes pairwise averages of pixels along each dimension. We rescale both operators by $1/(2\sqrt{2})$ for notational convenience. Note that combining these operators makes \mathbf{W} an invertible transform and it holds that $\mathbf{W}^T \mathbf{W} = \mathbf{I}$, which is not the case for \mathbf{D} alone. However, note that $\mathbf{W} \mathbf{W}^T \neq \mathbf{I}$ due to \mathbf{W} being redundant [29].

\mathbf{W} can be rewritten as a union of four orthogonal transforms $\{\mathbf{W}_k\}_{k \in [1..4]}$, allowing \mathbf{W} to be interpreted as the union of scaled and shifted Haar wavelets and scaling functions [30]. This viewpoint provides us with the central idea of FPPA, which recasts the TV regularizer by using the four orthogonal Haar transforms:

$$\mathcal{R}_{\text{TV}}(\mathbf{x}) = \tau\sqrt{2} \sum_{k=1}^4 \sum_{n \in \mathcal{H}_k} |[\mathbf{W}_k \mathbf{x}]_n|. \quad (3)$$

$\mathcal{H}_k \subset [1, \dots, N]$ is the set of all the detail (*i.e.* difference) coefficients of the transform \mathbf{W}_k . This relationship is then used to design the following updates at iteration t :

$$\mathbf{s}^t \leftarrow \mu_t \mathbf{x}^{t-1} + (1 - \mu_t) \mathbf{x}^{t-2} \quad (4a)$$

$$\mathbf{z}^t \leftarrow \mathbf{s}^t - \gamma_t \mathbf{H}^T (\mathbf{H} \mathbf{s}^t - \mathbf{y}) \quad (4b)$$

$$\mathbf{x}^t \leftarrow \mathbf{W}^T \mathcal{T}(\mathbf{W} \mathbf{z}^t, 2\sqrt{2}\tau\gamma_t), \quad (4c)$$

where the scalar soft-thresholding function

$$\mathcal{T}(z, \tau) \triangleq \text{sgn}(z) \max(|z| - \tau, 0), \quad (5)$$

is applied element-wise on the detail coefficients. As in the FISTA implementation of TV (TV-FISTA) [8], the parameters $\{\mu_t\}$ are set as [31]

$$\mu_t = 1 - \frac{1 - q_{t-1}}{q_t}, \text{ with } q_t = \frac{1}{2}(1 + \sqrt{1 + 4q_{t-1}^2}) \quad (6)$$

and $q_0 = 1$. Note that FPPA exploits the well-known connection between the Haar wavelet-transform and TV, and it is closely related to a technique called cycle spinning [32–35].

The convergence rate of FPPA is given by [28]

$$\mathcal{C}(\mathbf{x}^t) - \mathcal{C}(\mathbf{x}^*) \leq \frac{2}{\gamma(t+1)^2} \|\mathbf{x}^0 - \mathbf{x}^*\|_{\ell_2}^2 + 4\gamma G^2, \quad (7)$$

where $\{\mathbf{x}^t\}$ are the iterates from (4), \mathcal{C} is the true TV cost functional, and \mathbf{x}^* is a minimizer of \mathcal{C} . This means that for a constant step-size $\gamma > 0$, convergence can be established in the neighborhood of the optimum, which can be made arbitrarily close by letting $\gamma \rightarrow 0$. Additionally, the global convergence rate of FPPA $O(1/t^2)$ matches that of TV-FISTA [8].

FPPA that works with a fixed regularizer such as TV. The idea and convergence of FPPA can be generalized to regularizers beyond TV by using other wavelet transform and considering multiple resolutions.

3 Proposed Approach: Trainable Parallel Proximal Algorithm (TPPA)

We now present our method, which adapts the regularization to the data rather than being designed for a fixed one. Given \mathbf{H} and \mathbf{W} , we see that the shrinkage function solely determines the reconstruction. We have noted that the TV reconstruction is strictly linked to the soft-thresholding within the scheme outlined in (4). However, the efficiency of a shrinkage function varies with the type of object being imaged [36]. This necessitates revisiting FPPA to obtain a data-specific reconstruction algorithm.

Our model keeps $\{\mathbf{W}_k\}_{k \in [1..4]}$ as the pairwise averages and differences and considers an iteration-dependent sequence of shrinkage functions for each wavelet channel $k \in [1 \dots 4]$. We adopt the following parametrization:

$$\mathcal{T}_k^t(x) = \sum_{p=-P}^P c_{kp}^t \varphi\left(\frac{x}{\Delta} - p\right), \quad (8)$$

where $\{c_{kp}^t\}$ are the expansion coefficients and φ is a basis function positioned on the grid $\Delta[-P \dots P] \subset \Delta\mathbb{Z}$. We additionally reparametrize each step-size $\gamma_t > 0$ with a scalar $\alpha_t \in \mathbb{R}$ and a one-to-one function

$$\gamma = \phi(\alpha) = \begin{cases} e^{\alpha-1} & \text{if } \alpha \leq 1, \\ \alpha & \text{otherwise.} \end{cases} \quad (9)$$

This representation facilitates automatic tuning of the step-sizes $\{\gamma_t\}$ while ensuring their non-negativity. We note that the overall parametrization can be restricted to its iteration-independent counterpart (*i.e.* same set of parameters for each iteration). Moreover, by appropriately constraining the parameters to lie in a well-characterized subspace, the convergence rate given in (7) can be preserved. However, such constraints are potentially restrictive on the reconstruction performance [37].

At iteration t , the TPPA updates are

$$\mathbf{s}^t \leftarrow \mu_t \mathbf{x}^{t-1} + (1 - \mu_t) \mathbf{x}^{t-2} \quad (10a)$$

$$\mathbf{z}^t \leftarrow \mathbf{s}^t - \phi(\alpha_t) \mathbf{H}^T (\mathbf{H} \mathbf{s}^t - \mathbf{y}) \quad (10b)$$

$$\mathbf{x}^t \leftarrow \sum_{k=1}^4 \mathbf{W}_k^T \mathcal{T}_k^t (\mathbf{W}_k \mathbf{z}^t), \quad (10c)$$

where the scaling factors are absorbed into the coefficients. In contrast to (4), TPPA uses a sequence of adjustable shrinkage functions for each \mathbf{W}_k in addition to self-tuning the step-size. More importantly, compared to similar approaches based on ADMM [26, 37], TPPA does not rely on $\mathbf{H}^T \mathbf{H}$ being a structured matrix (such as block-circulant) for computational efficiency.

3.1 Training of Model Parameters

We now consider determining our model parameters (*i.e.* shrinkage functions and step-sizes) via an offline training. Through a collection of training pairs $\{(\mathbf{x}_\ell, \mathbf{y}_\ell)\}_{\ell \in [1 \dots L]}$, our goal is to learn

$$\boldsymbol{\theta} = \{\boldsymbol{\theta}^t\}_{t \in [1 \dots T]},$$

where $\boldsymbol{\theta}^t \triangleq \{\alpha_t, \mathbf{c}^t\}$, with $\mathbf{c}^t = \{c_{kp}^t\}_{k \in [1 \dots 4], p \in [-P \dots P]}$ denoting the vector of coefficients. The total number of trainable parameters is $\dim(\boldsymbol{\theta}) = T + TK(2P + 1)$. We define the cost for parameter learning to be the mean squared error (MSE) over the training data

$$\mathcal{E}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{\ell=1}^L \|\widehat{\mathbf{x}}(\boldsymbol{\theta}, \mathbf{y}_\ell) - \mathbf{x}_\ell\|_{\ell_2}^2, \quad (11)$$

where $\widehat{\mathbf{x}}(\boldsymbol{\theta}, \mathbf{y})$ is the output of (10) for a given measurement vector \mathbf{y}_ℓ and set of parameters $\boldsymbol{\theta}$ after a fixed number of iterations T . The learned parameters are thus obtained via

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \mathcal{E}(\boldsymbol{\theta}). \quad (12)$$

The implication of (12) is immediate: Given a fixed computation cost (that is expectedly cheaper than that for TV), the shrinkages are optimized to maximize the reconstruction accuracy over the training dataset.

Note that the optimization problem in (12) is smooth and hence first-order optimization methods are convenient. We use the gradient descent algorithm with Nesterov's acceleration scheme [31] (see Algorithm 1¹).

We now explain how the gradient of the cost function in (12) is derived. We denote the gradient by $\nabla \mathcal{E}$ and rely on backpropagation [38] for obtaining its analytical expression. Here, we point out the main aspects

¹Note that iterates of the parameters are represented by $\boldsymbol{\theta}^{(i)}$ to distinguish from $\boldsymbol{\theta}^t$.

Algorithm 1 Parameter training

Input: a training pair $(\mathbf{x}_\ell, \mathbf{y}_\ell)$, learning rate ν , number of training iterations I .

Output: optimized parameters $\boldsymbol{\theta}^*$.

Initialize: $\boldsymbol{\theta}^{(0)}$ and $\phi^{(0)}$

Set: $q_0 = 1$.

For $i = 1, 2, \dots, I$, compute

$$\boldsymbol{\theta}^{(i)} \leftarrow \phi^{(i-1)} - \nu \nabla \mathcal{E}(\phi^{(i-1)}) \text{ (use Algorithm 2),}$$

$$q_i \leftarrow \left(1 + \sqrt{1 + 4q_{i-1}^2}\right) / 2,$$

$$\phi^{(i)} \leftarrow \boldsymbol{\theta}^{(i)} + (q_{i-1} - 1/q_i)(\boldsymbol{\theta}^{(i)} - \boldsymbol{\theta}^{(i-1)}),$$

and return $\boldsymbol{\theta}^{(I)}$.

of our derivation since such calculations are lengthy. First, we define the residual term $\mathbf{r}^t \triangleq [\partial \mathcal{E} / \partial \mathbf{x}^t]^\top$ and use the chain rule to get:

$$\begin{aligned} \mathbf{r}^{t-2} &= \left[\frac{\partial \mathcal{E}}{\partial \mathbf{x}^{t-2}} \right]^\top \\ &= \left[\frac{\partial \mathbf{x}^{t-1}}{\partial \mathbf{x}^{t-2}} \right]^\top \mathbf{r}^{t-1} + \left[\frac{\partial \mathbf{x}^t}{\partial \mathbf{x}^{t-2}} \right]^\top \mathbf{r}^t. \end{aligned} \quad (13)$$

Using matrix calculus, the derivatives are as thus

$$\frac{\partial \mathbf{x}^{t-1}}{\partial \mathbf{x}^{t-2}} = \mu_{t-1} \mathbf{W}^\top \text{diag}(\mathcal{T}'_{t-1}(\mathbf{u}^t)) \mathbf{W} (\mathbf{I} - \gamma_{t-1} \mathbf{H}^\top \mathbf{H}),$$

where $\mathbf{u}^t = \mathbf{W} \mathbf{z}^t$. Similarly, we compute

$$\frac{\partial \mathbf{x}^t}{\partial \mathbf{x}^{t-2}} = (1 - \mu_t) \mathbf{W}^\top \text{diag}(\mathcal{T}'_t(\mathbf{u}^t)) \mathbf{W} (\mathbf{I} - \gamma_t \mathbf{H}^\top \mathbf{H}).$$

As for the derivatives of the training parameters, we use the chain rule once again to attain the following:

$$\begin{aligned} \left[\frac{\partial \mathcal{E}}{\partial \alpha^t} \right]^\top &= -\phi'(\alpha^t) (\mathbf{H} \mathbf{s}^t - \mathbf{y}) \mathbf{H} \mathbf{W}^\top \text{diag}(\mathcal{T}'_t(\mathbf{u}^t)) \mathbf{W} \mathbf{r}^t; \\ \left[\frac{\partial \mathcal{E}}{\partial \mathbf{c}_k^t} \right]^\top &= (\boldsymbol{\Phi}_k^t)^\top \mathbf{W}_k \mathbf{r}^t, \end{aligned}$$

where $\boldsymbol{\Phi}_k^t \in \mathbb{R}^{N \times (2P+1)}$ is the matrix representation of the basis functions in the sense that $\mathcal{T}'_k(\mathbf{u}_k^t) = \boldsymbol{\Phi}_k^t \mathbf{c}_k^t$. Considering these partial derivatives along with (13), we obtain the scheme described in Algorithm 2 to compute the backpropagation. Finally, we note that the reconstruction in (10) and Algorithms 1 and 2 essentially share the same structure (*i.e.* gradient descent with acceleration). This makes the proposed model convenient, since the computational implementation can be reused.

4 Numerical Results

We now present *in silico* experiments corroborating TPPA, with deconvolution of fluorescence microscopy images where the point spread function (PSF) of the microscope is approximated by a Gaussian kernel of variance 2 pixels. The imaging process is assumed not to be photon-limited; noise is modeled as additive white Gaussian noise (AWGN) of 30 dB SNR.

Algorithm 2 Backpropagation for Algorithm 1

Input: a training pair $(\mathbf{x}_\ell, \mathbf{y}_\ell)$, the set of parameters θ , number of TPPA iterations T .

Output: components of the gradient $\nabla \mathcal{E}$.

Set: $\mu^{T+1} = 1$, $\mathbf{v}^{T+1} = \mathbf{0}$, and $\mathbf{r}^T = (\hat{\mathbf{x}}(\theta, \mathbf{y}_\ell) - \mathbf{x}_\ell)$.

For $t = T, T-1, \dots, 1$, compute

$$\mathbf{b}^t \leftarrow \mathbf{W}^T \text{diag}(\mathcal{T}'_t(\mathbf{u}^t)) \mathbf{W} \mathbf{r}^t,$$

$$\mathbf{v}^t \leftarrow \mathbf{b}^t - \gamma_t \mathbf{H}^T \mathbf{H} \mathbf{b}^t,$$

$$\mathbf{r}^{t-1} \leftarrow \mu_t \mathbf{v}^t + (1 - \mu_{t+1}) \mathbf{v}^{t+1},$$

and store

$$\left[\frac{\partial \mathcal{E}}{\partial \alpha^t} \right]^T = -\phi'(\alpha^t) (\mathbf{H} \mathbf{s}^t - \mathbf{y}_\ell)^T \mathbf{H} \mathbf{b}^t,$$

$$\left[\frac{\partial \mathcal{E}}{\partial \mathbf{c}_k^t} \right]^T = (\Phi_k^t)^T \mathbf{W}_k \mathbf{r}^t \quad (k = 1, \dots, 4).$$

Table 1: Average deconvolution performance (on the validation set) of the methods considered in the experiments. Numbers indicate SNR in decibel units.

PSF Kernel Size	Deconvolution Algorithm		
	TV	PnP (using BM3D)	Proposed Method
5×5	21.99	24.69	24.96
9×9	20.77	22.33	22.57

Fluorescence microscopy images of human bone osteosarcoma epithelial cells (U2OS Line)², from [39] are used as our ground-truth data. All images' intensity are scaled between 0 and 1. To generate the training pairs, we use 100 images and apply the forward model to a single patch (per image) of size 64×64 extracted around the center of the field-of-view. Once the images chosen for training are excluded, we select a different 20 images of size 256×256 as a validation set.

Learning is carried out by using Algorithm 1 with 200 iterations (that is $I = 200$) with $\nu = 5 \times 10^{-4}$. We set the number of layers for TPPA as $T = 10$. The shrinkage functions are parametrized by 10^3 equally-spaced cubic B-splines over the dynamic range of $\mathbf{W}\mathbf{x}$. All shrinkages are initialized with the identity operator. Finally, we note that $\alpha_0 = 1/\|\mathbf{H}^T \mathbf{H}\|_2^2$ and $\mathbf{x}^0 = \mathbf{0} \in \mathbb{R}^N$.

As a baseline comparison, we consider TV regularization implemented using FPPA described in (4). The algorithm is run until either 100 iterations is reached or $\|\mathbf{x}_t - \mathbf{x}_{t-1}\|_2 / \|\mathbf{x}_{t-1}\|_2 \leq 10^{-6}$ is satisfied. We also compare against the PnP model where the proximal of TV is replaced by BM3D [17]. The latter is implemented using 10 FISTA iterations (same as the number of layers in TPPA) and all methods use a zero initialization. For each validation image, we optimize the regularization parameters for both algorithms (by using an oracle) for the best-possible SNR performance. Average SNRs of the reconstruction are reported in Table 1 for different sizes of the blur kernel.

The results show that the accuracy of our model is better than the other algorithms considered. In particular, the SNR performance provided by TPPA is significantly better than that of TV. Furthermore, visual inspection of the reconstructions reveals that the TV deconvolution creates the characteristic blocky artifacts at textured regions (see Figure 1). Since the successive sequence of shrinkage functions are adapted to the

²The dataset consists of multi-color images where each color channel depicts different flurophores targeted to different organelles or cellular components. In our simulations, we use the channel corresponding to the blue color which targets the cell nucleus.

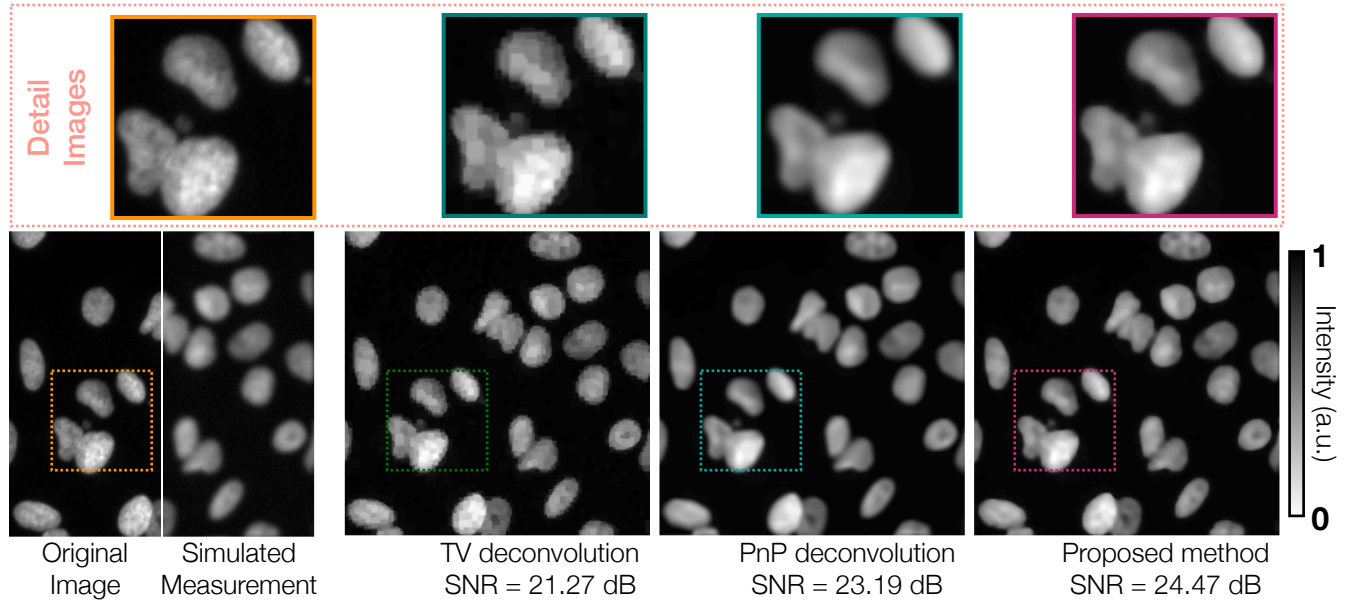


Figure 1: An example of 2D deconvolution with fluorescence microscopy images, given a known 9×9 Gaussian PSF: The ground-truth image (left) illustrates the nucleus of a group of U2OS cells. Our learning-based reconstruction preserves homogeneity of the background and more of the texture, increasing the SNR. See text for further details.

underlying features of the training data, one notices that these artifacts are reduced for TPPA. Our method also renders the boundary of the nucleus more faithfully and provides a homogeneous background. These observations confirm the efficiency of our data-specific approach (in terms of deconvolution quality) and highlight its potential importance in practical scenarios.

5 Conclusion

We developed a learning-based algorithm for linear inverse problems that is in the spirit of TV regularization. Our approach, TPPA, has enabled us to move away from the soft-thresholding operator (with a fixed threshold value at all iterations) to a collection of parametrized shrinkage functions that are optimized (in MSE sense) for a training set. Compared to TV regularization and PnP technique, our deconvolution simulations demonstrate that advantages of TPPA in terms of accuracy.

References

- [1] L. I. Rudin, S. Osher, and E. Fatemi, “Nonlinear total variation based noise removal algorithms,” *Physica D*, vol. 60, no. 1–4, pp. 259–268, November 1992.
- [2] D. L. Donoho, “Compressed sensing,” *IEEE Trans. Inf. Theory*, vol. 52, no. 4, pp. 1289–1306, April 2006.
- [3] E. J. Candès and T. Tao, “Decoding by linear programming,” *IEEE Trans. Inf. Theory*, vol. 51, no. 12, pp. 4203–4215, December 2005.
- [4] M. Persson, D. Bone, and H. Elmqvist, “Total variation norm for three-dimensional iterative reconstruction in limited view angle tomography,” *Phys. Med. Biol.*, vol. 46, no. 3, pp. 853–866, 2001.

- [5] M. M. Bronstein, A. M. Bronstein, M. Zibulevsky, and H. Azhari, "Reconstruction in diffraction ultrasound tomography using nonuniform FFT," *IEEE Trans. Med. Imag.*, vol. 21, no. 11, pp. 1395–1401, November 2002.
- [6] M. Lustig, D. L. Donoho, and J. M. Pauly, "Sparse MRI: The application of compressed sensing for rapid MR imaging," *Magn. Reson. Med.*, vol. 58, no. 6, pp. 1182–1195, December 2007.
- [7] U. S. Kamilov, I. N. Papadopoulos, M. H. Shoreh, A. Goy, C. Vonesch, M. Unser, and D. Psaltis, "Optical tomographic image reconstruction based on beam propagation and sparse regularization," *IEEE Trans. Comp. Imag.*, vol. 2, no. 1, pp. 59–70, March 2016.
- [8] A. Beck and M. Teboulle, "Fast gradient-based algorithm for constrained total variation image denoising and deblurring problems," *IEEE Trans. Image Process.*, vol. 18, no. 11, pp. 2419–2434, November 2009.
- [9] M. V. Afonso, J. M. Bioucas-Dias, and M. A. T. Figueiredo, "Fast image recovery using variable splitting and constrained optimization," *IEEE Trans. Image Process.*, vol. 19, no. 9, pp. 2345–2356, September 2010.
- [10] J. J. Moreau, "Proximité et dualité dans un espace hilbertien," *Bull. Soc. Math. France*, vol. 93, pp. 273–299, 1965.
- [11] J. Cheng and B. Hofmann, *Handbook of Mathematical Methods in Imaging*. Springer, 2011, ch. Chapter 3: Regularization Methods for Ill-Posed Problems, pp. 87–109.
- [12] S. V. Venkatakrishnan, C. A. Bouman, and B. Wohlberg, "Plug-and-play priors for model based reconstruction," in *Proc. IEEE Global Conf. Signal Process. and Inf. Process. (GlobalSIP)*, Austin, TX, USA, December 3-5, 2013, pp. 945–948.
- [13] S. Sreehari, S. V. Venkatakrishnan, B. Wohlberg, G. T. Buzzard, L. F. Drummy, J. P. Simmons, and C. A. Bouman, "Plug-and-play priors for bright field electron tomography and sparse interpolation," *IEEE Trans. Comp. Imag.*, vol. 2, no. 4, pp. 408–423, December 2016.
- [14] S. H. Chan, X. Wang, and O. A. Elgendy, "Plug-and-play ADMM for image restoration: Fixed-point convergence and applications," *IEEE Trans. Comp. Imag.*, vol. 3, no. 1, pp. 84–98, March 2017.
- [15] U. S. Kamilov, H. Mansour, and B. Wohlberg, "A plug-and-play priors approach for solving nonlinear imaging inverse problems," *IEEE Signal. Proc. Let.*, vol. 24, no. 12, pp. 1872–1876, December 2017.
- [16] Y. Romano, M. Elad, and P. Milanfar, "The little engine that could: Regularization by denoising (RED)," *SIAM Journal on Imaging Sciences*, vol. 10, no. 4, pp. 1804–1844, 2017.
- [17] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian, "Image denoising by sparse 3-D transform-domain collaborative filtering," *IEEE Trans. Image Process.*, vol. 16, no. 16, pp. 2080–2095, August 2007.
- [18] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, pp. 436–444, May 28, 2015.
- [19] M. T. McCann, K. H. Jin, and M. Unser, "Convolutional neural networks for inverse problems in imaging: A review," *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 85–95, 2017.
- [20] A. Barbu, "Training an active random field for real-time image denoising," *IEEE Trans. Image Process.*, vol. 18, no. 11, pp. 2451–2462, November 2009.
- [21] U. Schmidt and S. Roth, "Shrinkage fields for effective image restoration," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, Columbus, OH, USA, June 23-28, 2014, pp. 2774–2781.
- [22] Y. Chen, W. Yu, and T. Pock, "On learning optimized reaction diffusion processes for effective image restoration," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, Boston, MA, USA, June 8-10, 2015, pp. 5261–5269.

- [23] U. S. Kamilov and H. Mansour, “Learning optimal nonlinearities for iterative thresholding algorithms,” *IEEE Signal Process. Lett.*, vol. 23, no. 5, pp. 747–751, May 2016.
- [24] D. Mahapatra, S. Mukherjee, and C. S. Seelamantula, “Deep sparse coding using optimized linear expansion of thresholds,” 2017, arXiv:1705.07290 [cs.LG].
- [25] Z. Qin, D. Goldfarb, and S. Ma, “An alternating direction method for total variation denoising,” *Optim. Method Softw.*, vol. 30, no. 3, pp. 594–615, 2015.
- [26] Y. Yang, J. Sun, H. Li, and Z. Xu, “Deep ADMM-Net for compressive sensing MRI,” in *Adv. Neural Inf. Process. Syst. 29 (NIPS 2016)*, 2016, pp. 1–9.
- [27] N. Antipa, G. Kuo, R. Heckel, B. Mildenhall, E. Bostan, R. Ng, and L. Waller, “DiffuserCam: Lensless single-exposure 3D imaging,” *Optica*, vol. 5, no. 1, pp. 1–9, 2018.
- [28] U. S. Kamilov, “A parallel proximal algorithm for anisotropic total variation minimization,” *IEEE Trans. Image Process.*, vol. 26, no. 2, pp. 539–548, February 2017.
- [29] M. Elad, P. Milanfar, and R. Rubinstein, “Analysis versus synthesis in signal priors,” *Inverse Problems*, vol. 23, no. 3, pp. 947–968, 2007.
- [30] S. Mallat, *A Wavelet Tool of Signal Processing: The Sparse Way*, 3rd ed. San Diego: Academic Press, 2009.
- [31] Y. E. Nesterov, “A method for solving the convex programming problem with convergence rate $O(1/k^2)$,” *Dokl. Akad. Nauk SSSR*, vol. 269, pp. 543–547, 1983, (in Russian).
- [32] R. R. Coifman and D. L. Donoho, *Springer Lecture Notes in Statistics*. Springer-Verlag, 1995, ch. Translation-invariant de-noising, pp. 125–150.
- [33] A. K. Fletcher, K. Ramchandran, and V. K. Goyal, “Wavelet denoising by recursive cycle spinning,” in *Proc. IEEE Int. Conf. Image Proc. (ICIP’02)*, Rochester, NY, USA, September 22–25, 2002, pp. II.873–II.876.
- [34] G. Steidl, J. Weickert, T. Brox, P. Mrazek, and M. Welk, “On the equivalence of soft wavelet shrinkage, total variation diffusion, total variation regularization, and SIDEs,” *SIAM J. Numer. Anal.*, vol. 42, no. 2, pp. 686–713, 2004.
- [35] U. S. Kamilov, E. Bostan, and M. Unser, “Variational justification of cycle spinning for wavelet-based solutions of inverse problems,” *IEEE Signal Process. Lett.*, vol. 21, no. 11, pp. 1326–1330, November 2014.
- [36] E. Bostan, U. S. Kamilov, M. Nilchian, and M. Unser, “Sparse stochastic processes and discretization of linear inverse problems,” *IEEE Trans. Image Process.*, vol. 22, no. 7, pp. 2699–2710, July 2013.
- [37] H. Nguyen, E. Bostan, and M. Unser, “Learning convex regularizers for optimal Bayesian denoising,” *IEEE Transactions on Signal Processing*, to appear.
- [38] H. B. Demuth, M. H. Beale, J. A. Bittker, O. D. Jess, and M. T. Hagan, *Neural Network Design*. USA: Martin Hagan, 2014.
- [39] M.-A. B. *et al.*, “A dataset of images and morphological profiles of 30,000 small-molecule treatments using the Cell Painting assay,” *GigaScience*, p. giw014, 2017.