

# Frank-Wolfe Network: An Interpretable Deep Structure for Non-Sparse Coding

Dong Liu, *Senior Member, IEEE*, Ke Sun, Zhangyang Wang, *Member, IEEE*,  
Runsheng Liu, and Zheng-Jun Zha, *Member, IEEE*

**Abstract**—The problem of  $L_p$ -norm constrained coding is to convert signal into code that lies inside an  $L_p$ -ball and most faithfully reconstructs the signal. Previous works under the name of sparse coding considered the cases of  $L_0$  and  $L_1$  norms. The cases with  $p > 1$  values, i.e. non-sparse coding studied in this paper, remain a difficulty. We propose an interpretable deep structure namely Frank-Wolfe Network (F-W Net), whose architecture is inspired by unrolling and truncating the Frank-Wolfe algorithm for solving an  $L_p$ -norm constrained problem with  $p \geq 1$ . We show that the Frank-Wolfe solver for the  $L_p$ -norm constraint leads to a novel closed-form nonlinear unit, which is parameterized by  $p$  and termed  $pool_p$ . The  $pool_p$  unit links the conventional pooling, activation, and normalization operations, making F-W Net distinct from existing deep networks either heuristically designed or converted from projected gradient descent algorithms. We further show that the hyper-parameter  $p$  can be made learnable instead of pre-chosen in F-W Net, which gracefully solves the non-sparse coding problem even with unknown  $p$ . We evaluate the performance of F-W Net on an extensive range of simulations as well as the task of handwritten digit recognition, where F-W Net exhibits strong learning capability. We then propose a convolutional version of F-W Net, and apply the convolutional F-W Net into image denoising and super-resolution tasks, where F-W Net all demonstrates impressive effectiveness, flexibility, and robustness.

**Index Terms**—Convolutional network, Frank-Wolfe algorithm, Frank-Wolfe network, non-sparse coding.

## I. INTRODUCTION

### A. $L_p$ -Norm Constrained Coding

Assuming a set of  $n$ -dimensional vectors  $\{\mathbf{x}_i \in \mathbb{R}^n | i = 1, 2, \dots, N\}$ , we aim to encode each vector  $\mathbf{x}_i$  into a  $m$ -dimensional code  $\mathbf{z}_i \in \mathbb{R}^m$ , such that the code reconstructs the vector faithfully. When  $m > n$  is the situation of interest, the most faithful coding is not unique. We hence consider the code to bear some structure, or in the Bayesian language, we want to impose some prior on the code. One possible

structure/prior is reflected by the  $L_p$ -norm, i.e. by solving the following problem:

$$\{\{\mathbf{z}_i^*\}, D^*\} = \arg \min \sum_{i=1}^N \|\mathbf{x}_i - D\mathbf{z}_i\|_2^2 \quad (1)$$

subject to  $\forall i, \|\mathbf{z}_i\|_p \leq c$

where  $D \in \mathbb{R}^{n \times m}$  is a linear decoding matrix, and both  $p$  and  $c$  are constants. In real-world applications, we have difficulty in pre-determining the  $p$  value, and the flexibility to learn the prior from the data becomes more important, which can be formulated as the following problem:

$$\{\{\mathbf{z}_i^*\}, D^*, p^*\} = \arg \min \sum_{i=1}^N \|\mathbf{x}_i - D\mathbf{z}_i\|_2^2 \quad (2)$$

subject to  $p \in \mathcal{P}, \forall i, \|\mathbf{z}_i\|_p \leq c$

where  $\mathcal{P} \subset \mathbb{R}$  is a domain of  $p$  values of interest. For example,  $\mathcal{P} = \{p | p \geq 1\}$  defines a domain that ensures the constraint to be convex with regard to  $\mathbf{z}_i$ . In this paper, (1) and (2) are termed  $L_p$ -norm constrained coding, and for distinguishing purpose we refer to (1) and (2) as known  $p$  and unknown  $p$ , respectively.

### B. Motivation

For known  $p$ , if the decoding matrix is given a priori as  $D^* = D_0$ , then it is sufficient to encode each  $\mathbf{x}_i$  individually, by means of solving:

$$\mathbf{z}_i^* = \arg \min \|\mathbf{x}_i - D_0\mathbf{z}_i\|_2^2 \quad (3)$$

subject to  $\|\mathbf{z}_i\|_p \leq c$

or its equivalent, unconstrained form (given a properly chosen Lagrange multiplier  $\lambda$ ):

$$\mathbf{z}_i^* = \arg \min \|\mathbf{x}_i - D_0\mathbf{z}_i\|_2^2 + \lambda \|\mathbf{z}_i\|_p \quad (4)$$

Eq. (4) is well known as an  $L_p$ -norm regularized least squares problem, which arises in many disciplines of statistics, signal processing, and machine learning. Several special  $p$  values, such as 0, 1, 2, and  $\infty$ , have been well studied. For example, when  $p = 0$ ,  $\|\mathbf{z}_i\|_0$  measures the number of non-zero entries in  $\mathbf{z}_i$ , thus minimizing the term induces the code that is as *sparse* as possible. While sparsity is indeed a goal in several situations, the  $L_0$ -norm minimization is NP-hard [1]. Researchers then proposed to adopt  $L_1$ -norm, which gives rise to a convex and easier problem, to approach  $L_0$ -norm. It was shown that, under some mild conditions, the resulting code of using  $L_1$ -norm coincides with that of using  $L_0$ -norm [2].  $L_1$ -norm regularization was previously known in lasso regression [3] and basis pursuit [4]. Due to its enforced sparsity, it leads

Date of current version August 19, 2019. This work was supported by the National Key Research and Development Plan under Grant 2017YFB1002401, and by the Natural Science Foundation of China under Grant 61772483.

D. Liu, K. Sun, and Z.-J. Zha are with the CAS Key Laboratory of Technology in Geo-Spatial Information Processing and Application System, University of Science and Technology of China, Hefei, China (e-mail: dongeliu@ustc.edu.cn; sunk@mail.ustc.edu.cn; zhazj@ustc.edu.cn).

Z. Wang is with Department of Computer Science and Engineering, Texas A&M University, College Station, TX, USA (e-mail: atlaswang@tamu.edu).

R. Liu was with the University of Science and Technology of China, Hefei, China (e-mail: lrs1892@outlook.com).

Copyright © 2019 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending an email to pubs-permissions@ieee.org.

to the great success of compressive sensing [5]. The cases of  $p > 1$  lead to non-sparse codes that were also investigated.  $L_2$ -norm regularization was extensively adopted under the name of weight decay in machine learning [6].  $L_\infty$ -norm was claimed to help spread information evenly among the entries of the resultant code, which is known as democratic [7] or spread representations [8], and benefits vector quantization [9], hashing [10], and other applications.

Besides the above-mentioned special  $p$  values, other general  $p$  values were much less studied due to mathematical difficulty. Nonetheless, it was observed that general  $p$  indeed could help in specific application domains, where using the special  $p$  values might be overly simplified. For sparse coding,  $L_p$ -norms with  $p \in [0, 1]$  all enforce sparsity to some extent, yet with different effects. In compressive sensing, it was known that using  $L_p$ -norm with  $0 < p < 1$  needs fewer measurements to reconstruct sparse signal than using  $L_1$ -norm; regarding computational complexity, solving  $L_p$ -norm regularization is more difficult than solving  $L_1$ -norm but still easier than solving  $L_0$ -norm [11]. Further studies found that the choice of  $p$  crucially affected the quality of results and noise robustness [12]. Xu *et al.* endorsed the adoption of  $L_{1/2}$ -norm regularization and proposed an iterative half thresholding algorithm [13]. In image deconvolution, Krishnan and Fergus investigated  $L_{1/2}$  and  $L_{2/3}$  norms and claimed their advantages over  $L_1$ -norm [14]. For non-sparse coding,  $L_p$ -norms with  $p > 1$  and different  $p$ 's have distinct impact on the solution. For example, Kloft *et al.* argued for  $L_p$ -norm with  $p > 1$  in the task of multiple kernel learning [15].

Now let us return to the problem (1) where the decoding matrix  $D$  is *unknown*. As discussed above,  $L_0$ -norm induces sparse representations, thus the special case of  $L_p$ -norm constrained coding with  $p = 0$  (similarly  $p = 1$ ) is to pursue a sparse coding of the given data. While sparse coding has great interpretability and possible relation to visual cognition [16], [17], and was widely adopted in many applications [18]–[21], we are inspired by the studies showing that general  $p$  performs better than special  $p$ , and ask: is general  $L_p$ -norm able to outperform  $L_0/L_1/L_2/L_\infty$ -norm on a specific dataset for a specific task? If yes, then which  $p$  will be the best (i.e. the case of unknown  $p$ )? These questions seem not being investigated before, to the best of our knowledge. Especially, *non-sparse coding*, i.e.  $p > 1$ , is clearly different from sparse coding and is the major theme of our study.

### C. Outline of Solution

Analytically solving the  $L_p$ -norm constrained problems is very difficult as they are not convex optimization. When designing numerical solutions, note that there are two (resp. three) groups of variables in (1) (resp. (2)), and it is natural to perform alternate optimization over them iteratively. Previous methods for sparse coding mostly follow this approach, for example in [16] and in the well-known K-SVD method [22]. One clear drawback of this approach is the high computational complexity due to the iterative nature. A variant of using early-stopped iterative algorithms to provide fast solution approximations was discussed in [23], built on the overly

strong assumption of close-to-orthogonal bases [24]. Besides the high complexity, it is not straightforward to extend the sparse coding methods to the cases of non-sparse coding, since they usually leverage the premise of sparse code heuristically but for general  $p > 1$  it is hard to design such heuristics.

A different approach for sparse coding was proposed by Gregor and LeCun [25], where an iterative algorithm known as iterative shrinkage and thresholding (ISTA), that was previously used for  $L_1$ -norm regularized least squares, is unrolled and truncated to construct a multi-layer feed-forward network. The network can be trained end-to-end to act as a regressor from a vector to its corresponding sparse code. Note that truncation helps lower the computational cost of the network than the original algorithm, while training helps compensate the error due to truncation. The trained network known as learned ISTA (LISTA) is then a fast alternative to the original ISTA algorithm. Following the approach of LISTA, several recent works consider the cases of  $L_0$ -norm [26] and  $L_\infty$ -norm [27], as well as extend other iterative algorithms to their network versions [28]–[30]. However, LISTA and its following-up works are not applicable for solving the cases of general  $p$ , because they all refer to the same category of first-order iterative algorithms, i.e. the projected gradient descent (PGD) algorithms. For general  $p$ , the projection step in PGD is not analytically solvable. In addition, such works have more difficulty in solving the unknown  $p$  problem.

Beyond the family of PGD, another first-order iterative algorithm is the Frank-Wolfe algorithm [31], which is free of projection. This characteristic inspires us to adapt the Frank-Wolfe algorithm to solve the general  $p$  problem. Similar to LISTA, we unroll and truncate the Frank-Wolfe algorithm to construct a network, termed Frank-Wolfe network (F-W Net), which can be trained end-to-end. Although the convergence speed of the original Frank-Wolfe algorithm is slow, we will show that F-W Net can converge faster than not only the original Frank-Wolfe and ISTA algorithms, but also the LISTA. Moreover, F-W Net has a novel, closed-form and nonlinear computation unit that is parameterized by  $p$ , and as  $p$  varies, that unit displays the behaviors of several classic pooling operators, and can be naturally viewed as a cascade of a generalized normalization and a parametric activation. Due to the fact that  $p$  becomes a (hyper) parameter in F-W Net, we can either set  $p$  a priori or make  $p$  learnable when training F-W Net. Thus, F-W Net has higher learning flexibility than LISTA and training F-W Net gracefully solves the unknown  $p$  problem.

### D. Our Contributions

To the best of our knowledge, we are the first to extend the sparse coding problem into  $L_p$ -norm constrained coding with general  $p > 1$  and unknown  $p$ ; we are also the first to unroll-and-truncate the Frank-Wolfe algorithm to construct trainable network. Technically, we make the following contributions:

- We propose the Frank-Wolfe network (F-W Net), whose architecture is inspired by unrolling and truncating the Frank-Wolfe algorithm for solving  $L_p$ -norm regularized least squares problem. F-W Net features a novel nonlinear

unit that is parameterized by  $p$  and termed  $pool_p$ . The  $pool_p$  unit links the conventional pooling, activation, and normalization operations in deep networks. F-W Net is verified to solve non-sparse coding with general known  $p > 1$  better than the existing iterative algorithms. More importantly, F-W Net solves the non-sparse coding with unknown  $p$  at low computational costs.

- We propose a convolutional version of F-W Net, which extends the basic F-W Net by adding convolutions and utilizing the  $pool_p$  unit point by point across different channels. The convolutional (Conv) F-W Net can be readily applied into image-related tasks.
- We evaluate the performance of F-W Net on an extensive range of simulations as well as a handwritten digit recognition task, where F-W Net exhibits strong learning capability. We further apply Conv F-W Net into image denoising and super-resolution, where it also demonstrates impressive effectiveness, flexibility, and robustness.

### E. Paper Organization

The remainder of this paper is organized as follows. Section II reviews literatures from four folds:  $L_p$ -norm constrained coding and its applications, deep structured networks, the original Frank-Wolfe algorithm, and nonlinear units in networks, from which we can see that F-W Net connects and integrates these separate fields. Section III formulates F-W Net, with detailed discussions on its motivation, structure, interpretation, and implementation issues. Section IV validates the performance of F-W Net on synthetic data under an extensive range of settings, as well as on a toy problem, handwritten digit recognition with the MNIST dataset. Section V discusses the proposed convolutional version of F-W Net. Section VI then provides experimental results of using Conv F-W Net on image denoising and super-resolution, with comparison to some other CNN models. Section VII concludes this paper. For reproducible research, our code and pre-trained models have been published online<sup>1</sup>.

## II. RELATED WORK

### A. $L_p$ -Norm Constrained Coding and Its Applications

Sparse coding as a representative methodology of the linear representation methods, has been used widely in signal processing and computer vision, such as image denoising, deblurring, image restoration, super-resolution, and image classification [18], [19], [21], [32]. The sparse representation aims to preserve the principle component and reduces the redundancy in the original signal. From the viewpoint of different norm minimizations used in sparsity constraints, these methods can be roughly categorized into the following groups: 1)  $L_0$ -norm minimization; 2)  $L_p$ -norm ( $0 < p < 1$ ) minimization; 3)  $L_1$ -norm minimization; and 4)  $L_{2,1}$ -norm minimization. In addition,  $L_2$ -norm minimization is extensively used, but it does not lead to sparse solution.

Varieties of dictionary learning methods have been proposed and implemented based on sparse representation. K-SVD

[22] seeks an over-complete dictionary from given training samples under the  $L_0$ -norm constraint, which achieves good performance in image denoising. Wright *et al.* [32] proposed a general classification algorithm for face recognition based on  $L_1$ -norm minimization, which shows if the sparsity can be introduced into the recognition problem properly, the choice of features is not crucial. Krogh *et al.* [6] showed that limiting the growth of weights through  $L_2$ -norm penalty can improve generalization in a feed-forward neural network. The simple weight decay has been adopted widely in machine learning.

### B. Deep Structured Networks

Deep networks are typically stacked with off-the-shelf building blocks that are jointly trained with simple loss functions. Since many real-world problems involve predicting statistically dependent variables or related tasks, deep structured networks [33], [34] were proposed to model complex patterns by taking into account such dependencies. Among many efforts, a noticeable portion has been devoted to unrolling the traditional optimization and inference algorithms into their deep end-to-end trainable formats.

Gregor and LeCun [25] first leveraged the idea to construct feed-forward networks as fast trainable regressors to approximate the sparse code solutions, whose idea was expanded by many successors, e.g. [26], [27], [35]–[39]. Those works show the benefits of incorporating the problem structure into the design of deep architectures, in terms of both performance and interpretability [40]. A series of works [41]–[44] established the theoretical background of this methodology. Note that many previous works [25], [26], [35], [41] built their deep architectures based on the *projected gradient descent* type algorithms, e.g. the iterative shrinkage and thresholding algorithm (ISTA). The projection step turned into the nonlinear activation function. Wang *et al.* [29] converted proximal methods to deep networks with continuous output variables. More examples include the message-passing inference machine [28], shrinkage fields [45], CRF-RNN [46], and ADMM-net [30].

### C. Frank-Wolfe Algorithm

The Frank-Wolfe algorithm [31], also known as conditional gradient descent, is one of the simplest and earliest known iterative solver, for the generic constrained convex problem:

$$\min_{\mathbf{z}} f(\mathbf{z}) \quad \text{s.t. } \mathbf{z} \in \mathcal{Z} \quad (5)$$

where  $f$  is a convex and continuously differentiable objective function, and  $\mathcal{Z}$  is a convex and compact subset of a Hilbert space. At each step, the Frank-Wolfe algorithm first considers the linear approximation of  $f(\mathbf{z})$ , and then moves towards this linear minimizer that is taken over  $\mathcal{Z}$ . Section III presents a concrete example of applying the Frank-Wolfe algorithm.

The Frank-Wolfe algorithm has lately re-gained popularity due to its promising applicability in handling structural constraints, such as sparsity or low-rank. The Frank-Wolfe algorithm is projection-free: while competing methods such as the projected gradient descent and proximal algorithms need to take a projection step back to the feasible set per iteration,

<sup>1</sup><https://github.com/sunke123/FW-Net>

the Frank-Wolfe algorithm only solves a linear problem over the same set in each iteration, and automatically stays in the feasible set. For example, the sparsity regularized problems are commonly relaxed as convex optimization over convex hulls of atomic sets, especially  $L_p$ -norm constrained domains [47], which makes the Frank-Wolfe algorithm easily applicable. We refer the readers to the comprehensive review in [48], [49] for more details about the algorithm.

### D. Nonlinear Units in Networks

There have been blooming interests in designing novel nonlinear activation functions [50], a few of which have parametric and learnable forms, such as the parametric ReLU [51]. Among existing deep structured networks, e.g. [25], [26], [35], their activation functions usually took fixed forms (e.g. some variants of ReLU) that reflected the pre-chosen structural priors. A data-driven scheme is presented in [52] to learn optimal thresholding functions for ISTA. Their adopted parametric representations led to spline curve-type activation function, which reduced the estimation error compared to using the common (fixed) piece-wise linear functions.

As another major type of nonlinearity in deep networks, pooling was originally introduced as a dimension-reduction tool to aggregate a collection of inputs into low-dimensional outputs [53]. Other than the input-output dimensions, the difference between activation and pooling also lies in that activation is typically applied element wise, while pooling is on groups of hidden units, usually within a spatial neighborhood. It was proposed in [54] to learn task-dependent pooling and to adaptively reshape the pooling regions. More learnable pooling strategies are investigated in [55], via either mixing two different pooling types or a tree-structured fusion. Gulcehre *et al.* [56] introduced the  $L_p$  unit that computed a normalized  $L_p$  norm over the set of outputs, with the value of  $p$  learnable.

In addition, we find our proposed nonlinear unit inherently related to normalization techniques. Jarrett *et al.* [53] demonstrated that a combination of nonlinear activation, pooling, and normalization improved object recognition. Batch normalization (BN) [57] rescaled the summed inputs of neurons over training batches, and substantially accelerated training. Layer normalization (LN) [58] normalized the activations across all activities within a layer. Ren *et al.* [59] re-exploited the idea of divisive normalization (DN) [60], [61], a well-grounded transformation in real neural systems. The authors viewed both BN and LN as special cases of DN, and observed improvements by applying DN on a variety of tasks.

## III. FRANK-WOLFE NETWORK

### A. Frank-Wolfe Solver for $L_p$ -Norm Constrained Least Squares

We investigate the  $L_p$ -norm constrained least squares problem as a concrete example to illustrate the construction of F-W Net. The proposed methodology can certainly be extended to more generic problems (5). Let  $\mathbf{x} \in \mathbb{R}^n$  denote the input signal, and  $\mathbf{z} \in \mathbb{R}^m$  denote the code (a.k.a. representation, feature vector).  $D \in \mathbb{R}^{n \times m}$  is the decoding matrix (a.k.a.

dictionary, bases). The problem considers  $\mathcal{Z}$  to be an  $L_p$ -norm ball of radius  $c$ , with  $p \geq 1$  to ensure the convexity of  $\mathcal{Z}$ , and  $f(\mathbf{z})$  to be a least-squares function:

$$\mathbf{z}^* = \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{x} - D\mathbf{z}\|_2^2 \quad \text{s.t.} \quad \|\mathbf{z}\|_p \leq c. \quad (6)$$

We initialize  $\mathbf{z}^0 \in \mathcal{Z}$ . At iteration  $t = 0, 1, 2, \dots$ , the Frank-Wolfe algorithm iteratively updates two variables  $\mathbf{z}^t, \mathbf{s}^t \in \mathbb{R}^m$  to solve (6):

$$\begin{aligned} \mathbf{s}^t &:= \arg \min_{\mathbf{s} \in \mathcal{Z}} \mathbf{s}^\top \nabla f(\mathbf{z}^t) \\ \mathbf{z}^{t+1} &:= (1 - \gamma^t) \mathbf{z}^t + \gamma^t \mathbf{s}^t, \end{aligned} \quad (7)$$

where  $\nabla f(\mathbf{z}^t) = D^\top(D\mathbf{z}^t - \mathbf{x})$ .  $\gamma^t$  is the step size, which is typically set as  $\gamma^t := \frac{2}{t+2}$  or chosen via line search. By Hölder's inequality, the solution of  $\mathbf{s}^t$  is:

$$s_i^t = -c \cdot \text{sign}(\nabla_i f(\mathbf{z}^t)) \frac{|\nabla_i f(\mathbf{z}^t)|^{\frac{1}{p-1}}}{\left(\sum_{j=1}^m |\nabla_j f(\mathbf{z}^t)|^{\frac{p}{p-1}}\right)^{\frac{1}{p}}}, \quad (8)$$

$$i = 1, 2, \dots, m$$

where  $s_i^t, \nabla_i f(\mathbf{z}^t)$  denote the  $i$ -th entry of  $\mathbf{s}^t, \nabla f(\mathbf{z}^t)$ , respectively. It is interesting to examine a few special  $p$  values in (8) (ignoring the negative sign for simplicity):

- $p = 1$ ,  $\mathbf{s}^t$  selects the largest entry in  $\nabla f(\mathbf{z}^t)$ , while setting all other entries to zero<sup>2</sup>.
- $p = 2$ ,  $\mathbf{s}^t$  re-scales  $\nabla f(\mathbf{z}^t)$  by its root mean square.
- $p = \infty$ , all entries of  $\mathbf{s}^t$  have the equal magnitude  $c$ .

The three special cases easily remind the behaviors of *max pooling*, *root-mean-square pooling* and *average pooling* [21], [53], although the input and output are both  $\mathbb{R}^m$  and no dimensionality reduction is performed. For general  $p \in [1, \infty)$ , it is expected to exhibit more varied behaviors that can be interpretable from a pooling perspective. We thus denote by the function  $\mathbb{R}^m \rightarrow \mathbb{R}^m$ :  $\mathbf{s}^t = \text{pool}_p(\nabla f(\mathbf{z}^t))$ , to illustrate the operation (8) associated with a specific  $p$ .

### B. Constructing the Network

Following the well received unrolling-then-truncating methodology, e.g. [25], we represent the Frank-Wolfe solver for (6) that runs finite  $T$  iterations ( $t = 0, 1, \dots, T-1$ ), as a multi-layer feed-forward neural network. Fig. 1 depicts the resulting architecture, named the Frank-Wolfe Network (F-W Net). As a custom, we set  $\mathbf{z}^0 = \mathbf{0}$ , which is an interior point of  $\mathcal{Z}$  for any  $p \geq 1$ . The  $T$  layer-wise weights  $W_t$ 's can be analytically constructed. Specifically, note that  $\nabla f(\mathbf{z}^t) = D^\top(D\mathbf{z}^t - \mathbf{x}) = D^\top D\mathbf{z}^t - D^\top \mathbf{x}$ , if we define

$$W_0 = -D^\top; W_t = D^\top D, t = 1, 2, \dots, T-1. \quad (9)$$

Then we have  $\nabla f(\mathbf{z}^t) = W_t \mathbf{z}^t + W_0 \mathbf{x}$ , as depicted in Fig. 1. Given the pre-chosen hyper-parameters, and without any further tuning, F-W Net outputs a  $T$ -iteration approximation  $\mathbf{z}^T$  of the exact solution  $\mathbf{z}$  to (6). As marked out, the intermediate outputs of F-W Net in Fig. 1 are aligned with the

<sup>2</sup>This reminds us of the well-known matching pursuit algorithm. We noted that a recent work [62] has revealed a unified view between the Frank-Wolfe algorithm and the matching pursuit algorithm.

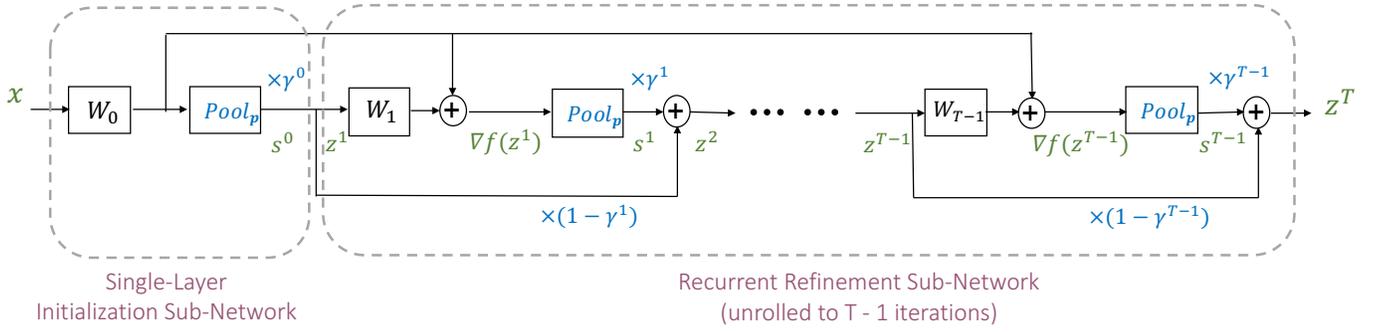


Fig. 1. (Please view in color)  $T$ -layer Frank-Wolfe network, consisting of two parts: (1) a single-layer initialization sub-network; and (2) a recurrent refinement sub-network, unrolled to  $T - 1$  iterations. We use the *Green* notations to remind to which variable in (7) the current intermediate output is corresponding. The layer-wise weights are denoted in *Black*, and the learnable units/hyper-parameters are in *Blue*.

iterative variables in the original Frank-Wolfe solver (7). Note that the two-variable update scheme in (7) naturally leads to two groups of “skip connections,” which might be reminiscent of the ResNet [63].

As many existing works did [25], [26], [35], the unrolled and truncated network in Fig. 1 could be alternatively treated as a trainable regressor to predict the exact  $\mathbf{z}$  from  $\mathbf{x}$ . We further view F-W Net as composed from two sub-networks: (1) a single-layer initialization sub-network, consisting of a linear layer  $W_0$  and a subsequent  $pool_p$  operator. It provides a rough estimate of the solution:  $\mathbf{z}_1 = \gamma^0 \mathbf{s}_0 = \gamma^0 pool_p(W_0 \mathbf{x})$ , which appears similar to typical sparse coding that often initializes  $\mathbf{z}$  with a thresholded  $D^T \mathbf{x}$  [64]. Note that  $\mathbf{z}_1$  has a direct shortcut path to the output, with the multiplicative weights  $\prod_{t=1}^{T-1} (1 - \gamma^t)$ ; (2) a recurrent refinement sub-network, that is unrolled to repeat the layer-wise transform for  $T - 1$  times to gradually refine the solution.

F-W Net is designed for large learning flexibility, by fitting almost all its parameters and hyper-parameters from training data:

- **Layer-wise weights.** Eq. (9) can be used to initialize the weights, but during training,  $W_0$  and  $W_t$  are untied with  $D$  and viewed as conventional fully-connected layers (without biases). All weights are learned with back-propagation from end to end. In some cases (e.g. Table III), we find that sharing the  $W_t$ 's (with  $t = 1, 2, \dots, T - 1$ ) is a choice worthy of consideration: it effectively reduces the actual number of parameters and makes the training to converge fast. But not sharing the weights is always helpful to improve the performance, as observed in our experiments. Thus, the weights are not shared unless otherwise noted. In addition, the relation between  $W_0$  and  $W_t$  (with  $t = 1, 2, \dots, T - 1$ ), as in (9), is not enforced during training.  $W_0$  is treated as a simple fully-connected layer without additional constraint.
- **Hyper-parameters.**  $p$  and  $\gamma^t$  were given or pre-chosen in (7). In F-W Net, they can be either pre-chosen or made learnable. For learnable hyper-parameters, we also compute the gradients w.r.t.  $p$  and  $\gamma^t$ , and update them via back-propagation. We adopt the same  $p$  throughout F-W Net as  $p$  implies the structural prior.  $\gamma^t$  is set to be

independent per layer. Learning  $p$  and  $\gamma^t$  also adaptively compensates for the truncation effect [52] of iterative algorithms. In addition, learning  $p$  gracefully solves the unknown  $p$  problem (2).

F-W Net can further be jointly tuned with a task-specific loss function, e.g. the softmax loss for classification, or the mean-squared-error loss and/or a semantic guidance loss [65] for denoising, in the form of an end-to-end network.

### C. Implementing the Network

1) *Reformulating  $pool_p$  as normalization plus neuron:* A closer inspection on the structure of the  $pool_p$  function (8) leads to a two-step decomposition (let  $\mathbf{u}^t = \nabla f(\mathbf{z}^t)$  for simplicity):

$$\mathbf{y}^t = \frac{\mathbf{u}^t}{\|\mathbf{u}^t\|_{\frac{p}{p-1}}} \quad // \text{ Step 1: } p\text{-conjugate normalization}$$

$$\mathbf{s}_i^t = -c \cdot \text{sign}(y_i^t) \cdot |y_i^t|^{\frac{1}{p-1}} \quad // \text{ Step 2: } p\text{-exponential neuron} \quad (10)$$

Step 1 performs a normalization step under the  $\frac{p}{p-1}$ -norm. Let  $q = \frac{p}{p-1}$ , which happens to be the Hölder conjugate of  $p$ : we thus call this step  $p$ -conjugate normalization. It coincides with a simplified form of DN [59], by setting all summation and suppression weights of DN to 1.

Step 2 takes the form of an exponential-type and non-saturated element-wise activation function [66], and is a learnable activation parameterized by  $p$  [50]. As the output range of Step 1 is  $[-1, 1]$ , the exponent displays suppression effect when  $p \in (1, 2)$ , and amplifies entries when  $p \in (2, \infty)$ .

While the decomposition of (8) is not unique, we carefully choose (10) due to its effectiveness and numerical stability. As a counterexample, if we adopt another more straightforward decomposition of (8):

$$y_i^t = \text{sign}(u_i^t) \cdot |u_i^t|^{\frac{1}{p-1}}; \quad \mathbf{s}^t = -\frac{c \cdot \mathbf{y}^t}{\|\mathbf{y}^t\|_p}, \quad (11)$$

then, large  $|u_i^t|$  values ( $> 1$ ) will be boosted by the power of  $\frac{1}{p-1}$  when  $p \in (1, 2)$ , and the second step may run the risk of explosion when  $p \rightarrow 1$ , in both feed-forward and back-propagation. In contrast, (10) first squashes each entry into

$[-1, 1]$  (Step 1) before feeding into the exponential neuron (Step 2), resolving the numerical issue well in practice.

The observation (10), called “*pool<sub>p</sub> = normalization + neuron*” for brevity, provides a new interesting insight into the connection between neuron, pooling and normalization, the three major types of nonlinear units in deep networks that were once considered separately. By splitting *pool<sub>p</sub>* into two modules sharing the parameter *p*, the back-propagation computation is also greatly simplified, as directly computing the gradient of *pool<sub>p</sub>* w.r.t. *p* can be quite involved, with more potential numerical problems.

2) *Network initialization and training*: The training of F-W Net evidently benefits from high-quality initializations. Although  $W_t$  and  $D$  are disentangled,  $W_t$  can be well initialized from the given or estimated  $D^3$  via (9).  $p$  is typically initialized with a large scalar, but its learning process is found to be quite insensitive to initialization. As observed in our experiments, no matter what  $p$  is initialized as, it will converge stably and smoothly to almost the same value<sup>4</sup>.  $\gamma^t$  is initialized with the rule  $\frac{2}{t+2}$ , and is re-parametrized using a sigmoid function to enforce  $[0, 1]$  range during training.  $c$  is the only hyper-parameter that needs to be manually chosen and fed into F-W Net. In experiments, we find that a good  $c$  choice could accelerate the convergence.

In the original  $p$ -NCLS,  $\|D\|_2 = 1$  is assumed to avoid the scale ambiguity, and is commonly enforced in dictionary learning [22]. Similarly, to restrain the magnitude growth of  $W$ , we adopt the  $L_2$ -norm weight decay regularizer, with the default coefficient  $2 \times 10^{-4}$  in experiments.

#### D. Interpretation of Frank-Wolfe Network as LSTM

The ISTA algorithm [40] has been interpreted as a stack of plain Recurrent Neural Networks (RNNs). We hereby show that F-W Net can be potentially viewed as a special case of the Long Short-Term Memory (LSTM) architecture [67], which incorporates a gating mechanism [68] and may thus capture long-term dependencies better than plain RNNs when  $T$  is large. Although we include no such experiment in this paper, we are interested in applying F-W Net as LSTM to model sequential data in future work.

Let us think  $\mathbf{x}$  as a constant *input*, and  $\mathbf{z}^t$  is the *previous hidden state* of the  $t$ -th step ( $t = 1, \dots, T - 1$ ). The current *candidate hidden state* is computed by  $\mathbf{s}^t = \text{pool}_p(W_0\mathbf{x} + W_t\mathbf{z}^t)$ , where the sophisticated *pool<sub>p</sub>* function replaces the common  $\tanh$  to be the activation function.  $\gamma_t$  and  $(1 - \gamma_t)$  each take the role of the *input gate* and *forget gate*, that control how much of the newly computed and previous state will go through, respectively.  $\gamma^t\mathbf{s}^t + (1 - \gamma^t)\mathbf{z}^t$  constitutes the *internal memory* of the current unit. Eventually, with a simple *output gate* equal to 1, the *new hidden state* is updated to  $\mathbf{z}^{t+1}$ .

<sup>3</sup>For example, when  $D$  is not given, we can estimate  $D$  using K-SVD [22] to initialize  $W$ .

<sup>4</sup>We are aware of the option to re-parameterize  $p$  to ensure  $p \geq 1$  [56]. We have not implemented it in our experiments, since we never encountered  $p < 1$  during learning. The re-parameterization trick can be done if necessary.

## IV. EVALUATION OF FRANK-WOLFE NETWORK

### A. Simulations

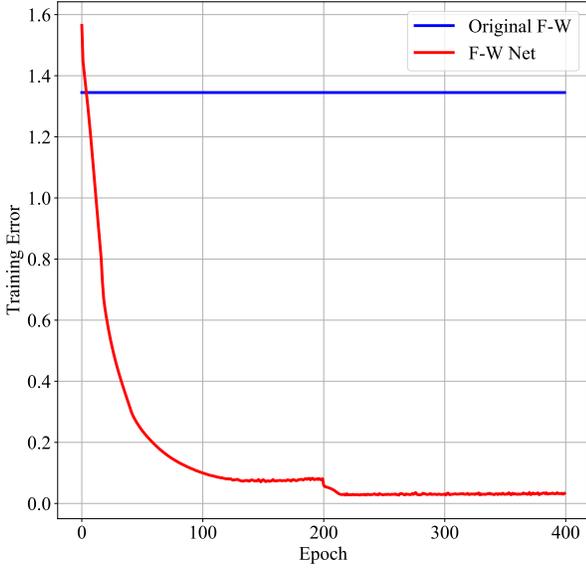
We generate synthetic data in the following steps. First we generate a random vector  $\mathbf{y} \in \mathbb{R}^m$  and then we project it to the  $L_p$ -ball of radius  $c$ . The projection is achieved by solving the problem:  $\mathbf{z} = \arg \min \frac{1}{2} \|\mathbf{y} - \mathbf{z}\|_2^2$  s.t.  $\|\mathbf{z}\|_p \leq c$ , using the original F-W algorithm until convergence. This  $p$  will be termed *real p* in the following. We then create a random matrix  $D \in \mathbb{R}^{n \times m}$ , and achieve  $\mathbf{x} := D\mathbf{z} + \mathbf{e}$ , where  $\mathbf{e} \in \mathbb{R}^n$  is additive white Gaussian noise with variance  $\sigma^2 = 0.01$ . We use the default values  $m = 100$ ,  $n = 50$ ,  $c = 5$ , and generate 15,000 samples for training. A testing set of 1,000 samples are generated separately in the identical manner. Our goal is to train a F-W Net to predict/regress  $\mathbf{z}$  from the given observation  $\mathbf{x}$ . The performance is measured by mean-squared-error (MSE) between predicted and ground-truth  $\mathbf{z}$ . All network models are implemented with Caffe [69] and trained by using MSE between ground-truth and network-output vectors as loss.

We first choose real  $p = 1.3$ , vary  $T = 2, 4, 6, 8$ , and compare the following methods:

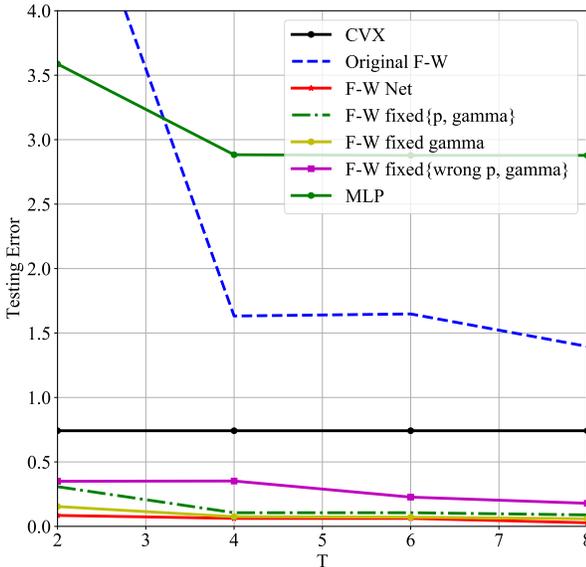
- *CVX*: solving the problem using the ground-truth  $D$  and real  $p$ . The CVX package [70] is employed to solve this convex problem. No training process is involved here.
- *Original F-W*: running the original Frank-Wolfe algorithm for  $T$  iterations, using the ground-truth  $D$  and real  $p$  and fixing  $\gamma^t = \frac{2}{t+2}$ . No training process is involved here.
- *MLP*: replacing the *pool<sub>p</sub>* in F-W Net with ReLU, having a feed-forward network of  $T$  fully-connected layers, which has the same number of parameters with F-W Net (except  $p$ ).
- *F-W Net*: the proposed network that jointly learns  $W_t$ ,  $p$ , and  $\gamma^t$ ,  $t = 0, 1, \dots, T - 1$ .
- *F-W fixed*  $\{p, \gamma\}$ : fixing  $p = 1.3$  and  $\gamma^t = \frac{2}{t+2}$ , learning  $W_t$  in F-W Net.
- *F-W fixed*  $\gamma$ : fixing  $\gamma^t = \frac{2}{t+2}$ , learning  $W_t$  and  $p$  in F-W Net.
- *F-W fixed*  $\{\text{wrong } p, \gamma\}$ : fixing  $p = 1$  (i.e. an incorrect structural prior) and  $\gamma^t = \frac{2}{t+2}$ , learning  $W_t$  in F-W Net.

Fig. 2 (a) depicts the training curve of F-W net at  $T = 8$ . We also run the *Original F-W* (for 8 iterations) for comparison. After a few epochs, F-W Net is capable in achieving significantly smaller error than the Original F-W, and converges stably. It shows the advantage of training a network to compensate for the truncation error caused by limited iterations. It is noteworthy that F-W Net has the ability to adjust the dictionary  $D$ , prior  $p$  and step size  $\gamma^t$  at the same time, all of which are learned from training data. Comparing the test error of F-W Net and Original F-W at different  $T$ 's (Fig. 2 (b)), F-W Net shows the superiority of flexibility, especially when  $T$  is small or even the  $p$  is wrong. F-W Net learns to find a group parameters ( $D$ ,  $p$  and  $\gamma^t$ ) during the training process, so that these parameters coordinate each other to lead to better performance.

Fig. 2 (b) compares the testing performance of different methods at different  $T$ 's. F-W Net with learnable parameters achieves the best performance. The *Original F-W* performs



(a) Training Error - Epoch



(b) Testing Error - T

Fig. 2. (a) The plot of average training error per sample w.r.t. the epochs at  $T = 8$ . Result of the *Original F-W* is also plotted; (b) Average test error per sample comparison on the testing set, at different  $T$ 's.

poorly at  $T = 2$ , and its error barely decreases as  $T$  increases to 4 and 8, since the original F-W algorithm is known to converge slowly in practice. CVX does not perform well as this problem is quite difficult ( $p = 1.3$ ). Though having the same number of parameters, MLP performs not well, which indicates that this synthetic task is not trivial. Especially, the original Frank-Wolfe algorithm also significantly outperforms MLP after 4 iterations.

Furthermore, we reveal the effect of each component ( $D$ ,  $p$  and  $\gamma^t$ ). Firstly, we fix  $p$  and  $\gamma^t$ , i.e. we learn the  $W_t$ 's compared to the original F-W algorithm. By observing the two curves of testing performance, F-W Net improves the performance with learnable  $W_t$ 's which coordinate the fixed  $p$  and  $\gamma^t$ , and achieves better approximation of the target  $\mathbf{z}$

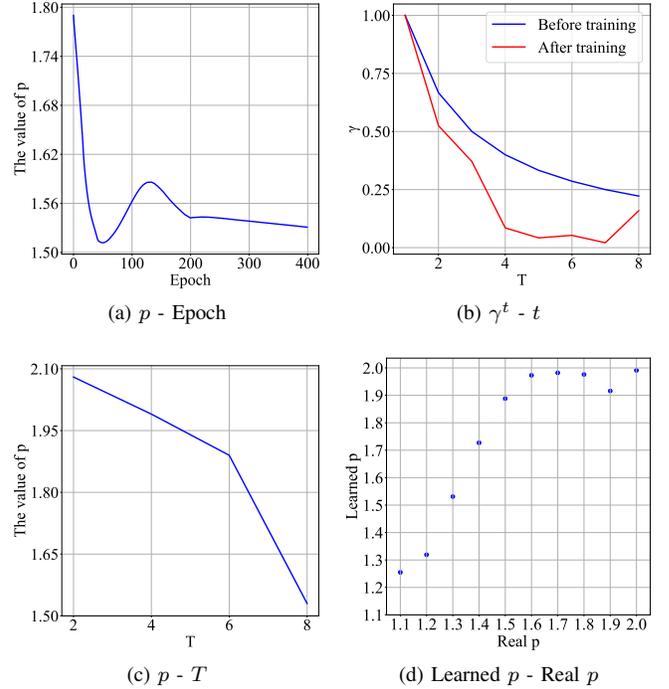


Fig. 3. (a) The plot of  $p$  during training w.r.t. the epochs at  $T = 8$  (real  $p = 1.3$ ); (b) The plot of  $\gamma^t$  values before and after training at  $T = 8$  (real  $p = 1.3$ ); (c) The plot of learned  $p$  at different  $T$ 's (real  $p = 1.3$ ); (d) The plot of learned  $p$  at  $T = 8$ , when the real  $p$  value is from 1.1 to 2.0.

in a few steps. Then, we let  $p$  be learnable and only fix  $\gamma^t$ , which is slightly superior to *F-W fixed*  $\{p, \gamma\}$ . F-W Net also maintains a smaller, but consistent margin over *F-W fixed*  $\gamma$ . Those three comparisons confirm that except for  $W_t$ , learning  $p$  and  $\gamma$  are both useful. Finally, we give the wrong  $p$  to measure the influence on F-W Net. It is noteworthy that *F-W fixed*  $\{wrong\ p, \gamma\}$  suffers from much larger error than other F-W networks. That demonstrates the huge damage that an incorrect or inaccurate structural prior can cause to the learning process. But, as mentioned before, F-W Net has the high flexibility to adjust the other parameters. Even though under the wrong  $p$  condition, F-W Net still outperforms the original F-W algorithm in a few steps.

Fig. 3 inspects the learning of  $p$  and  $\gamma^t$ . As seen from Fig. 3 (a), the  $p$  value fluctuates in the middle of training, but ends up converging stably (initial  $p = 2$  i.e.  $L_2$ -norm). In Fig. 3 (c), as  $T$  goes up, the learned  $p$  by F-W net approaches the real  $p = 1.3$  gradually. This phenomenon can be interpreted by that the original F-W algorithm cannot solve the problem well in only a few steps, and thus F-W Net adaptively controls each component through learning them from training data to approximate the distribution of  $\mathbf{z}$  as much as possible. To understand why the learned  $p$  is usually larger than the real  $p$ , we may intuitively think that  $pool_p$  will “compress” the input more heavily as  $p$  gets down closer to 1. To predict  $\mathbf{z}$ , while the original Frank-Wolfe algorithm may run many iterations, F-W Net has to achieve within a much smaller, fixed number of iterations. Thus, each  $pool_p$  has to let more energy “pass through,” and the learning result has larger  $p$ . Fig. 3 (b) observes the change of  $\gamma^t$  before and after training, at  $T = 8$ .

While  $\gamma^0$  remains almost unchanged,  $\gamma^t$  ( $t \geq 1$ ) all decreases. As a possible reason, F-W Net might try to compensate the truncation error by raising the weight of the initialization  $\mathbf{z}^1$  in the final output  $\mathbf{z}^T$ .

We then look into  $p = 1$  case, and re-generate synthetic data. We compare three methods as defined before: CVX, F-W Net, and F-W fixed  $p$ . In addition, we add LISTA [25] into comparison, because LISTA is designed for  $p = 1$ . The depth, layer-wise dimensions, and weight-sharing of LISTA are configured identically to F-W Net. Note that LISTA is dedicated to the  $L_1$  case and cannot be easily extended for general  $L_p$  cases. We re-tune all the parameters to get the best performance with LISTA to ensure a fair comparison. Table I compares their results at  $T = 2, 4, 6$ . The CVX is able to solve  $L_1$  problems to a much better accuracy than the case of  $p = 1.3$ , and F-W Net still outperforms F-W fixed  $p$ . More interesting is the comparison between F-W Net and LISTA: F-W Net is outperformed by LISTA at  $T = 2$ , then reaches a draw at  $T = 4$ , and eventually outperforms LISTA by a large margin at  $T = 6$ . Increasing  $T$  demonstrates a more substantial boost on the performance of F-W Net than that of LISTA, which can be interpreted as we have discussed in Section III-D that F-W Net is a special LSTM, but LISTA is a vanilla RNN [40]. Note that the real  $p$  is 1 (corresponding to sparse), but the learned  $p$  in F-W Net is larger than 1 (corresponding to non-sparse). The success of F-W Net does not imply that the problem itself is a non-sparse coding one. This is also true for all the following experiments.

We also simulate with other real  $p$  values. Fig. 3 (d) shows the learned  $p$  values with respect to different real  $p$ 's. When the real  $p$  approaches 1 or 2, F-W Net is able to estimate the value of  $p$  more accurately, probably because the convex problems with  $L_1$ -norm and  $L_2$ -norm minimization can be solved more easily.

### B. Handwritten Digit Recognition

Similar to what has been done in [25], we adopt F-W Net as a feature extractor and then use logistic regression to classify the features for the task of handwritten digit recognition. We use the MNIST dataset [71] to experiment. We design the following procedure to pre-train the F-W Net as a feature extractor. The original images are dimensionality reduced to 128-dim by PCA for input to F-W Net. Then we further perform PCA on each digit separately to reduce the dimension to 15. We construct a 150-dim sparse code for each image, whose 150 dimensions are divided into 10 groups to correspond to 10 digits, only 15 dimensions of which are filled by the corresponding PCA result, whereas the other dimensions are all zero. This sparse code is regarded as the ‘‘ground-truth’’ for F-W Net in training. Accordingly, the transformation matrices in the second-step PCA are concatenated to serve as  $D \in \mathbb{R}^{128 \times 150}$ , which is used to initialize the fully-connected layers in F-W Net according to (9). In this experiment, we use stochastic gradient descent (SGD), a momentum of 0.9 and a mini-batch size of 100. The F-W Net is pre-trained for 200 epochs, and the initial learning rate is 0.1, decayed to 0.01 at 100 epochs. Then the pre-trained F-W Net is augmented

with a fully-connected layer with softmax that is randomly initialized, resulting a network that can be trained end-to-end for classification. We observe that the performance benefits from joint training marginally but consistently.

If we formulate this feature extraction problem as  $L_p$ -norm constrained, then the  $p$  is unknown, and most of previous methods adopt  $L_1$ -norm or  $L_2$ -norm. Different from them, F-W Net tries to attain a  $p$  from training data, which suits the real data better. The results are shown in Table II, comparing F-W Net with simple feed-forward fully-connected networks (MLP) [71] and LCoD [25]. F-W Net achieves lower error rate than the others. Especially, it takes 50 iterations for LCoD to achieve an error rate of 1.39, but only 3 layers for F-W Net to achieve 1.34, where the numbers of parameters are similar. Moreover, with the increasing number of layers, F-W Net makes a continuous improvement, which is consistent with the observation in the simulation.

Table III provides the results of different initializations of the hyper-parameter  $p$ , showing that F-W Net is insensitive to the initialization of  $p$  and can converge to the learned  $p \approx 1.6$ . However, fixing the  $p$  value is not good for F-W Net even fixing to the finally learned  $p = 1.6$ . This is interesting as it seems the learnable  $p$  provides advantages not only for the final trained model but also for training itself. An adjustable  $p$  value may suit for the evolving parameters during training F-W Net, which we plan to study further.

## V. CONVOLUTIONAL FRANK-WOLFE NETWORK

Previous similar works [25], [26], [35] mostly result in fully-connected networks, as they are unrolled and truncated from linear sparse coding models. Nonetheless, fully-connected networks are less effective than convolutional neural networks (CNNs) when tackling structured multi-dimensional signal such as images. A natural idea to extend this type of works to convolutional cases seems to be convolutional sparse coding [72], which also admits iterative solvers. However, the re-formulation will be inefficient both memory-wise and computation-wise, as discussed in [73].

We therefore seek a simpler procedure to build the convolutional F-W Net: in Fig. 1, we replace the fully-connected layers with convolutional layers, and operate  $pool_p$  across all the output feature maps for each location individually. The latter is inspired by the well-known conversion between convolutional and fully-connected layers by reshaping inputs<sup>5</sup>, and turns  $pool_p$  into a form of cross-channel parametric pooling [74]–[76].

The convolutional F-W Net then bears the similar flexibility to jointly learn weights and hyper-parameters ( $p$  and  $\gamma^t$ ). Yet different from the original version of F-W Net, the  $pool_p$  in convolutional F-W Net reflects a diversified treatment of different convolutional filters at the same location, and should be differentiated from pooling over multiple points:

- $p = 1$ , only the channel containing the strongest response will be preserved at each location, which is reduced to max-out [74].

<sup>5</sup><http://cs231n.github.io/convolutional-networks/#convert>

TABLE I  
AVERAGE ERROR PER TESTING SAMPLE AT DIFFERENT  $T$ 'S. THE REAL  $p$  IS 1.

|     | CVX       | F-W Net |        |        | F-W Net with fixed $p$ |           |           | LISTA     |           |           |
|-----|-----------|---------|--------|--------|------------------------|-----------|-----------|-----------|-----------|-----------|
|     |           | $T=2$   | $T=4$  | $T=6$  | $T=2$                  | $T=4$     | $T=6$     | $T=2$     | $T=4$     | $T=6$     |
| MSE | 0.0036    | 0.5641  | 0.3480 | 0.1961 | 1.2076                 | 0.8604    | 0.7358    | 0.4157    | 0.3481    | 0.3053    |
| $p$ | 1 (fixed) | 1.360   | 1.254  | 1.222  | 1 (fixed)              | 1 (fixed) | 1 (fixed) | 1 (fixed) | 1 (fixed) | 1 (fixed) |

TABLE II  
RESULTS ON THE MNIST TEST SET.

|                       | # Params | Error rate (%) |
|-----------------------|----------|----------------|
| 3-layer: 300+100 [71] | 266,200  | 3.05           |
| 3-layer: 500+150      | 468,500  | 2.95           |
| 3-layer: 500+300      | 545,000  | 1.53           |
| 1-iter LCoD [25]      | 65,536   | 1.60           |
| 5-iter LCoD           | 65,536   | 1.47           |
| 50-iter LCoD          | 65,536   | 1.39           |
| 2-layer F-W Net       | 43,350   | 2.20           |
| 3-layer F-W Net       | 65,850   | 1.34           |
| 4-layer F-W Net       | 88,200   | <b>1.25</b>    |

TABLE III  
RESULTS ON THE MNIST TEST SET OF F-W NETS WITH FIXED  $p$  AND LEARNABLE  $p$ . THE  $W_t$  PARAMETERS ARE SHARED ACROSS  $t = 1, \dots, T - 1$  IN F-W NET.

| Initial $p$ | Learnable $p$ |       | Fixed $p$    |     |
|-------------|---------------|-------|--------------|-----|
|             | Accuracy (%)  | $p$   | Accuracy (%) | $p$ |
| 1.1         | 98.51         | 1.601 | 97.35        | 1.1 |
| 1.3         | 98.45         | 1.600 | 97.70        | 1.3 |
| 1.5         | 98.43         | 1.595 | 98.30        | 1.5 |
| 1.6         | 98.66         | 1.614 | 98.26        | 1.6 |
| 1.8         | 98.38         | 1.601 | 97.90        | 1.8 |
| 2.0         | 98.50         | 1.585 | 97.78        | 2.0 |

- $p = 2$ ,  $pool_p$  re-scales the feature maps across the channels by its root-mean-square.
- $p = \infty$  denotes the equal importance of all channels and leads to the cross-channel average.

By involving  $p$  into optimization,  $pool_p$  essentially learns to re-scale local responses, as advocated by the neural lateral inhibition mechanism. The learned  $p$  will indicate the relative importance of different channels, and can potentially be useful for network compression [77].

The convolutional kernels in the convolutional F-W Net are not directly tied with any dictionary, thus we have not seen an initialization strategy straightforwardly available. In this paper, we use random initialization for the convolutional filters, but we initialize  $p$  and  $\gamma^t$  in the same way as mentioned before. As a result, convolutional F-W Nets often converge slower than fully-connected ones. We notice some recent works that construct convolutional filters [78], [79] from data in an explicit unsupervised manner, and plan to exploit them as future options to initialize convolutional F-W Nets.

To demonstrate its capability, we apply convolutional F-W Net to low-level image processing tasks, including image denoising and image super-resolution in this paper. Our focus lies on building light-weight compact models and verifying the benefit of introducing the learnable  $p$ .

## VI. EXPERIMENTS OF CONVOLUTIONAL FRANK-WOLFE NETWORK

### A. Image Denoising

We investigate two versions of F-W Net for image denoising. The first version is the fully-connected (FC) F-W Net as used in the previous simulation. Here, the basic FC F-W Net (Fig. 1) is augmented with one extra fully-connected layer, whose parameters are denoted by  $W_R \in \mathbb{R}^{n \times m}$  to reconstruct  $\hat{\mathbf{x}} = W_R \mathbf{z}$ .  $W_R$  is naturally initialized by  $D$ . The network output is compared against the original/clean image to calculate MSE loss. To train this FC F-W Net, note that one strategy to image denoising is to split a noisy image into small (like  $8 \times 8$ ) overlapping patches, process the patches individually, and compose the patches back to a complete image. Then, we re-shape  $8 \times 8$  blocks into 64-dim samples (i.e.  $n = 64$ ). We adopt  $m = 512$ ,  $T = 8$ , and  $c = 1$ . The network is trained with a number of  $8 \times 8$  noisy blocks as well as their noise-free counterparts. The second version of F-W Net is the proposed convolutional (Conv) F-W Net as discussed in Section V. The adopted configurations are:  $3 \times 3$  filters, 64 feature maps per layer,  $T = 4$ , and  $c = 1$ .

We use the BSD-500 dataset [80] for experiment. The BSD-68 dataset is used as testing data, and the remaining 432 images are converted to grayscale and added with white Gaussian noise  $\mathcal{N}(0, \sigma^2)$  for training. Several competitive baselines are chosen: a FC LISTA [25] network that is configured identically to our FC F-W Net; a Conv LISTA network that is configured identically to our Conv F-W Net; KSVD + OMP that represents the  $L_0$ -norm optimization; BM3D [81] using the dictionary size of 512; and DnCNN [82] which is a recently developed method based on CNN, here our re-trained DnCNN includes 4 convolutional layers followed by BN and ReLU. We consider three noise levels:  $\sigma = 15$ ,  $\sigma = 25$ , and  $\sigma = 50$ .

Table IV provides the results of different image denoising methods on the BSD-68 dataset. FC F-W Net is better than LISTA in all cases. We also study the effect of  $p$ , as shown in Table V. F-W Net with learnable  $p$  outperforms F-W Net with fixed  $p = 2$  by a large margin, and is slightly better than F-W Net with fixed  $p = 1.4$ . Here, learnable  $p$  seems benefiting not only the final model but also the training itself, as has been observed in IV-B. BM3D is a strong baseline, which the deep networks cannot beat easily. As seen from Table IV, BM3D outperforms DnCNN (4 layers), but Conv F-W Net (4 layers) is better than BM3D.

It is worth noting that the original DnCNN in [82] has 20 layers and much more parameters and outperforms BM3D. We conduct an experiment to compare DnCNN with Conv F-W Net at different number of layers. The results are shown in Fig. 4 (a). We observe that for shallow networks, Conv F-W Net outperforms DnCNN significantly. As the network

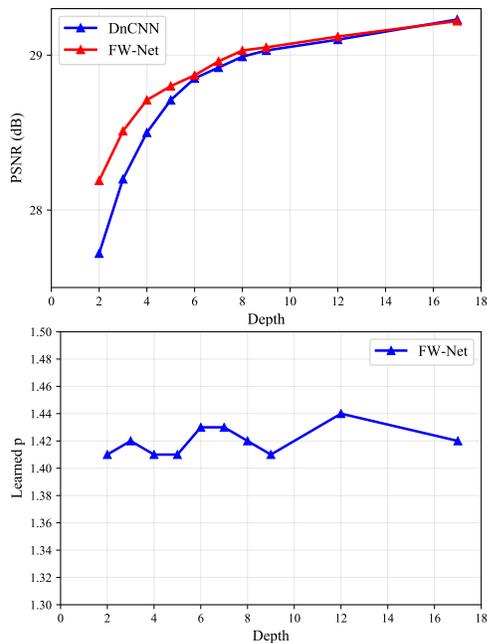


Fig. 4. Top: image denoising results on the (gray) BSD-68 dataset when  $\sigma = 25$ , with DnCNN and (Conv) F-W Net at different depths. Bottom: learned  $p$  value of the Conv F-W Net at different depths.

TABLE IV  
IMAGE DENOISING RESULTS (AVERAGE PSNR/DB) ON THE (GRAY) BSD-68 DATASET. ALL THE NETWORKS HAVE 4 FULLY-CONNECTED OR CONVOLUTIONAL LAYERS FOR FAIR COMPARISON.

|                                | $\sigma = 15$ | $\sigma = 25$ | $\sigma = 50$ |
|--------------------------------|---------------|---------------|---------------|
| FC LISTA                       | 29.20         | 27.63         | 24.33         |
| FC F-W Net, learned $p = 1.41$ | 29.35         | 27.71         | 24.52         |
| KSVD + OMP                     | 30.82         | 27.97         | 23.97         |
| BM3D                           | 31.07         | 28.57         | 25.62         |
| DnCNN (4 layers)               | 30.89         | 28.42         | 25.42         |
| DnCNN (4 layers, w/o BN)       | 30.85         | 28.43         | 25.36         |
| Conv LISTA (4 layers)          | 30.90         | 28.50         | 25.55         |
| Conv F-W Net (4 layers)        | <b>31.27</b>  | <b>28.71</b>  | <b>25.66</b>  |

depth increases, both Conv F-W Net and DnCNN perform similarly. This may be attributed to the learnable parameter  $p$  in the Conv F-W Net, which helps much when the network parameters are less. Thus,  $pool_p$  may be favorable if we want to build a compact network. The learned  $p$  values are shown in Fig. 4 (b). We observe that the learned  $p$  value is stable across networks with different depths. It is also similar to the learned  $p$  value in the FC F-W Net ( $p = 1.41$ ). Thus, we consider that the  $p$  value is determined by the data, and F-W Net can effectively identify the  $p$  value regardless of FC or Conv structures.

TABLE V  
RESULTS OF FC F-W NETS WITH FIXED  $p$  AND LEARNABLE  $p$  FOR IMAGE DENOISING ON THE BSD-68 DATASET.

|                             | $\sigma = 15$ | $\sigma = 25$ | $\sigma = 50$ |
|-----------------------------|---------------|---------------|---------------|
| F-W Net, fixed $p = 2$      | 28.35         | 26.41         | 23.68         |
| F-W Net, fixed $p = 1.4$    | 29.25         | 27.62         | 24.45         |
| F-W Net, learned $p = 1.41$ | <b>29.35</b>  | <b>27.71</b>  | <b>24.52</b>  |

TABLE VI  
 $3\times$  IMAGE SUPER-RESOLUTION RESULTS (AVERAGE PSNR/DB) ON THE (GRAY) SET-5 AND SET-14 DATASETS.

| Method               | Set-5        | Set-14       |
|----------------------|--------------|--------------|
| 3-layer SRCNN        | 32.34        | 28.64        |
| 4-layer VDSR         | 32.51        | 28.71        |
| 4-layer VDSR (+BN)   | 32.55        | 28.69        |
| 4-layer Conv F-W Net | <b>32.85</b> | <b>28.76</b> |

TABLE VII  
ABLATION STUDY RESULTS OF CONV F-W NET FOR IMAGE SUPER-RESOLUTION ON THE SET-5 AND SET-14 DATASETS.

| Method                         | Set-5        | Set-14       |
|--------------------------------|--------------|--------------|
| F-W Net (No Skip)              | 31.52        | -            |
| F-W Net (Fixed $\gamma = 1$ )  | 32.21        | -            |
| F-W Net (ReLU)                 | 32.57        | 28.66        |
| F-W Net                        | <b>32.85</b> | <b>28.76</b> |
| F-W Net (Fixed $p = 1.3$ )     | 32.49        | 28.42        |
| F-W Net (Fixed $p = 1.5$ )     | 32.81        | 28.63        |
| F-W Net (Fixed $p = 1.8$ )     | 32.73        | 28.69        |
| F-W Net (Fixed $p = 2.3$ )     | 32.59        | 28.58        |
| F-W Net (Fixed $p = 2.5$ )     | 32.65        | 28.67        |
| F-W Net (Learned $p = 1.489$ ) | <b>32.85</b> | <b>28.76</b> |

## B. Image Super-Resolution

For the experiments about single image super-resolution (SR), we compare Conv F-W Net with baselines SRCNN [83] and VDSR [84], and all methods are trained on the 91-image standard set and evaluated on Set-5/Set-14 with a scaling factor of 3. We train a 4-layer Conv F-W Net and a 4-layer VDSR (4 convolutional layers equipped with ReLU), both of which have the same number of convolutional kernels. We adopt the same training configurations as presented in [84]. For SRCNN we directly use the model released by the authors.

Table VI compares the results of these methods. Our Conv F-W Net achieves the best performance among the three. To further understand the influence of each component in F-W Net, we experimentally compare the following settings:

- F-W Net (No Skip): removing the top skip connections in Fig. 1;
- F-W Net (Fixed  $\gamma = 1$ ): setting  $\gamma^t = 1, t = 1, 2, \dots, T-1$ , which is equivalent to removing the bottom skip connections in Fig. 1;
- F-W Net (ReLU): replacing all the  $pool_p$  units with ReLU.

Both *F-W Net (No Skip)* and *F-W Net (Fixed  $\gamma = 1$ )* break the original structure introduced by the Frank-Wolfe algorithm, and incur severe performance drop. *F-W Net (ReLU)* performs similarly to 4-layer VDSR and is worse than F-W Net. These results further demonstrate that each component in Conv F-W Net contributes to the final performance.

We also measure the effect of the hyper-parameter  $p$ . The results are shown in Table VII. Different  $p$  values indeed influence on the final performance significantly, and F-W Net with learnable  $p$  achieves the best performance, which again verifies the advantage of attaining the prior from the data.

## VII. CONCLUSION

We have studied the general non-sparse coding problem, i.e. the  $L_p$ -norm constrained coding problem with general

$p > 1$ . We have proposed the Frank-Wolfe network, whose architecture is carefully designed by referring to the Frank-Wolfe algorithm. Many aspects of F-W Net are inherently connected to the existing success of deep learning. F-W Net has gained impressive effectiveness, flexibility, and robustness in our conducted simulation and experiments. Results show that learning the hyper-parameter  $p$  is beneficial especially in real-data experiments, which highlights the necessity of introducing general  $L_p$ -norm and the advantage of F-W Net in learning the  $p$  during the end-to-end training.

Since the original Frank-Wolfe algorithm deals with convex optimization only, the proposed F-W Net can handle  $p \geq 1$  cases, but not  $p < 1$  cases. Thus, F-W Net is good at solving non-sparse coding problems.  $p = 1$  is quite special, as it usually leads to sparse solution [2], thus, F-W Net with fixed  $p = 1$  can solve sparse coding, too, but then its efficiency seems inferior to LISTA as observed in our experiments. For a real-world problem, is sparse coding or non-sparse coding better? This is an open problem and calls for future research. In addition, a number of promising directions have emerged as our future work, including handling more constraints other than the  $L_p$ -norm, and the customization of F-W Net for more real-world applications.

## REFERENCES

- [1] B. K. Natarajan, "Sparse approximate solutions to linear systems," *SIAM Journal on Computing*, vol. 24, no. 2, pp. 227–234, 1995.
- [2] D. L. Donoho, "For most large underdetermined systems of linear equations the minimal  $\ell_1$ -norm solution is also the sparsest solution," *Communications on Pure and Applied Mathematics*, vol. 59, no. 6, pp. 797–829, 2006.
- [3] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 267–288, 1996.
- [4] S. S. Chen, D. L. Donoho, and M. A. Saunders, "Atomic decomposition by basis pursuit," *SIAM Review*, vol. 43, no. 1, pp. 129–159, 2001.
- [5] D. L. Donoho, "Compressed sensing," *IEEE Transactions on Information Theory*, vol. 52, no. 4, pp. 1289–1306, 2006.
- [6] A. Krogh and J. A. Hertz, "A simple weight decay can improve generalization," in *NIPS*, 1992, pp. 950–957.
- [7] C. Studer, T. Goldstein, W. Yin, and R. G. Baraniuk, "Democratic representations," *arXiv preprint arXiv:1401.3420*, 2014.
- [8] J.-J. Fuchs, "Spread representations," in *ASILOMAR*, 2011, pp. 814–817.
- [9] Y. Lyubarskii and R. Vershynin, "Uncertainty principles and vector quantization," *IEEE Transactions on Information Theory*, vol. 56, no. 7, pp. 3491–3501, 2010.
- [10] Z. Wang, J. Liu, S. Huang, X. Wang, and S. Chang, "Transformed anti-sparse learning for unsupervised hashing," in *BMVC*, 2017, pp. 1–12.
- [11] R. Chartrand, "Exact reconstruction of sparse signals via nonconvex minimization," *IEEE Signal Processing Letters*, vol. 14, no. 10, pp. 707–710, 2007.
- [12] R. Chartrand and V. Staneva, "Restricted isometry properties and nonconvex compressive sensing," *Inverse Problems*, vol. 24, no. 3, p. 035020, 2008.
- [13] Z. Xu, X. Chang, F. Xu, and H. Zhang, " $L_{1/2}$  regularization: A thresholding representation theory and a fast solver," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 23, no. 7, pp. 1013–1027, 2012.
- [14] D. Krishnan and R. Fergus, "Fast image deconvolution using hyper-Laplacian priors," in *NIPS*, 2009, pp. 1033–1041.
- [15] M. Kloft, U. Brefeld, S. Sonnenburg, and A. Zien, " $l_p$ -norm multiple kernel learning," *Journal of Machine Learning Research*, vol. 12, pp. 953–997, 2011.
- [16] B. A. Olshausen and D. J. Field, "Sparse coding with an overcomplete basis set: A strategy employed by V1?" *Vision Research*, vol. 37, no. 23, pp. 3311–3325, 1997.
- [17] H. Lee, C. Ekanadham, and A. Y. Ng, "Sparse deep belief net model for visual area V2," in *NIPS*, 2008, pp. 873–880.
- [18] M. Elad and M. Aharon, "Image denoising via learned dictionaries and sparse representation," in *CVPR*, vol. 1, 2006, pp. 895–900.
- [19] J. Mairal, M. Elad, and G. Sapiro, "Sparse representation for color image restoration," *IEEE Transactions on Image Processing*, vol. 17, no. 1, pp. 53–69, 2008.
- [20] M. Ranzato, F.-J. Huang, Y.-L. Boureau, and Y. LeCun, "Unsupervised learning of invariant feature hierarchies with applications to object recognition," in *CVPR*, 2007, pp. 1–8.
- [21] J. Yang, K. Yu, Y. Gong, and T. Huang, "Linear spatial pyramid matching using sparse coding for image classification," in *CVPR*, 2009, pp. 1794–1801.
- [22] M. Aharon, M. Elad, and A. Bruckstein, "K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation," *IEEE Transactions on Signal Processing*, vol. 54, no. 11, pp. 4311–4322, 2006.
- [23] H. Xu, Z. Wang, H. Yang, D. Liu, and J. Liu, "Learning simple thresholded features with sparse support recovery," *IEEE Transactions on Circuits and Systems for Video Technology*, DOI: 10.1109/TCSVT.2019.2901713, 2019.
- [24] N. Bansal, X. Chen, and Z. Wang, "Can we gain more from orthogonality regularizations in training deep networks?" in *NIPS*, 2018, pp. 4261–4271.
- [25] K. Gregor and Y. LeCun, "Learning fast approximations of sparse coding," in *ICML*, 2010, pp. 399–406.
- [26] Z. Wang, Q. Ling, and T. S. Huang, "Learning deep  $\ell_0$  encoders," in *AAAI*, 2016, pp. 2194–2200.
- [27] Z. Wang, Y. Yang, S. Chang, Q. Ling, and T. S. Huang, "Learning a deep  $\ell_\infty$  encoder for hashing," in *IJCAI*, 2016, pp. 2174–2180.
- [28] S. Ross, D. Munoz, M. Hebert, and J. A. Bagnell, "Learning message-passing inference machines for structured prediction," in *CVPR*, 2011, pp. 2737–2744.
- [29] S. Wang, S. Fidler, and R. Urtasun, "Proximal deep structured models," in *NIPS*, 2016, pp. 865–873.
- [30] J. Sun, H. Li, and Z. Xu, "Deep ADMM-net for compressive sensing MRI," in *NIPS*, 2016, pp. 10–18.
- [31] M. Frank and P. Wolfe, "An algorithm for quadratic programming," *Naval Research Logistics*, vol. 3, no. 1-2, pp. 95–110, 1956.
- [32] J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry, and Y. Ma, "Robust face recognition via sparse representation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 2, pp. 210–227, 2009.
- [33] L.-C. Chen, A. Schwing, A. Yuille, and R. Urtasun, "Learning deep structured models," in *ICML*, 2015, pp. 1785–1794.
- [34] A. G. Schwing and R. Urtasun, "Fully connected deep structured networks," *arXiv preprint arXiv:1503.02351*, 2015.
- [35] P. Sprechmann, A. M. Bronstein, and G. Sapiro, "Learning efficient sparse and low rank models," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 9, pp. 1821–1833, 2015.
- [36] Z. Wang, S. Chang, J. Zhou, M. Wang, and T. S. Huang, "Learning a task-specific deep architecture for clustering," in *SDM*. SIAM, 2016, pp. 369–377.
- [37] Z. Wang, D. Liu, S. Chang, Q. Ling, Y. Yang, and T. S. Huang, "D3: Deep dual-domain based fast restoration of JPEG-compressed images," in *CVPR*, 2016, pp. 2764–2772.
- [38] Z. Wang, S. Huang, J. Zhou, and T. S. Huang, "Doubly sparsifying network," in *IJCAI*. AAAI Press, 2017, pp. 3020–3026.
- [39] J. Zhang and B. Ghanem, "ISTA-Net: Interpretable optimization-inspired deep network for image compressive sensing," in *CVPR*, 2018, pp. 1828–1837.
- [40] S. Wisdom, T. Powers, J. Pitton, and L. Atlas, "Interpretable recurrent neural networks using sequential sparse recovery," *arXiv preprint arXiv:1611.07252*, 2016.
- [41] B. Xin, Y. Wang, W. Gao, D. Wipf, and B. Wang, "Maximal sparsity with deep networks?" in *NIPS*, 2016, pp. 4340–4348.
- [42] T. Moreau and J. Bruna, "Understanding trainable sparse coding via matrix factorization," *arXiv preprint arXiv:1609.00285*, 2016.
- [43] X. Chen, J. Liu, Z. Wang, and W. Yin, "Theoretical linear convergence of unfolded ISTA and its practical weights and thresholds," in *NIPS*, 2018, pp. 9061–9071.
- [44] J. Liu, X. Chen, Z. Wang, and W. Yin, "ALISTA: Analytic weights are as good as learned weights in LISTA," *ICLR 2019*, <https://openreview.net/forum?id=B1lnzn0ctQ>, 2019.
- [45] U. Schmidt and S. Roth, "Shrinkage fields for effective image restoration," in *CVPR*, 2014, pp. 2774–2781.
- [46] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. Torr, "Conditional random fields as recurrent neural networks," in *ICCV*, 2015, pp. 1529–1537.

- [47] V. Chandrasekaran, B. Recht, P. A. Parrilo, and A. S. Willsky, "The convex geometry of linear inverse problems," *Foundations of Computational Mathematics*, vol. 12, no. 6, pp. 805–849, 2012.
- [48] M. Jaggi, "Revisiting Frank-Wolfe: Projection-free sparse convex optimization," in *ICML*, 2013, pp. 427–435.
- [49] L. Zhang, G. Wang, D. Romero, and G. B. Giannakis, "Randomized block Frank-Wolfe for convergent large-scale learning," *IEEE Transactions on Signal Processing*, vol. 65, no. 24, pp. 6448–6461, 2017.
- [50] F. Agostinelli, M. Hoffman, P. Sadowski, and P. Baldi, "Learning activation functions to improve deep neural networks," *arXiv preprint arXiv:1412.6830*, 2014.
- [51] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *ICCV*, 2015, pp. 1026–1034.
- [52] U. S. Kamilov and H. Mansour, "Learning optimal nonlinearities for iterative thresholding algorithms," *IEEE Signal Processing Letters*, vol. 23, no. 5, pp. 747–751, 2016.
- [53] K. Jarrett, K. Kavukcuoglu, and Y. LeCun, "What is the best multi-stage architecture for object recognition?" in *ICCV*. IEEE, 2009, pp. 2146–2153.
- [54] M. Malinowski and M. Fritz, "Learnable pooling regions for image classification," *arXiv preprint arXiv:1301.3516*, 2013.
- [55] C.-Y. Lee, P. W. Gallagher, and Z. Tu, "Generalizing pooling functions in convolutional neural networks: Mixed, gated, and tree," in *Artificial Intelligence and Statistics*, 2016, pp. 464–472.
- [56] C. Gulcehre, K. Cho, R. Pascanu, and Y. Bengio, "Learned-norm pooling for deep feedforward and recurrent neural networks," in *ECML-PKDD*. Springer, 2014, pp. 530–546.
- [57] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *ICML*, 2015, pp. 448–456.
- [58] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv preprint arXiv:1607.06450*, 2016.
- [59] M. Ren, R. Liao, R. Urtasun, F. H. Sinz, and R. S. Zemel, "Normalizing the normalizers: Comparing and extending network normalization schemes," *arXiv preprint arXiv:1611.04520*, 2016.
- [60] S. Lyu, "Divisive normalization: Justification and effectiveness as efficient coding transform," in *NIPS*, 2010, pp. 1522–1530.
- [61] J. Ballé, V. Laparra, and E. P. Simoncelli, "Density modeling of images using a generalized normalization transformation," *arXiv preprint arXiv:1511.06281*, 2015.
- [62] F. Locatello, R. Khanna, M. Tschannen, and M. Jaggi, "A unified optimization view on generalized matching pursuit and Frank-Wolfe," *arXiv preprint arXiv:1702.06457*, 2017.
- [63] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.
- [64] Z. Wang, J. Yang, H. Zhang, Z. Wang, Y. Yang, D. Liu, and T. S. Huang, *Sparse coding and its applications in computer vision*. World Scientific, 2016.
- [65] D. Liu, B. Wen, X. Liu, Z. Wang, and T. S. Huang, "When image denoising meets high-level vision tasks: A deep learning approach," in *IJCAI*. AAAI Press, 2018, pp. 842–848.
- [66] D.-A. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and accurate deep network learning by exponential linear units (ELUs)," *arXiv preprint arXiv:1511.07289*, 2015.
- [67] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [68] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *arXiv preprint arXiv:1412.3555*, 2014.
- [69] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *ACM Multimedia*. ACM, 2014, pp. 675–678.
- [70] M. Grant and S. Boyd, "CVX: Matlab software for disciplined convex programming," <http://cvxr.com/cvx>, Mar. 2014.
- [71] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [72] H. Bristow, A. Eriksson, and S. Lucey, "Fast convolutional sparse coding," in *CVPR*. IEEE, 2013, pp. 391–398.
- [73] H. Sreter and R. Gyryes, "Learned convolutional sparse coding," in *ICASSP*, 2018, pp. 2191–2195.
- [74] I. J. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio, "Maxout networks," *arXiv preprint arXiv:1302.4389*, 2013.
- [75] M. Lin, Q. Chen, and S. Yan, "Network in network," *arXiv preprint arXiv:1312.4400*, 2013.
- [76] T. Pfister, J. Charles, and A. Zisserman, "Flowing convnets for human pose estimation in videos," in *ICCV*, 2015, pp. 1913–1921.
- [77] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding," *arXiv preprint arXiv:1510.00149*, 2015.
- [78] T.-H. Chan, K. Jia, S. Gao, J. Lu, Z. Zeng, and Y. Ma, "PCANet: A simple deep learning baseline for image classification?" *IEEE Transactions on Image Processing*, vol. 24, no. 12, pp. 5017–5032, 2015.
- [79] J. Bruna and S. Mallat, "Invariant scattering convolution networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1872–1886, 2013.
- [80] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik, "Contour detection and hierarchical image segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, no. 5, pp. 898–916, 2011.
- [81] H. C. Burger, C. J. Schuler, and S. Harmeling, "Image denoising: Can plain neural networks compete with BM3D?" in *CVPR*. IEEE, 2012, pp. 2392–2399.
- [82] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a Gaussian denoiser: Residual learning of deep CNN for image denoising," *IEEE Transactions on Image Processing*, vol. 26, no. 7, pp. 3142–3155, 2017.
- [83] C. Dong, C. C. Loy, K. He, and X. Tang, "Image super-resolution using deep convolutional networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 2, pp. 295–307, 2016.
- [84] J. Kim, J. K. Lee, and K. M. Lee, "Accurate image super-resolution using very deep convolutional networks," in *CVPR*, 2016, pp. 1646–1654.