

Learned Extragradient ISTA with Interpretable Residual Structures for Sparse Coding

Paper ID: 8787

Abstract

Recently, the study on learned iterative shrinkage thresholding algorithm (LISTA) has attracted increasing attentions. A large number of experiments as well as some theories have proved the high efficiency of LISTA for solving sparse coding problems. However, existing LISTA methods are all serial connection. To address this issue, we propose a novel extragradient based LISTA (ELISTA), which has a residual structure and theoretical guarantees. Moreover, most LISTA methods use the soft thresholding function, which has been found to cause large estimation bias. Therefore, we propose a thresholding function for ELISTA instead of soft thresholding. In the theoretical aspect, we prove that our method attains linear convergence. In addition, through ablation experiments, the improvements of our method on the network structure and the thresholding function are verified. Extensive empirical results verify the advantages of our method.

1 Introduction

In this paper, we mainly consider the following problem, which is to recover a sparse vector $x^* \in \mathbb{R}^n$ from an observation vector $y \in \mathbb{R}^m$ with noise $\varepsilon \in \mathbb{R}^m$ (e.g., additive Gaussian white noise):

$$y = Ax^* + \varepsilon, \quad (1)$$

where $A \in \mathbb{R}^{m \times n}$ ($m \ll n$) is the dictionary matrix. Generally, (1) is an ill-posed problem. Therefore, some prior information such as sparsity or low-rankness needs to be incorporated, for example, x^* is sparse, i.e., the number of elements of the support set of x^* , $S = \{i|x_i^* \neq 0\}$, is much smaller than the dimension n . A common way to estimate x^* is to solve the Lasso problem (Tibshirani 1996):

$$\min_{x \in \mathbb{R}^n} P(x) = f(x) + g(x) = \frac{1}{2}\|y - Ax\|_2^2 + \lambda\|x\|_1, \quad (2)$$

where $\lambda \geq 0$ is a regularization parameter. There are many methods for solving the problem of sparse coding, such as least angle regression (Efron et al. 2004), iterative shrinkage thresholding algorithm (ISTA) (Daubechies, Defrise, and De Mol 2004; Blumensath and Davies 2008) and approxi-

mate message passing (AMP) (Donoho, Maleki, and Montanari 2009). The update rule of ISTA is

$$x^{t+1} = \text{ST}\left(x^t + \frac{1}{L}A^T(y - Ax^t), \frac{\lambda}{L}\right), \quad t = 0, 1, 2, \dots, \quad (3)$$

where $\text{ST}(\cdot, \theta)$ is the soft-thresholding function (ST) with the threshold θ , $\frac{1}{L}$ is the step size which should be taken in $(0, \frac{2}{L})$, where L is usually taken as the largest eigenvalue of $A^T A$, and $A^T(Ax^t - y)$ is actually equal to $\nabla f(x^t)$.

ISTA converges slowly with only a sublinear rate (Beck and Teboulle 2009). Inspired by ISTA and Deep Neural Networks (DNNs) (LeCun, Bengio, and Hinton 2015), Gregor and LeCun (2010) viewed ISTA as a recurrent neural network (RNN) and proposed a learning-based model named Learned ISTA (LISTA):

$$x^{t+1} = \text{ST}(W_1^t y + W_2^t x^t, \theta^t), \quad t = 0, 1, 2, \dots, \quad (4)$$

where W_1^t , W_2^t and θ^t are initialized as $\frac{1}{L}A^T$, $I - \frac{1}{L}A^T A$ and $\frac{\lambda}{L}$, respectively. All the parameters $\Theta = \{W_1^t, W_2^t, \theta^t\}$ are learnable and data-driven. Many empirical and theoretical results (Aberdam, Golts, and Elad 2020; Giryes et al. 2018) have shown that T -layer LISTA or its variants can recover x^* from y more accurately and use one or two orders-of-magnitude fewer iterations than the original ISTA. Moreover, the CSC version of LISTA can be used to explain the CNN in series (Papyan, Romano, and Elad 2017).

On one hand, inspired by (Gregor and LeCun 2010), many learnable network methods such as (Wang, Ling, and Huang 2016; Sprechmann, Bronstein, and Sapiro 2015; Ito, Takabe, and Wadayama 2019; Borgerding, Schniter, and Rangan 2017; Sreret and Giryes 2018) have been proposed and successfully used in different fields, and got satisfactory experimental results.

On the other hand, many works (Xin et al. 2016b; Giryes et al. 2018; Moreau and Bruna 2017; Chen et al. 2018; Liu et al. 2019; Wu et al. 2020; Ablin et al. 2019) discussed LISTA and its variants from a theoretical point of view. Among them, Xin et al. (2016b) first discussed LIHT (Wang, Ling, and Huang 2016), which was obtained by unfolding the iterative hard thresholding (IHT) (Blumensath and Davies 2009) inspired by (Gregor and LeCun 2010), in terms of improving the restricted isometry property (RIP) constant. Inspired by (Xin et al. 2016b), He et al. (2017)

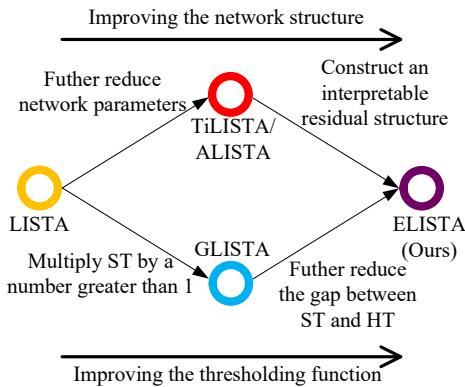


Figure 1: Subsequent improvements on LISTA

connected sparse Bayesian learning (SBL) (Tipping 2001) with long short-term memory (LSTM) (Gers, Schraudolph, and Schmidhuber 2002), and Moreau and Bruna (2017) explained the mechanism of LISTA by re-factorizing the Gram matrix of dictionary. Other works (Chen et al. 2018; Liu et al. 2019; Wu et al. 2020; Ablin et al. 2019) related to this paper will be detailed in Section 1.1.

A series of studies on LISTA have attracted increasing attentions and inspired many subsequent works in different aspects, including learning based optimization (Xie et al. 2019; Sun et al. 2016), design of DNNs (Metzler, Mousavi, and Baraniuk 2017; Zhang and Ghanem 2018; Zhou et al. 2018; Chen et al. 2020; Rick Chang et al. 2017; Zhang et al. 2020; Simon and Elad 2019) and interpreting the DNNs (Zarka et al. 2020; Popyan, Romano, and Elad 2017; Ablardam, Sulam, and Elad 2019; Sulam et al. 2018, 2019).

1.1 Related Works

Chen et al. (2018) proved the coupling relationship between W_1^t and W_2^t , i.e., $W_2^t \rightarrow (I - W_1^t A)$ when $t \rightarrow \infty$, which greatly reduced the number of learnable parameters of LISTA. They also first provided the rigorous proof of the linear convergence of LISTA, which is the basis of the subsequent works. Moreover, the subsequent improvements of LISTA can be divided into two categories: improvements of the network structure and the thresholding function.

For the improvement of the network structure, Liu et al. (2019) further reduced the number of learnable parameters by proposing a novel algorithm, whose update rule is $x^{t+1} = ST(x^t - \alpha^t W(Ax^t - y), \theta^t)$, where α^t is a learnable scaler. They proposed TiLISTA when W is a learnable parameter and ALISTA when W is obtained by solving a data-independent optimization problem. For the improvement of the thresholding function, Wu et al. (2020) argued that the code components in LISTA estimations may be lower than expected, i.e., the algorithms require gains. Inspired by gated recurrent unit (GRU) (Cho et al. 2014a; Chung et al. 2015), Wu et al. (2020) proposed GLISTA, which can be viewed as multiplying ST by a coefficient greater than 1 to reduce the gap between ST and HT. All the improvements of LISTA in different aspects above are shown in Figure 1, where ELISTA is an innovative algorithm proposed in this paper,

which will be described in detail in Section 2.

Moreover, Ablin et al. (2019) also discussed LISTA from the theoretical aspect. They proposed a simple step size strategy which can improve the convergence rate of ISTA by leveraging the space of the iterates, and presented a network named SLISTA to learn only the step size of ISTA for unsupervised training.

1.2 Motivations and Main Contributions

We attempt to answer the following questions, which are not fully addressed in literature yet:

- All the existing variants of LISTA with convergence proofs are serial, the residual network (Res-Net) (He et al. 2016), which is influential in deep learning has not been introduced into LISTA. An important reason is that changing the original structure of LISTA will destroy its excellent mathematical interpretability. Can we get a LISTA with a interpretable residual structure, which has a convergence guarantee?

- Recent studies (Fan and Li 2001; Gu, Wang, and Liu 2014; Xu and Gu 2016; Zhu and Gu 2015; Lederer 2013; Deledalle et al. 2017) have shown that ST may cause large estimation bias, and incurs worse empirical performance than the hard-thresholding function (HT), which means there are some limitations by using ST for sparse coding. Can we improve the thresholding function to reduce the gap between ST and HT?

Our Main Contributions: The main contributions of this paper are as follows:

- We propose a novel variant of LISTA with residual structure by introducing the idea of extragradient into LISTA and establishing the relationship with Res-Net, which is an improvement about the network structure for solving the sparse coding problem. To the best of our knowledge, this is the first residual structure LISTA with theoretical guarantee.

- We design a new thresholding function, called Multistage-Thresholding function (MT), to reduce the gap between ST and HT. A large number of experiments show that MT can ensure the sparsity of the representation as low as possible and obtain effective sparse representation.

- Using extragradient and the MT operator, we propose a novel algorithm, named Extragradient based LISTA (ELISTA), and prove the convergence of ELISTA. Moreover, we conduct ablation experiments to verify the effectiveness of each of our improvements. Extensive experimental results show our ELISTA is superior to the state-of-the-art methods.

2 Extragradient Based LISTA and Multistage-thresholding

In this section, we first introduce the idea of extragradient into LISTA. Then we propose a new multistage-thresholding function (MT) and analyze its advantages. Finally, by combining the idea of extragradient and MT, we propose an innovative algorithm, named *Extragradient based LISTA* (ELISTA), and depict it in detail. Moreover, we also establish the relationship between ELISTA and Res-Net, which is one of the reasons why our algorithm is advantageous.

169 2.1 Extragradient Method

170 We note that iterative algorithms, such as ISTA, can actually be treated as a proximal gradient descent method, which
 171 is a first-order optimization algorithm, for special objective
 172 functions. Thus, we want to introduce the idea of extragradient
 173 into the related iterative algorithms. The extragradient
 174 method was first proposed by (Korpelevich 1976), which is a
 175 classical method for variational inequality problems. For opti-
 176 mization problems, the idea of extragradient was first used
 177 in (Nguyen et al. 2018), which proposed an extended extra-
 178 gradient method (EEG) by combining this idea with some
 179 first-order descent methods. In the t -th iteration of EEG, it
 180 first calculates the gradient at x^t , and updates x^t according
 181 to the gradient to get a middle point $x^{t+\frac{1}{2}}$, then calculates the
 182 gradient at $x^{t+\frac{1}{2}}$, and updates the original point x^t accord-
 183 ing to the gradient at the middle point $x^{t+\frac{1}{2}}$ to obtain x^{t+1} ,
 184 which is the key idea of extragradient. Intuitively, the addi-
 185 tional step in each iteration of EEG allows us to examine the
 186 geometry of the problem and consider its curvature informa-
 187 tion, which is one of the most important bottlenecks for first-
 188 order methods. Thus, by using the idea of extragradient, we
 189 can get a better result after each iteration. The update rules
 190 of EEG for Problem (2) can be rewritten as follows:
 191

$$\begin{aligned} x^{t+\frac{1}{2}} &= \text{ST}\left(x^t - \frac{1}{L} A^T (Ax^t - y), \frac{\lambda}{L}\right), \\ x^{t+1} &= \text{ST}\left(x^t - \frac{1}{L} A^T (Ax^{t+\frac{1}{2}} - y), \frac{\lambda}{L}\right). \end{aligned} \quad (5)$$

192 This form of EEG is similar to ISTA, thus it can be regarded
 193 as a generalization of ISTA.

194 2.2 Multistage-thresholding

The nonlinear transformations in most LISTA related algorithms are realized by the standard ST. However, according to its definition, we know that ST has a weakness, i.e., $|x_i^t|$ obtained from the algorithms with ST is actually smaller than the real $|x_i^*|$, which was described by Proposition 1 in (Wu et al. 2020) and alleviated by (Wu et al. 2020) with the proposal of a gain gate (GG) and an algorithm called GLISTA, whose update rule is as follows:

$$x^{t+1} = \text{ST}(W^t(g_t(x^t, y|\Lambda_g^t) \odot x^t) + U^t y, b^t),$$

where $g_t(x^t, y|\Lambda_g^t)$ is the gate function and greater than 1, and Λ_g^t is the set of its parameters to learn. Besides, W^t , U^t and b^t are also learnable parameters. We define $\tilde{x}^t \triangleq g_t(x^t, y|\Lambda_g^t) \odot x^t$, and obtain

$$\tilde{x}^{t+1} = g_{t+1}(x^{t+1}, y|\Lambda_g^{t+1}) \odot \text{ST}(W^t \tilde{x}^t + U^t y, b^t),$$

which means that GLISTA multiplies ST by a number greater than 1, thus reducing the gap between ST and HT. Therefore, GLISTA can be treated as an improvement of ST. However, the proposal of GG in (Wu et al. 2020) is based on the assumption that there is no "false positive", which is not always true in reality. Therefore, GLISTA will increase some values that should be decreased, which will bring bad results. To address the issue, we design and propose an innovative thresholding function called *Multistage-Thresholding*

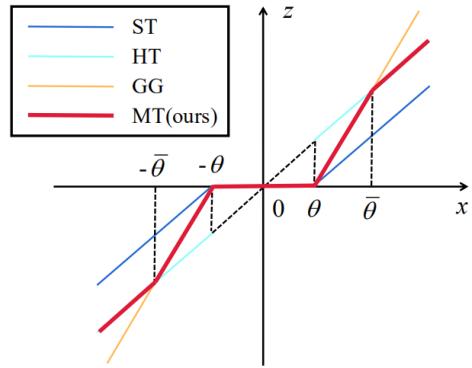


Figure 2: Different thresholding functions

function (MT), which is defined as follows:

$$z = \text{MT}(x, \theta, \bar{\theta}) \stackrel{\Delta}{=} \begin{cases} 0, & 0 \leq |x| < \theta, \\ \frac{\bar{\theta}}{\theta - \bar{\theta}} \text{sign}(x)(|x| - \theta), & \theta \leq |x| < \bar{\theta}, \\ x, & |x| \geq \bar{\theta}. \end{cases} \quad (6)$$

Different thresholding functions are shown in Figure 2, from which we know that MT is equal to GG when $0 \leq |x| < \bar{\theta}$, which plays the role of gain to ST, and when $|x| \geq \bar{\theta}$, it is equal to HT, which makes the result more accurate. Therefore, compared with other thresholding functions, MT can get a better result at each layer.

Our MT is similar to the function $\text{HELU}_\sigma(\cdot)$ proposed in (Wang, Ling, and Huang 2016). However, the motivation of its proposal and the internal mathematical mechanism are different from those of MT. We will give detailed explanations and verifications in the Supplementary Material.

2.3 Extragradient Based LISTA and the Relationship with Res-Net

In order to speed up the convergence of EEG, we combine the algorithm with deep networks and regard $\frac{1}{L} A^T$ and two thresholds of two steps in (5) as learnable parameters, and get the following update rules:

$$\begin{aligned} x^{t+\frac{1}{2}} &= \text{ST}(x^t - W_1^t(Ax^t - y), \theta_1^t), \\ x^{t+1} &= \text{ST}(x^t - W_2^t(Ax^{t+\frac{1}{2}} - y), \theta_2^t). \end{aligned} \quad (7)$$

However, since the above scheme has two different matrices W_1^t and W_2^t to learn in each layer, the number of network parameters greatly increases and the training of the network slows down significantly. Therefore, to address this issue and further establish the connection between the two steps of (7), we convert W_1^t and W_2^t into $\alpha_1^t W^t$ and $\alpha_2^t W^t$, respectively, where α_1^t and α_2^t are two scalars to learn. Then, inspired by (Liu et al. 2019), we change the W^t of each layer into the same W and get a tied algorithm, which can significantly reduce the number of learnable parameters. By replacing ST with MT, we finally obtain the following update rules for our *Extragradient Based LISTA* (ELISTA):

$$\begin{aligned} x^{t+\frac{1}{2}} &= \text{MT}(x^t - \alpha_1^t W(Ax^t - y), \theta_1^t, \bar{\theta}_1^t), \\ x^{t+1} &= \text{MT}(x^t - \alpha_2^t W(Ax^{t+\frac{1}{2}} - y), \theta_2^t, \bar{\theta}_2^t), \end{aligned} \quad (8)$$

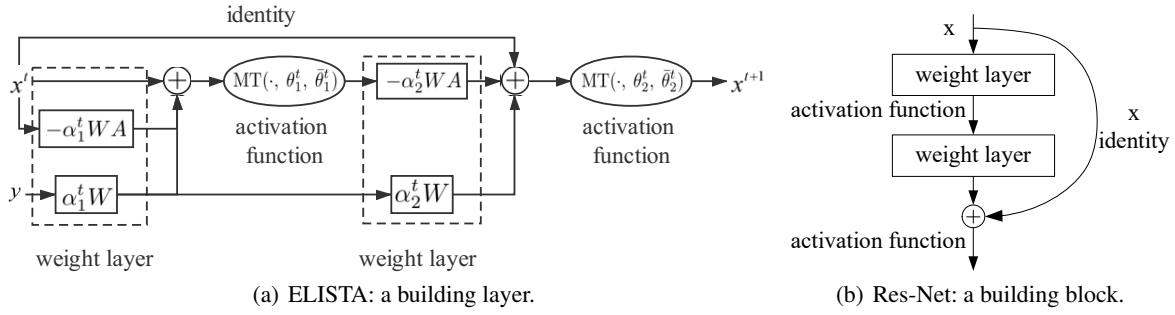


Figure 3: Comparison of the network structures of ELISTA and Res-Net.

Table 1: Comparison of the number of parameters to learn in different methods.

| LISTA | LAMP | GLISTA | ELISTA-m-t | ELISTA-m | ELISTA-t | ELISTA |
|------------------------|------------------------|------------------------|------------------------|-----------------------|------------------------|-----------------------|
| $\mathcal{O}(TMN + T)$ | $\mathcal{O}(TMN + T)$ | $\mathcal{O}(TMN + T)$ | $\mathcal{O}(TMN + T)$ | $\mathcal{O}(MN + T)$ | $\mathcal{O}(TMN + T)$ | $\mathcal{O}(MN + T)$ |

234 where $\bar{\theta}_1^t$ and $\bar{\theta}_2^t$ are also learnable parameters.

235 In order to make the algorithms in this paper easy to distinguish, we present the following naming system:

236 *ELISTA is our main algorithm, which is obtained by introducing the idea of extragradient into LISTA and using MT, and it is a tied algorithm. It should be emphasized that we use +m or -m to represent using MT or not, and -t to indicate that the algorithm is untied. For example, ELISTA-m means ELISTA using ST instead of MT.*

237 Besides, according to (8), we can get the network structure diagram of ELISTA. Through our observation and comparison, we find that the network structure of ELISTA is corresponding to the Res-Net. Since y is already given, we can regard y as a bias. Thus, from Figure 3, we can see that the structure of the network obtained by ELISTA is the same as that of Res-Net, including weight layer, activation function and identity. As we all know, Res-Net can obtain a better performance by improving the network structure. Therefore, it is meaningful to discuss and study the explanation for the internal mathematical mechanism of Res-Net. On the one hand, to some extend, our algorithm may be regarded as a mathematical explanation of the reason for the superiority of Res-Net. On the other hand, the connection and combination of ELISTA and Res-Net might be able to explain why our algorithm has better performance than existing methods.

238 Moreover, the comparison on the number of parameters of the network corresponding to different algorithms is shown in Table 1, where LAMP (Borgerding, Schniter, and Rangan 2017) is an algorithm to transform AMP (Donoho, Maleki, and Montanari 2009) into a neural network inspired by (Gregor and LeCun 2010).

3 Convergence Analysis

239 In this section, we provide the convergence analysis of our algorithms. We first give a basic assumption and two useful definitions. Then we provide the convergence property of 240 ELISTA, and that of ELISTA-t is similar. We note that our 241 analysis, like that of Theorems 3 and 4 of (Wu et al. 2020),

242 is proved under the existence of “false positive”, while the 243 theoretical analysis of (Chen et al. 2018; Liu et al. 2019) was 244 provided under the assumption of no “false positive”, which 245 is difficult to satisfy in reality.

246 **Assumption 1** (Basic assumption). *The signal x^* is sampled from the following set:*

$$x^* \in \mathcal{X}(B, s) \stackrel{\Delta}{=} \{x^* | |x_i^*| \leq B, \forall i, \|x^*\|_0 \leq s\}.$$

247 In other words, x^* is bounded and s -sparse ($s \geq 2$). Furthermore, we assume $\varepsilon = 0$.

248 We note that this assumption is a basic assumption for this 249 class of algorithms. Almost all the related algorithms need 250 to satisfy this assumption, for example (Liu et al. 2019; Wu 251 et al. 2020).

252 **Definition 1** (Liu et al. 2019). *Given a matrix $A \in \mathbb{R}^{m \times n}$, its generalized mutual coherence is defined as follows:*

$$\mu(A) = \inf_{W \in \mathbb{R}^{n \times m}, W_{i,:} A_{:,i} = 1, \forall i} \left\{ \max_{i \neq j, 1 \leq i, j \leq n} W_{i,:} A_{:,j} \right\}.$$

253 We let $\mathcal{W}(A)$ denote a set of all matrices with the generalized mutual coherence $\mu(A)$, which means that

$$\begin{aligned} & \mathcal{W}(A) \\ &= \left\{ W \mid \max_{i \neq j, 1 \leq i, j \leq n} W_{i,:} A_{:,j} = \mu(A), W_{i,:} A_{:,i} = 1, \forall i \right\}. \end{aligned}$$

254 A weight matrix W is “good” if $W \in \mathcal{W}(A)$.

255 This definition is also described in Definition 1 in (Liu 256 et al. 2019). From Lemma 1 in (Chen et al. 2018), we know 257 $\mathcal{W}(A) \neq \emptyset$.

258 **Definition 2.** *Given a model with Θ , in which*

$$\theta_1^t = \Gamma \mu(A) \sup_{x^*} \|x^t - x^*\|_1, \quad \theta_2^t = \Gamma \mu(A) \sup_{x^*} \|x^{t+\frac{1}{2}} - x^*\|_1,$$

Table 2: The results of ablation experiments. We use +m or -m to represent using MT or not, and -t to indicate that the algorithm is untied.

| | Verify the network structure | | | | Verify the thresholding function | | Ours | |
|-------|------------------------------|---------|------------|----------|----------------------------------|---------|----------|----------------|
| | LISTA | TiLISTA | ELISTA-m-t | ELISTA-m | GLISTA | LISTA+m | ELISTA-t | ELISTA |
| NMSE | -36.01 | -50.28 | -51.82 | -65.66 | -63.73 | -62.21 | -77.03 | -107.48 |
| FLSNE | 0.16 | 0.02 | 0.10 | 0.02 | 0.02 | 0.12 | 0.04 | 0.00 |
| SPERR | 147.12 | 46.26 | 3.23 | 2.35 | 57.22 | 0.80 | 0.15 | 0.01 |

we use $\omega_{t+\frac{1}{2}}(k_{t+\frac{1}{2}}|\Theta)$ and $\omega_{t+1}(k_{t+1}|\Theta)$ to characterize its relationship with the “false positive” rate, which is

$$\begin{aligned} & \omega_{k+\frac{1}{2}}(k_{t+\frac{1}{2}}|\Theta) \\ &= \sup_{\forall x^*, |supp(\tilde{x}^{t+\frac{1}{2}}) \cup supp(x^*)| \leq |supp(x^*)| + k_{t+\frac{1}{2}}} \Gamma, \\ & \omega_{k+1}(k_{t+1}|\Theta) \\ &= \sup_{\forall x^*, |supp(\tilde{x}^{t+1}) \cup supp(x^*)| \leq |supp(x^*)| + k_{t+1}} \Gamma, \end{aligned}$$

where $\tilde{x}^{t+\frac{1}{2}} := \text{MT}((I - \alpha_1^t W A)(x^{t+\frac{1}{2}} - x^*), \theta_1^t)$, $\tilde{x}^{t+1} := \text{MT}((I - \alpha_2^t W A)(x^{t+1} - x^*), \theta_2^t)$ and $k_{t+\frac{1}{2}}$ and k_{t+1} are the desired maximal number of “false positive” of $x^{t+\frac{1}{2}}$ and x^{t+1} , respectively.

This definition is similar to Definition 2 in (Wu et al. 2020). Besides, this definition is only an example for ELISTA. For our ELISTA-t, we can also easily get a similar definition.

Based on the assumption and these two definitions, we can get the linear convergence of ELISTA, which can be given by the following theorem.

Theorem 1 (Linear Convergence for ELISTA). *If Assumption 1 holds, $W \in \mathcal{W}(A)$ can be satisfied by selecting W properly,*

$$\begin{aligned} \theta_1^t &= \alpha_1^t \omega_{t+\frac{1}{2}}(k_{t+\frac{1}{2}}|\Theta) \mu(A) \sup_{x^*} \|x^t - x^*\|_1, \\ \theta_2^t &= \alpha_2^t \omega_{t+1}(k_{t+1}|\Theta) \mu(A) \sup_{x^*} \|x^{t+\frac{1}{2}} - x^*\|_1, \end{aligned} \quad (9)$$

$\bar{\theta}_1^t \geq \theta_1^t + |\tilde{x}_i^{t+\frac{1}{2}}|$, $\bar{\theta}_2^t \geq \theta_2^t + |\tilde{x}_i^{t+1}|$ are achieved, α_1^t and α_2^t belong to specific bounded intervals for different cases, and s is small enough, then for the sequences generated by ELISTA, the following result holds

$$\|x^t - x^*\|_2 \leq sB \exp \left(\sum_{i=1}^t c'_i \right) < sB \exp(c' t), \quad (10)$$

where $c' = \max_{i=1,2,\dots,t} \{c'_i\}$. $\exists t_0 = \lceil -\log(\frac{sB}{\sigma})/c \rceil$, for $i \geq t_0$, $0 < k_{i-\frac{1}{2}}, k_i < s$, if $\gamma^{i-\frac{1}{2}} = \gamma^i = 0$, then $c'_i < 0$, and for $i > t_0$, $k_{i-\frac{1}{2}} = k_i = 0$, if $1 - \omega_{i-\frac{1}{2}}(s|\Theta) < \gamma^{i-\frac{1}{2}} \leq 1$ and $1 - \omega_i(s|\Theta) < \gamma^i \leq 1$, then $c'_i < 0$. Thus, $c' < 0$.

The definitions of $\gamma^{i-\frac{1}{2}}$ and γ^i are given in the detailed proof of this theorem in the Supplementary Material. Here we give a simple sketch of the full proof:

To prove Theorem 1, we first need to obtain the relationship between $\|x^{t+1} - x^*\|_2$ and $\|x^t - x^*\|_2$. To calculate all non-zero elements of $x^{t+\frac{1}{2}} - x^*$, we divide them into three parts: $i \in \bar{S}^{(t+\frac{1}{2})}$, $i \in S \setminus \bar{S}^{(t+\frac{1}{2})}$ and $i \in$

$S^{(t+\frac{1}{2})}$, where $S \triangleq \text{supp}(x^*)$, $\bar{S}^{(t+\frac{1}{2})} \triangleq S \cap \text{supp}(x^{t+\frac{1}{2}})$ and $S^{(t+\frac{1}{2})} \triangleq \{i|i \in \text{supp}(x^{t+\frac{1}{2}}), i \notin S\}$, and then sum the results to obtain the relationship between $\|x^{t+\frac{1}{2}} - x^*\|_1$ and $\|x^t - x^*\|_1$. In a similar way, we can get the relationship between $\|x^{t+1} - x^*\|_1$, $\|x^{t+\frac{1}{2}} - x^*\|_1$ and $\|x^t - x^*\|_1$. Then, we can obtain the relationship between $\|x^{t+1} - x^*\|_1$ and $\|x^t - x^*\|_1$, and thus the relationship between $\|x^{t+1} - x^*\|_2$ and $\|x^t - x^*\|_2$. Finally, Theorem 1 can be proved by the recursion in terms of t .

4 Numerical Results

In this section, we first perform ablation experiments to verify the effectiveness of our method and provide the justification of some parameters in the algorithms and the verification of an assumption. Then we evaluate our ELISTA and ELISTA-t in terms of sparse representation performance, natural image inpainting, 3D geometry recovery via photometric stereo, support set accuracy and unsupervised experiment as in (Ablin et al. 2019). All the experimental settings are the same as previous works (Chen et al. 2018; Liu et al. 2019; Wu et al. 2020; Borgerding, Schniter, and Rangan 2017). We find that Support Selection (SS) (Chen et al. 2018) can generally improve the performance of related networks including ours. However, the performance of SS is greatly affected by the hyper parameters, and it is necessary to know the sparsity of x^* in advance to set the hyper parameters, which is difficult to get in real situations. Thus, in order to more fairly compare the impact of the network itself on performance, all the networks do not use SS. All training follows (Chen et al. 2018) (The details are provided in the Supplementary Material). For all our methods, α_1^t and α_2^t are initialized as 1.0. θ_1^t and θ_2^t are initialized as $\frac{\lambda}{L}$ when using ST, while θ_1^t and θ_2^t are initialized as $\frac{\lambda}{L} - 0.1$, $\bar{\theta}_1^t$ and $\bar{\theta}_2^t$ are initialized as $\frac{\lambda}{L}$ when using MT. All the results are obtained by running ten times and averaged. Verification of the parameters and the assumption, support set accuracy and unsupervised experiment are presented in the Supplementary Material.

4.1 Ablation Experiments

In this subsection, by controlling variables, we compare our ELISTA-m with LISTA (Gregor and LeCun 2010; Chen et al. 2018) and TiLISTA (Liu et al. 2019), and compare LISTA+m¹ with LISTA (Gregor and LeCun 2010; Chen

¹LISTA+m is an algorithm which replaces ST in LISTA with MT.

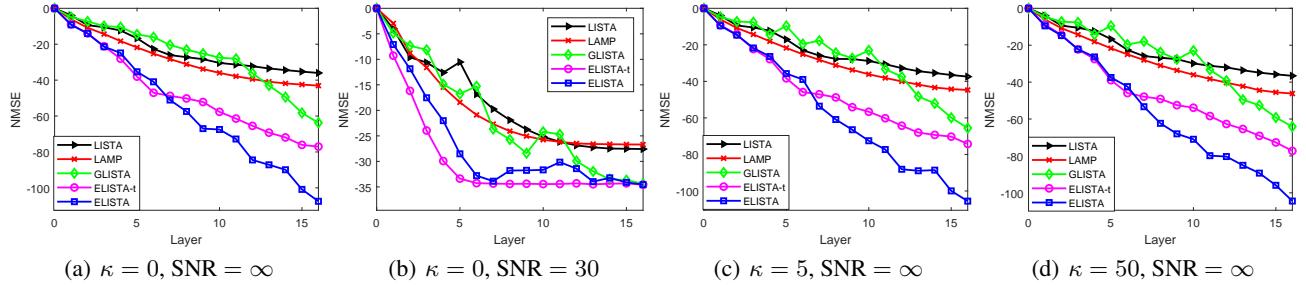


Figure 4: Comparison of sparse representation with different layers under different SNR and κ .

Table 3: The PSNR of natural image inpainting

| | Barbara | Boat | House | Lena | Peppers | C.man | Couple | Finger | Hill | Man | Montage |
|--------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| ISTA | 23.51 | 25.38 | 26.88 | 26.11 | 23.53 | 22.73 | 25.33 | 20.64 | 27.28 | 24.25 | 21.29 |
| LISTA | 24.52 | 27.29 | 29.50 | 27.84 | 25.78 | 24.51 | 27.20 | 23.60 | 28.92 | 26.32 | 22.50 |
| GLISTA | 25.30 | 28.95 | 30.95 | 29.97 | 27.64 | 25.76 | 27.48 | 26.29 | 29.53 | 28.14 | 24.31 |
| LFISTA | 26.01 | 29.68 | 32.06 | 32.12 | 28.57 | 26.77 | 29.77 | 28.10 | 30.69 | 30.22 | 26.94 |
| ELISTA | 26.60 | 30.33 | 32.76 | 32.75 | 29.61 | 27.67 | 30.09 | 28.20 | 30.41 | 30.36 | 28.49 |

et al. 2018) and GLISTA (Wu et al. 2020) in the noiseless condition to verify the improvement of the network structure and that of the thresholding function, respectively. For TiLISTA, we set

$$\begin{aligned} x^{t+\frac{1}{2}} &= \text{ST}(x^t - \alpha_1^t W(Ax^t - y), \theta_1^t) \\ x^{t+1} &= \text{ST}(x^{t+\frac{1}{2}} - \alpha_2^t W(Ax^{t+\frac{1}{2}} - y), \theta_2^t). \end{aligned} \quad (11)$$

as one layer². We set $m = 250$, $n = 500$ and $T = 16$. α_1^t and α_2^t in TiLISTA are also initialized as 1.0. We sample the elements of the dictionary matrix A randomly from a standard Gaussian distribution in simulations, the ground-truth x^* is also generated by the standard Gaussian distribution and we use Bernoulli distribution with a probability of 0.1 to ensure the sparsity. y is produced by A , x^* and noise ε . All experimental results are on the test set. The sparse representation performance is evaluated by NMSE (in dB):

$$\text{NMSE}(\hat{x}, x^*) = 10 \log_{10} \left(\frac{\mathbb{E} \|\hat{x} - x^*\|^2}{\mathbb{E} \|x^*\|^2} \right). \quad (12)$$

We use NMSE, FLSNE and SPERR as indicators to evaluate the networks, where NMSE is defined in (12), FLSNE is the number of “false negative” and SPERR denotes the number of support error.

From Table 2, we can find that: (i) Because of the residual structure brought by the extragradient, ELISTA-m is superior to LISTA and TiLISTA in terms of NMSE and SPERR, where the two latter are serial connection. (ii) ST tends to expand the size of the support set to get a smaller FLSNE, however this also leads to a very large SPERR and a worse NMSE. GG can obtain better results than ST by narrowing

the gap between ST and HT, but the SPERR of GLISTA is still large. That is, ST and GG expand the size of the support set in order to obtain a better sparse representation, so as to obtain a sparse representation that is not sparse. The residual structure induced by the extragradient can alleviate the problem of ST. Since MT is closer to HT, it can obtain a more sparse representation, which in turn enhances NMSE. Because our ELISTA is an improved algorithm combining these two improvements, it outperforms all the other algorithms, which also shows the effectiveness of the residual structure and the improvement of our thresholding function.

4.2 Sparse Representation Performance

In this subsection, we compare our ELISTA and ELISTA-t with the state-of-the-art methods: LISTA (Gregor and Le-Cun 2010; Chen et al. 2018), LAMP (Borgerding, Schniter, and Rangan 2017) and GLISTA (Wu et al. 2020). We train the networks with three different noise levels: SNR (Signal-to-Noise Ratio) = 30, 40, ∞ and three different ill conditioned matrices A with condition numbers $\kappa = 5, 30, 50$.

Figure 4 shows that our methods are obviously better than the compared methods in terms of both convergence speed and accuracy in the noiseless case. Especially, compared with LISTA, the NMSE performance of our methods is nearly twice better than that of LISTA. In the presence of noise, our methods achieve the state-of-the-art convergence accuracy and are obviously better than other methods in terms of convergence speed. We note that due to the limitation of space, only part of the results are given here, and more results are reported in the Supplementary Material.

4.3 Natural Image Inpainting

In this subsection, we apply our algorithm to solve the natural image inpainting problem, and comparing it with LISTA

²The definition of one layer is different from that of (Liu et al. 2019). The purpose of this change is to control variables to verify the validity of our ELISTA.

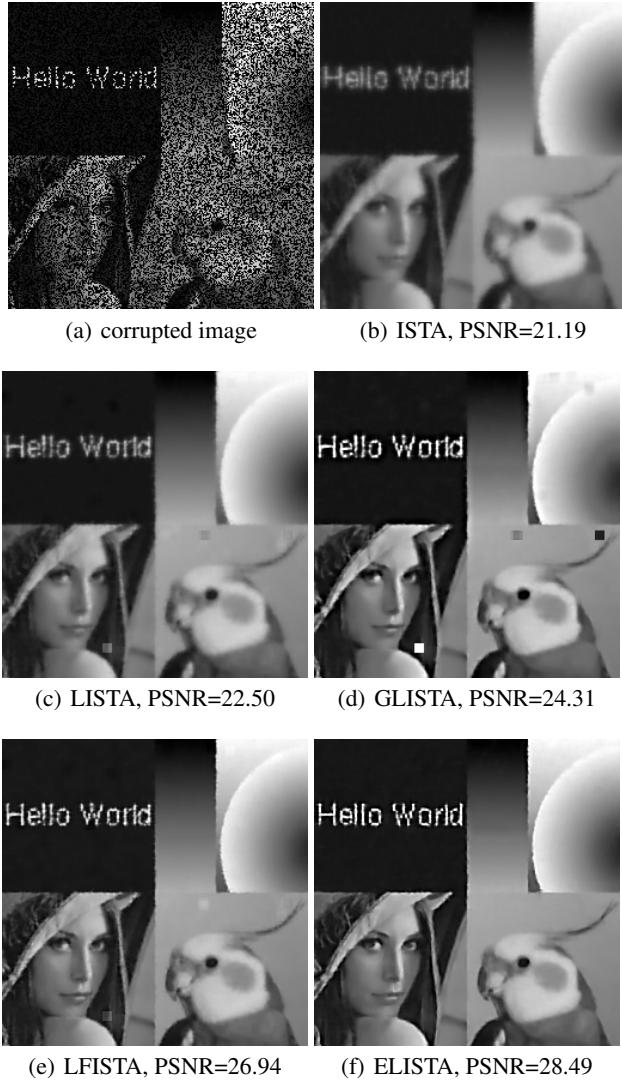


Figure 5: Image inpainting with 50% missing pixels on Montage.

Table 4: The mean angular error of 3D geometry recovery via photometric stereo

| q | LISTA | GLISTA | ELISTA-t | ELISTA |
|-----|---------|---------|----------------|----------------|
| 35 | 0.06836 | 0.06249 | 0.03534 | 0.02754 |
| 25 | 0.09664 | 0.10033 | 0.05885 | 0.04947 |
| 15 | 0.69334 | 0.63967 | 0.47569 | 0.60010 |

4.4 3D Geometry Recovery via Photometric Stereo

In this subsection, we compare our ELISTA and ELISTA-t with the state-of-the-art methods: LISTA (Gregor and Le-Cun 2010; Chen et al. 2018) and GLISTA (Wu et al. 2020) for 3D Geometry Recovery via Photometric Stereo. Photometric stereovision is a powerful technique used to recover high resolution surface normals from a 3D scene using appearance changes of 2D images in different lighting (Woodham 1980). In practice, however, the estimation process is often interrupted by non-lambert effects, such as highlights, shadows, or image noise. This problem can be solved by decomposing the observation matrix of the superimposed image under different lighting conditions into ideal lambert components and sparse error terms (Wu et al. 2010; Ikehata et al. 2012), i.e., $o = \rho Ln + e$, where $o \in \mathbb{R}^q$ denotes the resulting measurements, $n \in \mathbb{R}^3$ denotes the true surface normal, $L \in \mathbb{R}^{q \times 3}$ defines a lighting direction, ρ is the diffuse albedo, acting here as a scalar multiplier and $e \in \mathbb{R}^q$ is an unknown sparse vector. By multiplying both sides of $o = \rho Ln + e$ by the orthogonal complement to L , we can get $Proj_{null_{[L^\top]}}(o) = Proj_{null_{[L^\top]}}(e)$. Let $Proj_{null_{[L^\top]}}(o)$ be y and $Proj_{null_{[L^\top]}}(e)$ be Ax , e can be obtained by solving the sparse coding problem. Then we can use $L^\dagger(o - e)$ to recover n . The main experimental settings follow (Xin et al. 2016b; Wu et al. 2020; He et al. 2017). Tests are performed using the 32-bit HDR gray-scale images of objects ‘‘Bunny’’ as in (Xin et al. 2016b; Wu et al. 2020; He et al. 2017) with $q = 35, 25, 15$ and 40% of the elements of the sparse noise e are non-zero. From Table 4, we can find that our methods perform much better than LISTA and GLISTA, which is similar to the conclusion we came to in Section 4.2.

5 Conclusions

In this paper, we consider a sparse representation problem. We proposed an innovative algorithm called ELISTA with interpretable residual structure and a better thresholding function. Moreover, we proved that ELISTA can achieve linear convergence in theory. Extensive empirical results verified the high efficiency of our method. One limitation of this paper is that in the theoretical analysis, we use the same assumption as in the previous work (Chen et al. 2018; Liu et al. 2019; Wu et al. 2020), that s , i.e., the sparsity of x^* , is small enough. Removing this common assumption of the related algorithms is our future work.

412 (Chen et al. 2018), LFISTA (Moreau and Bruna 2017; Ab-
413 erdam, Golts, and Elad 2020) and GLISTA (Wu et al. 2020).
414 The training dataset is BSDS500 and the test dataset is Set
415 11. For LFISTA (Aberdam, Golts, and Elad 2020), we use
416 the code provided by this work and for the other algorithms,
417 we implement them ourselves. The PSNR of different algo-
418 rithms are shown in Table 3, the qualitative results on the
419 Montage image are shown in Figure 5 and the other qual-
420 itative results are shown in the Supplementary Material. In
421 addition, detailed experimental setup and other details are
422 also given in the Supplementary Material.

423 From Table 3, Figure 5 and all the other qualitative results
424 in the Supplementary Material, we can see that our ELISTA
425 outperforms other algorithms in most cases.

References

- 470
- 471 Aberdam, A.; Golts, A.; and Elad, M. 2020. Ada-LISTA:
472 Learned Solvers Adaptive to Varying Models. *arXiv preprint
473 arXiv:2001.08456*.
- 474 Aberdam, A.; Sulam, J.; and Elad, M. 2019. Multi-layer
475 sparse coding: The holistic way. *SIAM Journal on Mathe-
476 matics of Data Science* 1(1): 46–77.
- 477 Ablin, P.; Moreau, T.; Massias, M.; and Gramfort, A. 2019.
478 Learning step sizes for unfolded sparse coding. In *Advances
479 in Neural Information Processing Systems*, 13100–13110.
- 480 Beck, A.; and Teboulle, M. 2009. A fast iterative shrinkage-
481 thresholding algorithm for linear inverse problems. *SIAM
482 Journal on Imaging Sciences* 2(1): 183–202.
- 483 Blumensath, T.; and Davies, M. E. 2008. Iterative threshold-
484 ing for sparse approximations. *Journal of Fourier analysis
485 and Applications* 14(5-6): 629–654.
- 486 Blumensath, T.; and Davies, M. E. 2009. Iterative hard
487 thresholding for compressed sensing. *Applied and Compu-
488 tational Harmonic Analysis* 27(3): 265–274.
- 489 Borgerding, M.; Schniter, P.; and Rangan, S. 2017. AMP-
490 inspired deep networks for sparse linear inverse problems.
491 *IEEE Transactions on Signal Processing* 65(16): 4293–
492 4308.
- 493 Chen, X.; Li, Y.; Umarov, R.; Gao, X.; and Song, L. 2020.
494 RNA secondary structure prediction by learning unrolled al-
495 gorithms. In *Proceedings of the International Conference
496 on Learning Representations*.
- 497 Chen, X.; Liu, J.; Wang, Z.; and Yin, W. 2018. Theoretical
498 linear convergence of unfolded ISTA and its practical
499 weights and thresholds. In *Advances in Neural Information
500 Processing Systems*, 9061–9071.
- 501 Cho, K.; Van Merriënboer, B.; Gulcehre, C.; Bahdanau, D.;
502 Bougares, F.; Schwenk, H.; and Bengio, Y. 2014a. Learning
503 phrase representations using RNN encoder-decoder for sta-
504 tistical machine translation. *arXiv preprint arXiv:1406.1078*
505 .
- 506 Cho, K.; Van Merriënboer, B.; Gulcehre, C.; Bahdanau, D.;
507 Bougares, F.; Schwenk, H.; and Bengio, Y. 2014b. Learning
508 phrase representations using RNN encoder-decoder for sta-
509 tistical machine translation. *arXiv preprint arXiv:1406.1078*
510 .
- 511 Chung, J.; Gulcehre, C.; Cho, K.; and Bengio, Y. 2015.
512 Gated feedback recurrent neural networks. In *International
513 Conference on Machine Learning*, 2067–2075.
- 514 Daubechies, I.; Defrise, M.; and De Mol, C. 2004. An itera-
515 tive thresholding algorithm for linear inverse problems with
516 a sparsity constraint. *Communications on Pure and Applied
517 Mathematics: A Journal Issued by the Courant Institute of
518 Mathematical Sciences* 57(11): 1413–1457.
- 519 Deledalle, C.-A.; Papadakis, N.; Salmon, J.; and Vaiter, S.
520 2017. Clear: Covariant least-square refitting with applica-
521 tions to image restoration. *SIAM Journal on Imaging Sci-
522 ences* 10(1): 243–284.
- 523 Donoho, D. L.; Maleki, A.; and Montanari, A. 2009. 523
524 Message-passing algorithms for compressed sensing. *Pro-
525 ceedings of the National Academy of Sciences* 106(45):
526 18914–18919.
- 527 Efron, B.; Hastie, T.; Johnstone, I.; Tibshirani, R.; et al. 527
528 2004. Least angle regression. *The Annals of statistics* 32(2):
529 407–499.
- 530 Fan, J.; and Li, R. 2001. Variable selection via nonconcave 530
531 penalized likelihood and its oracle properties. *Journal of the
532 American Statistical Association* 96(456): 1348–1360.
- 533 Gers, F. A.; Schraudolph, N. N.; and Schmidhuber, J. 2002. 533
534 Learning precise timing with LSTM recurrent networks.
535 *Journal of Machine Learning Research* 3(Aug): 115–143.
- 536 Giryes, R.; Eldar, Y. C.; Bronstein, A. M.; and Sapiro, G. 536
537 2018. Tradeoffs between convergence speed and reconstruc-
538 tion accuracy in inverse problems. *IEEE Transactions on
539 Signal Processing* 66(7): 1676–1690.
- 540 Gregor, K.; and LeCun, Y. 2010. Learning fast approxima-
541 tions of sparse coding. In *Proceedings of the 27th Interna-
542 tional Conference on International Conference on Machine
543 Learning*, 399–406.
- 544 Gu, Q.; Wang, Z.; and Liu, H. 2014. Sparse pca with ora-
545 cle property. In *Advances in Neural Information Processing
546 Systems*, 1529–1537.
- 547 He, H.; Xin, B.; Ikehata, S.; and Wipf, D. 2017. From
548 Bayesian sparsity to gated recurrent nets. In *Advances in
549 Neural Information Processing Systems*, 5554–5564.
- 550 He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual
551 learning for image recognition. In *Proceedings of the IEEE
552 Conference on Computer Vision and Pattern Recognition*,
553 770–778.
- 554 Ikehata, S.; Wipf, D.; Matsushita, Y.; and Aizawa, K. 2012.
555 Robust photometric stereo using sparse regression. In *2012
556 IEEE Conference on Computer Vision and Pattern Recog-
557 nition*, 318–325. IEEE.
- 558 Ito, D.; Takabe, S.; and Wadayama, T. 2019. Trainable ISTA
559 for sparse signal recovery. *IEEE Transactions on Signal Pro-
560 cessing* 67(12): 3113–3125.
- 561 Korpelevich, G. 1976. The extragradient method for finding
562 saddle points and other problems. *Matecon* 12: 747–756.
- 563 LeCun, Y.; Bengio, Y.; and Hinton, G. 2015. Deep learning.
564 *Nature* 521(7553): 436–444.
- 565 Lederer, J. 2013. Trust, but verify: benefits and pitfalls of
566 least-squares refitting in high dimensions. *arXiv preprint
567 arXiv:1306.0113*.
- 568 Liu, J.; Chen, X.; Wang, Z.; and Yin, W. 2019. Alista: Ana-
569 lytic weights are as good as learned weights in lista. In *Pro-
570 ceedings of the International Conference on Learning Rep-
571 resentations*.
- 572 Metzler, C.; Mousavi, A.; and Baraniuk, R. 2017. Learned
573 D-AMP: Principled neural network based compressive im-
574 age recovery. In *Advances in Neural Information Processing
575 Systems*, 1772–1783.

- 576 Moreau, T.; and Bruna, J. 2017. Understanding trainable
 577 sparse coding via matrix factorization. In *Proceedings of*
 578 *the International Conference on Learning Representations*.
 579 Nguyen, T. P.; Pauwels, E.; Richard, E.; and Suter, B. W.
 580 2018. Extragradient method in optimization: Convergence
 581 and complexity. *Journal of Optimization Theory and Appli-*
 582 *cations* 176(1): 137–162.
- 583 Petyan, V.; Romano, Y.; and Elad, M. 2017. Convolutional
 584 neural networks analyzed via convolutional sparse coding.
 585 *The Journal of Machine Learning Research* 18(1): 2887–
 586 2938.
- 587 Rick Chang, J.; Li, C.-L.; Poczos, B.; Vijaya Kumar, B.; and
 588 Sankaranarayanan, A. C. 2017. One Network to Solve Them
 589 All—Solving Linear Inverse Problems Using Deep Projection
 590 Models. In *Proceedings of the IEEE International Confer-*
 591 *ence on Computer Vision*, 5888–5897.
- 592 Simon, D.; and Elad, M. 2019. Rethinking the csc model for
 593 natural images. In *Advances in Neural Information Process-*
 594 *ing Systems*, 2271–2281.
- 595 Sprechmann, P.; Bronstein, A. M.; and Sapiro, G. 2015.
 596 Learning efficient sparse and low rank models. *IEEE Trans-*
 597 *actions on Pattern Analysis and Machine Intelligence* 37(9):
 598 1821–1833.
- 599 Sreter, H.; and Giryes, R. 2018. Learned convolutional
 600 sparse coding. In *2018 IEEE International Conference on*
 601 *Acoustics, Speech and Signal Processing (ICASSP)*, 2191–
 602 2195. IEEE.
- 603 Sulam, J.; Aberdam, A.; Beck, A.; and Elad, M. 2019. On
 604 multi-layer basis pursuit, efficient algorithms and convolu-
 605 tional neural networks. *IEEE transactions on pattern anal-*
 606 *ysis and machine intelligence*.
- 607 Sulam, J.; Petyan, V.; Romano, Y.; and Elad, M. 2018.
 608 Multilayer convolutional sparse modeling: Pursuit and dic-
 609 tionary learning. *IEEE Transactions on Signal Processing*
 610 66(15): 4090–4104.
- 611 Sun, J.; Li, H.; Xu, Z.; et al. 2016. Deep ADMM-Net for
 612 compressive sensing MRI. In *Advances in Neural Informa-*
 613 *tion Processing Systems*, 10–18.
- 614 Tibshirani, R. 1996. Regression shrinkage and selection via
 615 the lasso. *Journal of the Royal Statistical Society: Series B*
 616 (*Methodological*) 58(1): 267–288.
- 617 Tipping, M. E. 2001. Sparse Bayesian learning and the rel-
 618 evance vector machine. *Journal of Machine Learning Re-*
 619 *search* 1(Jun): 211–244.
- 620 Wang, Z.; Ling, Q.; and Huang, T. S. 2016. Learning deep
 621 ℓ_0 encoders. In *Thirtieth AAAI Conference on Artificial In-*
 622 *telligence*.
- 623 Woodham, R. J. 1980. Photometric method for determining
 624 surface orientation from multiple images. *Optical engineer-*
 625 *ing* 19(1): 191139.
- 626 Wu, K.; Guo, Y.; Li, Z.; and Zhang, C. 2020. SPARSE COD-
 627 ING WITH GATED LEARNED ISTA. In *Proceedings of*
 628 *the International Conference on Learning Representations*.
- Wu, L.; Ganesh, A.; Shi, B.; Matsushita, Y.; Wang, Y.; and
 629 Ma, Y. 2010. Robust photometric stereo via low-rank matrix
 630 completion and recovery. In *Asian Conference on Computer*
 631 *Vision*, 703–717. Springer.
- Xie, X.; Wu, J.; Zhong, Z.; Liu, G.; and Lin, Z. 2019. Differ-
 633 entiable linearized ADMM. In *Proceedings of the 27th Inter-*
 634 *national Conference on International Conference on Ma-*
 635 *chine Learning*.
- Xin, B.; Wang, Y.; Gao, W.; and Wipf, D. 2016a. Max-
 637 imal Sparsity with Deep Networks? *arXiv preprint*
 638 *arXiv:1605.01636*.
- Xin, B.; Wang, Y.; Gao, W.; Wipf, D.; and Wang, B. 2016b.
 640 Maximal sparsity with deep networks? In *Advances in Neu-*
 641 *ral Information Processing Systems*, 4340–4348.
- Xu, P.; and Gu, Q. 2016. Semiparametric differential graph
 643 models. In *Advances in Neural Information Processing Sys-*
 644 *tems*, 1064–1072.
- Zarka, J.; Thiry, L.; Angles, T.; and Mallat, S. 2020. Deep
 646 Network classification by Scattering and Homotopy dictio-
 647 nary learning. In *Proceedings of the International Confer-*
 648 *ence on Learning Representations*.
- Zhang, J.; and Ghanem, B. 2018. ISTA-Net: Interpretable
 650 optimization-inspired deep network for image compressive
 651 sensing. In *Proceedings of the IEEE Conference on Com-*
 652 *puter Vision and Pattern Recognition*, 1828–1837.
- Zhang, Q.; Ye, X.; Liu, H.; and Chen, Y. 2020. A
 654 Novel Learnable Gradient Descent Type Algorithm for Non-
 655 convex Non-smooth Inverse Problems. *arXiv preprint*
 656 *arXiv:2003.06748*.
- Zhou, J. T.; Di, K.; Du, J.; Peng, X.; Yang, H.; Pan, S. J.;
 658 Tsang, I. W.; Liu, Y.; Qin, Z.; and Goh, R. S. M. 2018.
 659 SC2Net: Sparse LSTMs for sparse coding. In *Thirty-Second*
 660 *AAAI Conference on Artificial Intelligence*.
- Zhu, R.; and Gu, Q. 2015. Towards a lower sample complex-
 662 ity for robust one-bit compressed sensing. In *International*
 663 *Conference on Machine Learning*, 739–747.

665 Supplementary Material for “Learned 666 Extragradient ISTA with Interpretable 667 Residual Structures for Sparse Coding”

668 In this supplementary material, we first provide a lot of de-
669 tails of the experiments and more experimental results.

670 Moreover, we give the detailed proof of Theorem 1
671 by proposing and proving some lemmas (i.e., Lemma 1 -
672 Lemma 4).

673 Finally, we explain the differences between $\text{HELU}_\sigma(\cdot)$
674 proposed in (Wang, Ling, and Huang 2016) and our MT in
675 detail, including two parts: the difference of motivation and
676 the differences of some specific details. Besides, we give the
677 relevant theoretical analysis and the experimental verifica-
678 tion to support our statements.

679 S1. Experimental Details and More Experimental 680 Results

681 In this part, we first provide the detailed training process of
682 the networks. Next we verify some assumptions, including
683 the setting of parameters in Theorem 1 and an assumption
684 about support sets needed in the detailed proof. Moreover,
685 we also give more experimental results of sparse representa-
686 tion performance (Section 4.2) and the qualitative results for
687 natural image inpainting (Section 4.3). Finally, we provide
688 the experiments of our algorithms in support set accuracy
689 and unsupervised experiment.

690 • Details of Training Process

691 The training batch size is 64, and the size of the verifica-
692 tion set and test set is 1000. We use Adam (Cho et al. 2014b)
693 to train networks and let $\beta_1 = 0.9$, $\beta_2 = 0.999$. All training
694 follows (Chen et al. 2018), i.e., by gradually training the
695 sparse coding network to update more layers, we reduce the
696 learning rate of the current optimization layer and the basic
697 learning rate is 0.0005 when the verification loss of 4000 it-
698 erations is not reduced. When the validation loss is no longer
699 reduced and the learning rate drops to 0.00001, the training
700 on the current level stops. All the experiments were run on
701 two NVIDIA GeForce RTX 2080Ti.

702 • Verification of Some Assumptions

703 Here, we give the experimental verification of
704 some assumptions used in this paper. Firstly, in Fig-
705 ure 6, we verify the assumptions about parameters
706 α_1^t , α_2^t , θ_1^t and θ_2^t for our algorithms ELISTA-t and
707 ELISTA: α_1^t and α_2^t are bounded, θ_1^t , θ_2^t and α_1^t , α_2^t
708 are proportional to $\omega_{t+1/2} \sup_{x^*} \|x^t(x^*) - x^*\|_1$ and
709 $\omega_{t+1} \sup_{x^*} \|x^{t+1/2}(x^*) - x^*\|_1$, respectively.

710 Next, we verify an assumption used in the proof of The-
711 rem 1 (We assume $\text{supp}(x^t) \subset \text{supp}(x^{t+\frac{1}{2}}) \subset \text{supp}(x^{t+1})$
712 in the proof of Lemma 3), which is shown in Figure 7, from
713 which we can see that this assumption is not satisfied at the
714 beginning but will be satisfied in later layers soon.

715 • More Experimental Results of Sparse Representation 716 Performance

717 More experimental results for Section 4.2 are shown in
718 Figure 8. From Figure 8, we can see that our methods are
719 obviously better than the compared methods in terms of both
720 convergence speed and accuracy in the noiseless case (i.e.,

721 $\kappa = 30$, $\text{SNR} = \infty$). Especially, compared with LISTA,
722 the NMSE performance of our methods is nearly twice
723 better than that of LISTA. In the presence of noise (i.e.,
724 $\kappa = 0$, $\text{SNR} = 40$), our methods achieve the state-of-the-
725 art convergence accuracy and are obviously better than the
726 compared methods in terms of convergence speed.
727

• Natural Image Inpainting

728 Moreover, we apply our ELISTA to the task of natural
729 image inpainting as in (Aberdam, Golts, and Elad 2020). We
730 assume that the image is corrupted by a known mask with a
731 missing pixel ratio of p . Therefore, the problem that we need
732 to solve becomes

$$\min_{x \in \mathbb{R}^n} P(x) = \frac{1}{2} \|y - MAx\|_2^2 + \lambda \|x\|_1, \quad (13)$$

733 where $M \in \mathbb{R}^{m \times m}$ represents the mask, being an identity
734 matrix with a percentage of p diagonal elements equal to
735 0, $y \in \mathbb{R}^n$ is a corrupt image patch which has the same
736 size as the clean one, and $A \in \mathbb{R}^{m \times n}$ is a dictionary ob-
737 tained by training on clean image patches. Thus, although
738 the dictionary is given, there is a different and known mask
739 for each patch, which makes the final effective dictionary
740 A_{eff} changes for each signal. Then the corresponding update
741 rules for our ELISTA change to the following form.

$$\begin{aligned} x^{t+\frac{1}{2}} &= \text{MT}(x^t - \alpha_1^t (W_1)^T M^T M W_1 x^t \\ &\quad + \alpha_1^t (W_2)^T M^T y), \theta_1^t, \bar{\theta}_1^t), \\ x^{t+1} &= \text{MT}(x^t - \alpha_2^t (W_1)^T M^T M W_1 x^{t+\frac{1}{2}} \\ &\quad + \alpha_2^t (W_2)^T M^T y), \theta_2^t, \bar{\theta}_2^t). \end{aligned} \quad (14)$$

742 Since the original update rules make LISTA and some of
743 its variants unable to face different dictionaries and cannot
744 be used to solve the inpainting problem, for our algorithm,
745 besides learning some scalar parameters in each layer of the
746 network, it also learns two different matrices W_1 and W_2 as
747 in (Aberdam, Golts, and Elad 2020), which is different from
748 the original ELISTA. And for the fairness of experimental
749 comparison, we change all the compared algorithms to this
750 form.

751 In addition, we also present the qualitative results on the
752 test data set, i.e., Set 11 (except the results on the Montage
753 image, which are shown in the main paper), which are pro-
754 vided in Figures 9-18.

• Support Set Accuracy

755 We compare our ELISTA and ELISTA-t with the state-
756 of-the-art methods: MaxSparseNet (Xin et al. 2016b) and
757 GFLSTM (He et al. 2017). We set $m = 20$, $n = 100$ and
758 different sparsity levels s . We use strict accuracy and loose
759 accuracy as indicators to evaluate the networks, where strict
760 accuracy is percentage of trials with exact support recovery
761 and loose accuracy is percentage of “true positive”. MaxS-
762 parseNet and GFLSTM treat support set recovery as mul-
763 tiple dichotomy problems and only output the index of the
764 support set at the last layer. Different from them, we still
765 output an x^t at each layer, and an index is considered as an
766 index of the support set when it is selected by more than half
767 of x^t .

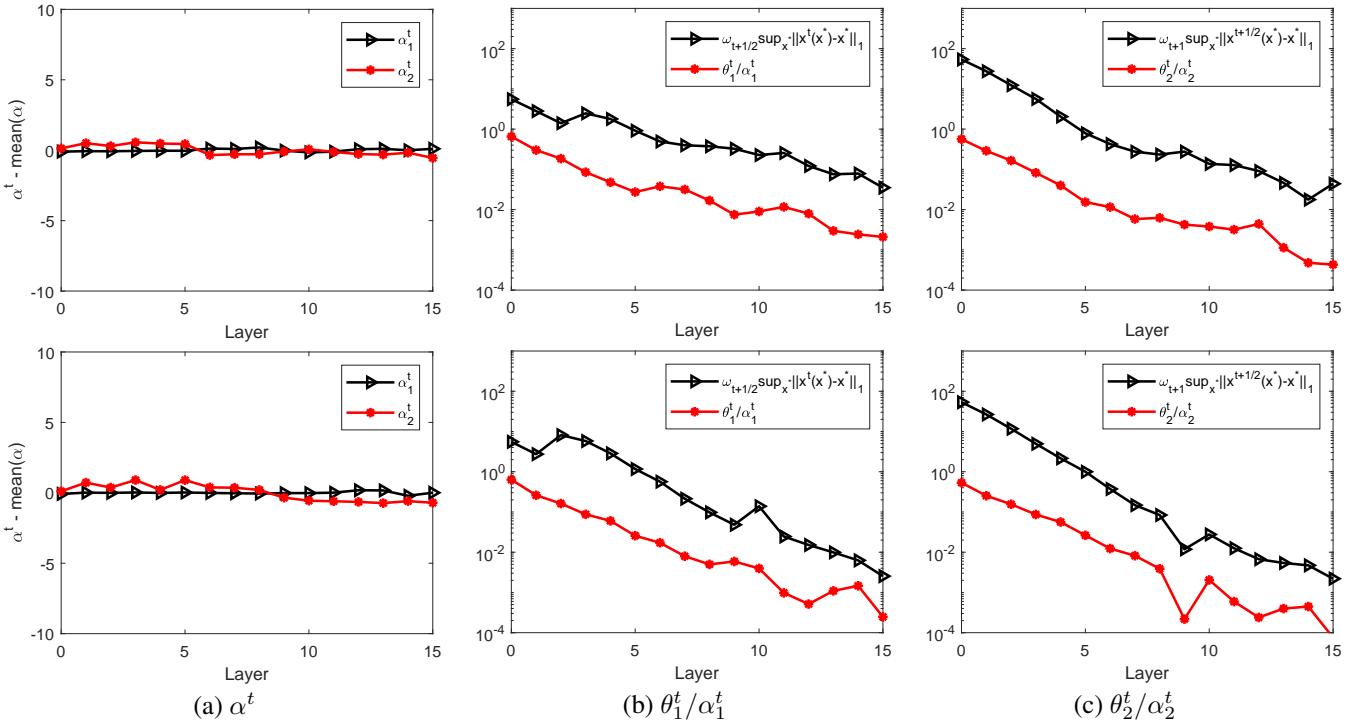


Figure 6: Justification of ELISTA-t (top) and ELISTA (bottom) (noiseless case): parameters obtained by training satisfy (9).

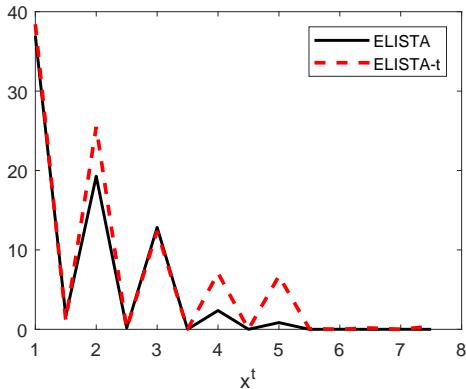


Figure 7: Verification for Support Set Assumption

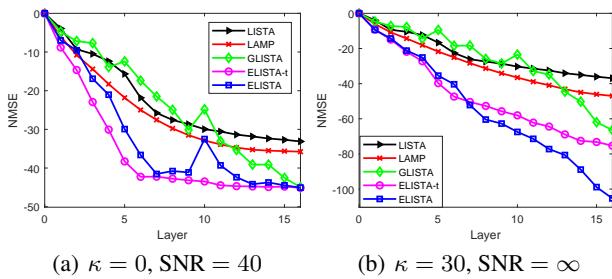


Figure 8: Comparison of sparse representation with different layers under different SNR and κ .

From Figure 19, we can find that ELISTA and ELISTA-t far exceed MaxSparseNet and GFLSTM in terms of strict accuracy, which is the same as our conclusion in Section 4.1: MT tends to recover more sparse representation. In terms of loose accuracy, our methods lag slightly behind GFLSTM which specializes in recovering the support set but still far outperform MaxSparseNet.

• Unsupervised Experiment

Moreover, we compare ELISTA to other learned methods by unsupervised experiments on MNIST as in (Ablin et al. 2019). The result for $\lambda = 0.5$ is shown in Figure 20.

S2. Linear Convergence for ELISTA

In this section, we prove the convergence of our main algorithm, ELISTA.

We first give and prove a key lemma, whose main content is that the no ‘‘false positive’’ assumption holds under the special choices of θ_1^t and θ_2^t . Moreover, this lemma plays an important role in our theoretical analysis.

Lemma 1. *There is no ‘‘false positive’’ when*

$$\theta_1^t = \alpha_1^t \mu(A) \sup_{x^*} \|x^t - x^*\|_1, \quad \theta_2^t = \alpha_2^t \mu(A) \sup_{x^*} \|x^{t+\frac{1}{2}} - x^*\|_1.$$

That is, $\forall t$, $\text{supp}(x^t) \subset S$ and $\text{supp}(x^{t+\frac{1}{2}}) \subset S$, where $S = \text{supp}(x^)$.*

Proof. We use Mathematical Induction to finish this proof, i.e., we assume $\text{supp}(x^t) \subset S$ to prove $\text{supp}(x^{t+\frac{1}{2}}) \subset S$ and $\text{supp}(x^{t+1}) \subset S$. Recall that the update rules of

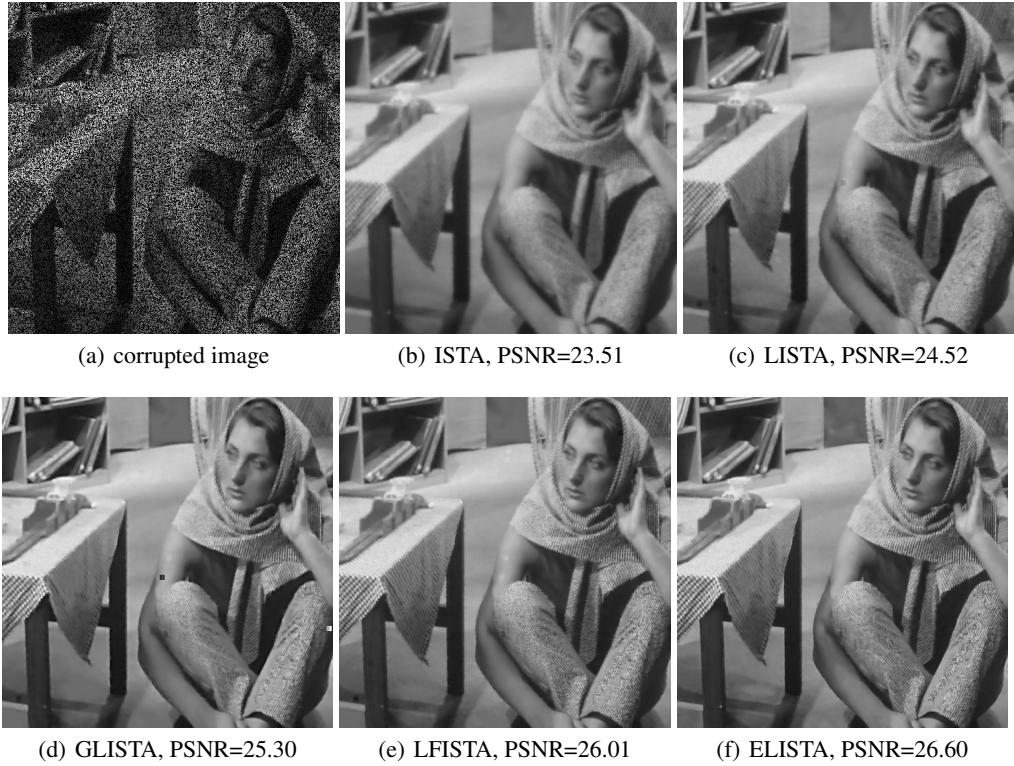


Figure 9: Image inpainting with 50% missing pixels on Barbara.

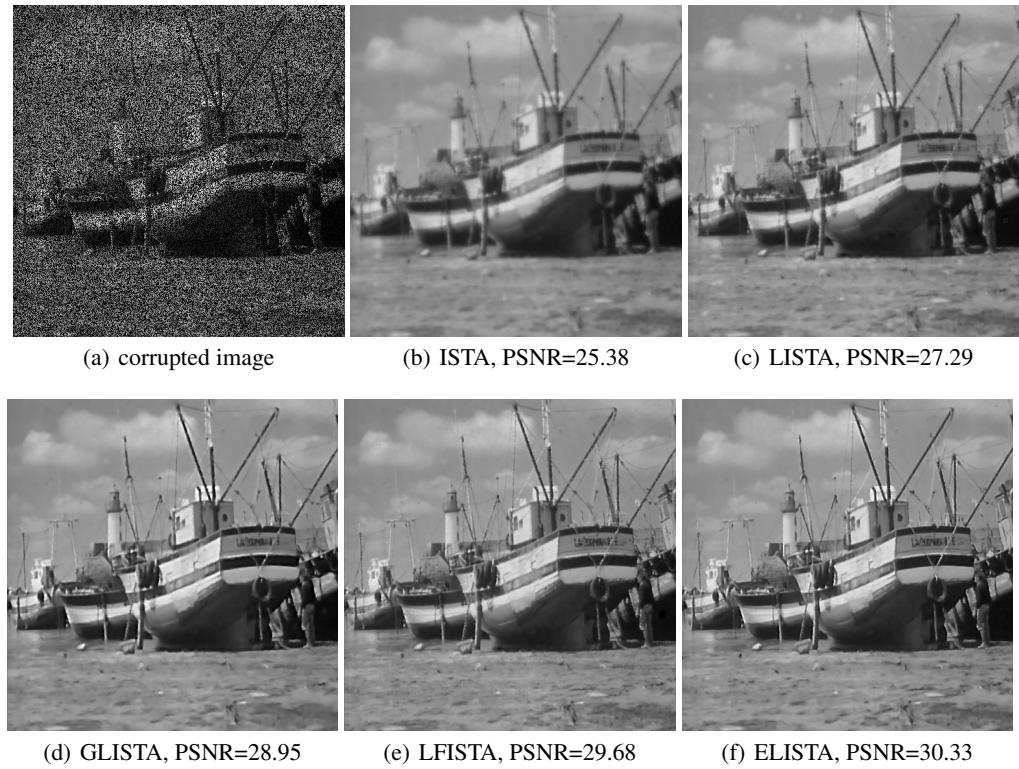


Figure 10: Image inpainting with 50% missing pixels on Boat.

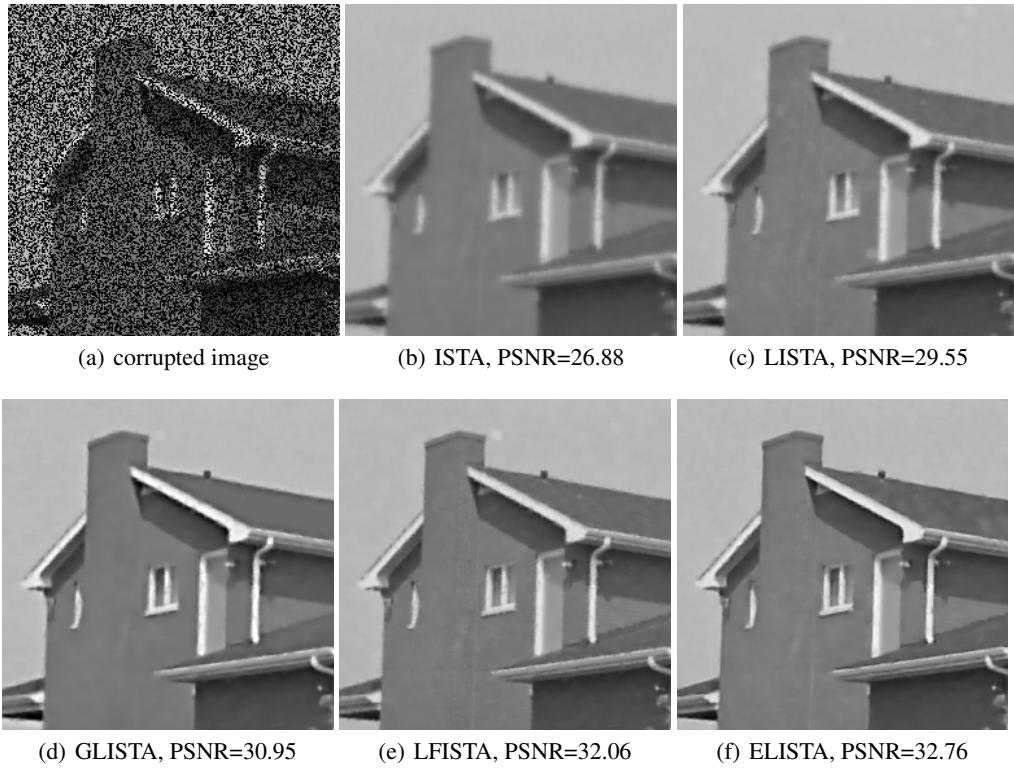


Figure 11: Image inpainting with 50% missing pixels on House.

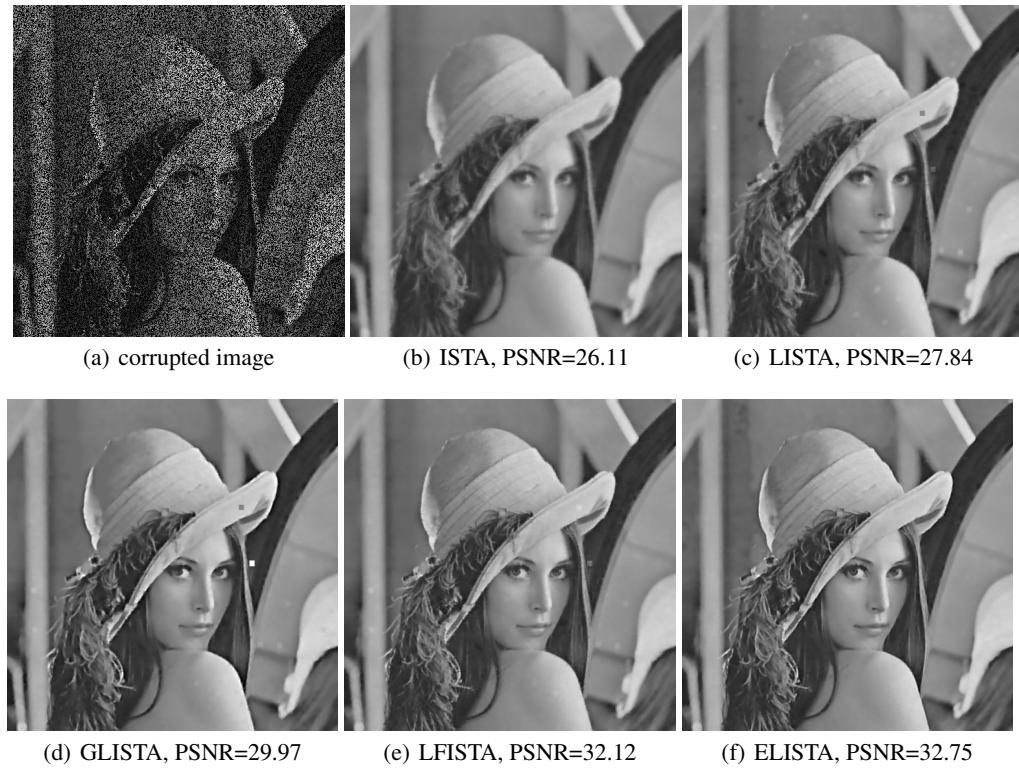


Figure 12: Image inpainting with 50% missing pixels on Lena.

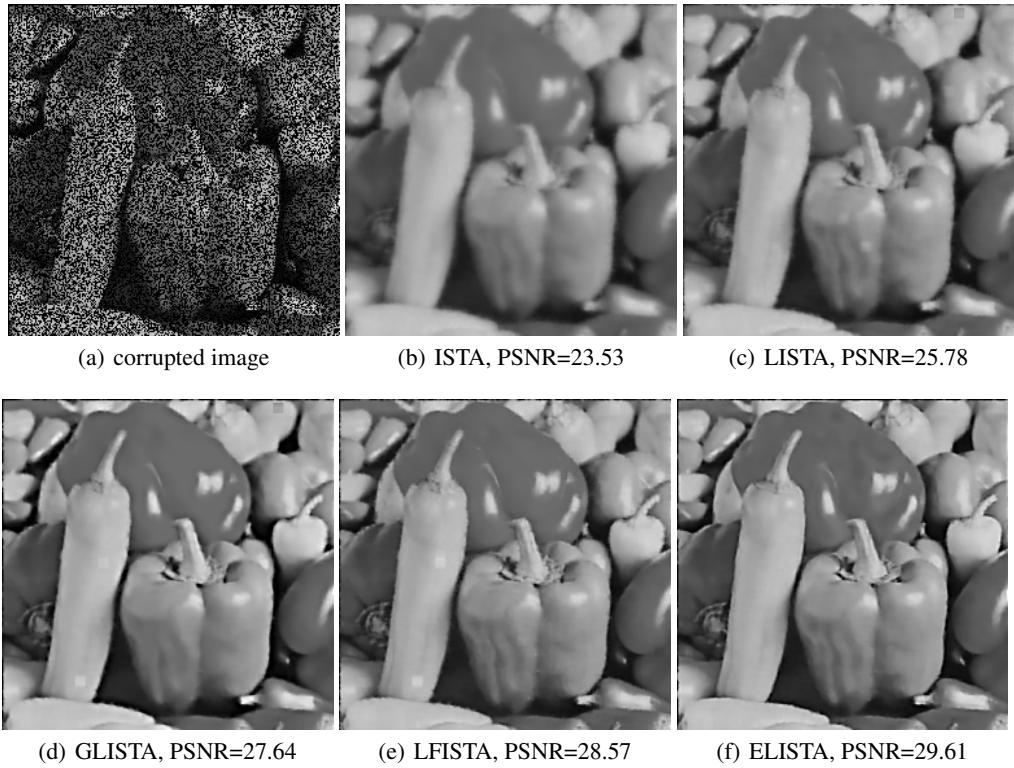


Figure 13: Image inpainting with 50% missing pixels on Peppers.

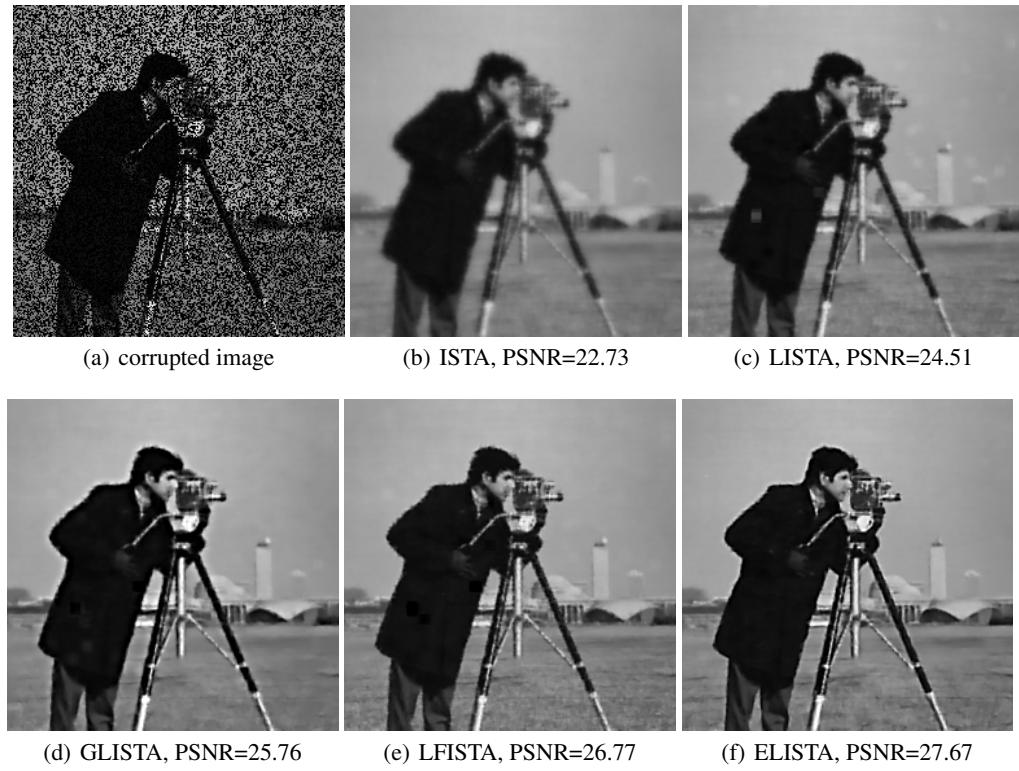


Figure 14: Image inpainting with 50% missing pixels on Cameraman.

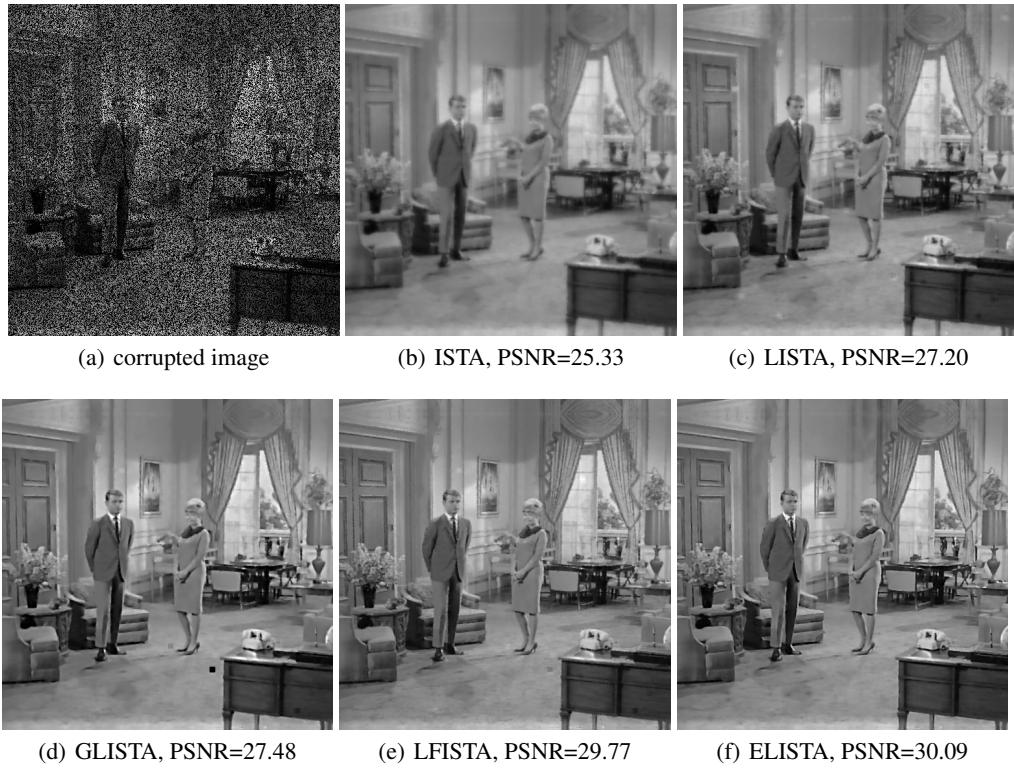


Figure 15: Image inpainting with 50% missing pixels on Couple.

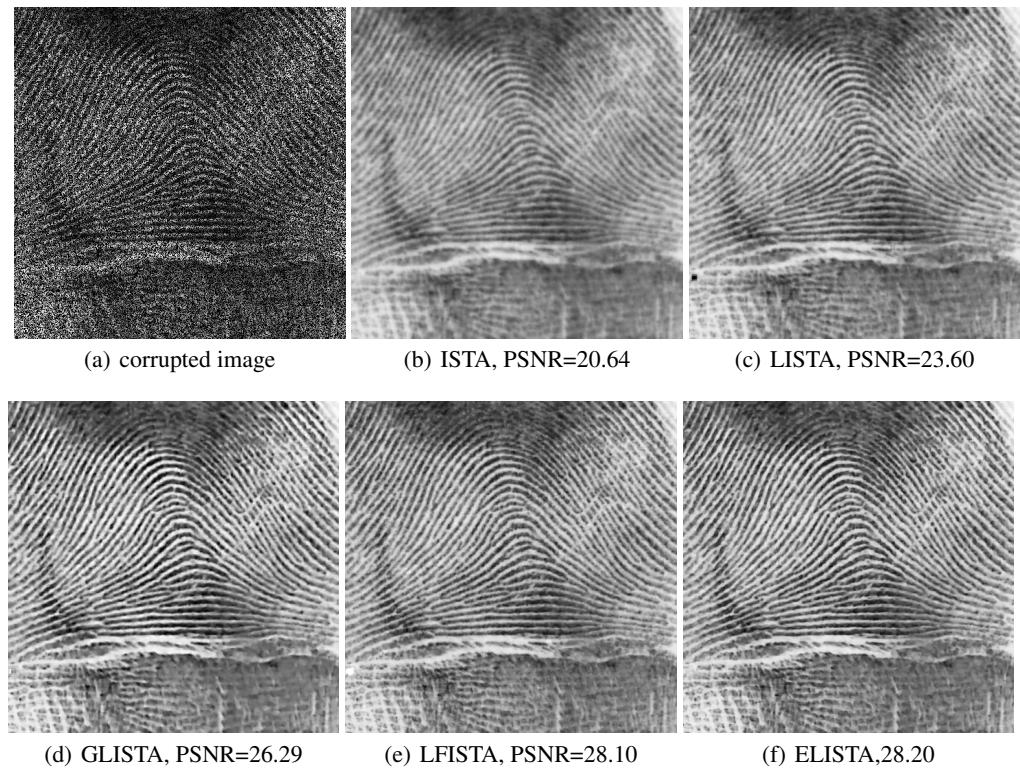


Figure 16: Image inpainting with 50% missing pixels on Fingerprint.

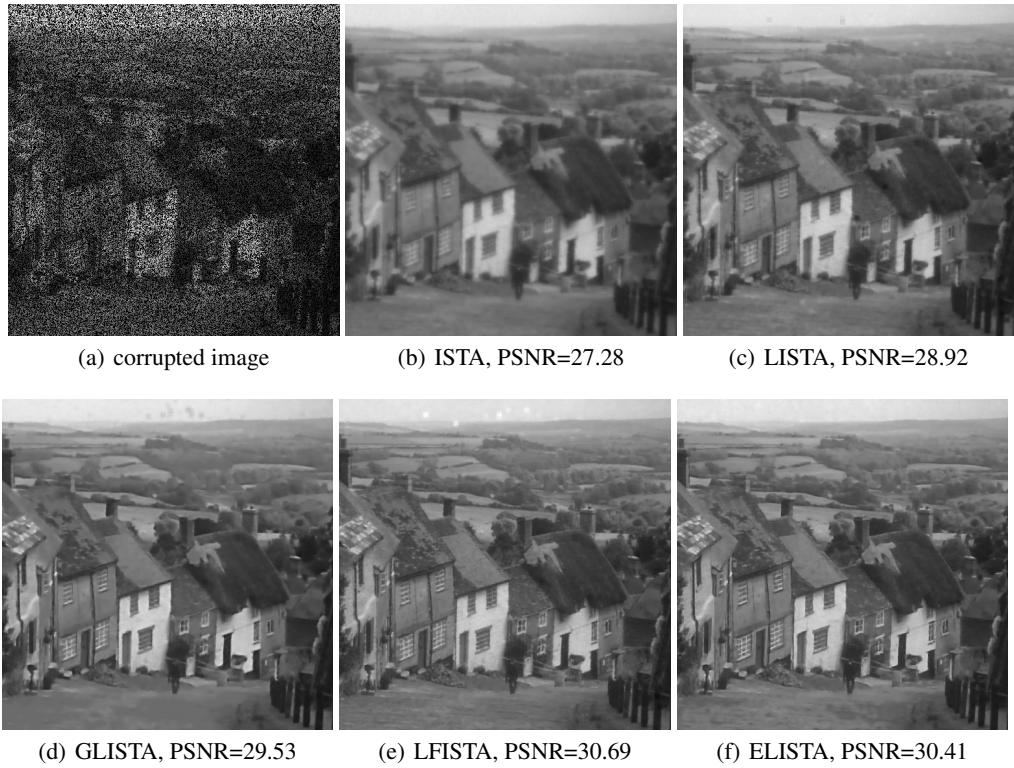


Figure 17: Image inpainting with 50% missing pixels on Hill.

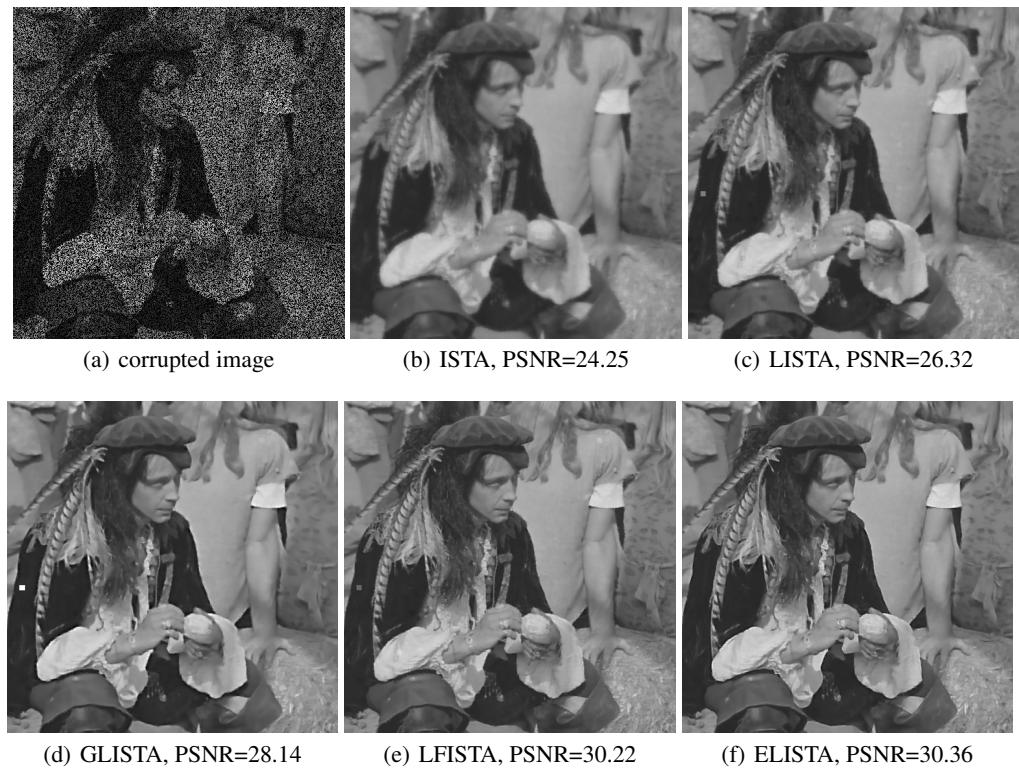


Figure 18: Image inpainting with 50% missing pixels on Man.

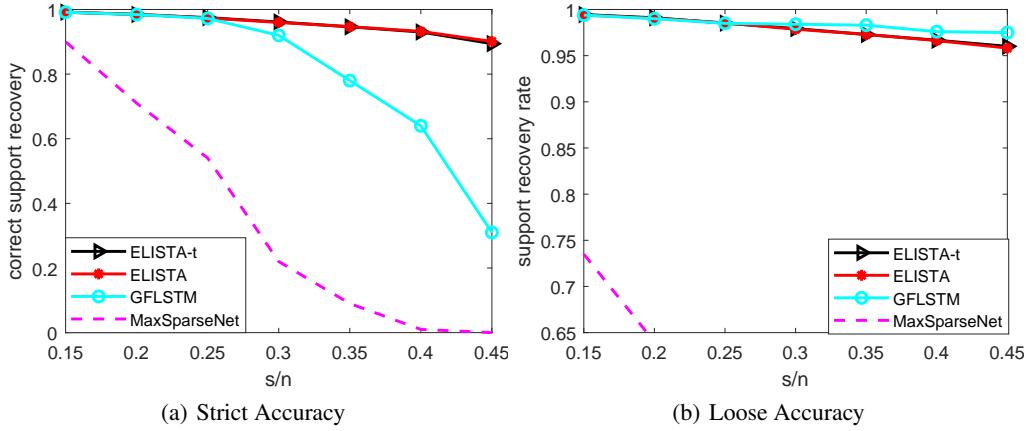


Figure 19: Support set accuracy.

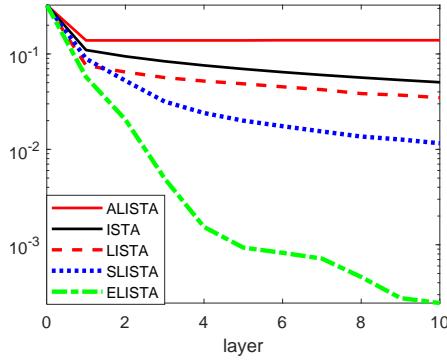


Figure 20: Unsupervised Experiment for $\lambda = 0.5$

ELISTA can be formulated as follows:

$$\begin{cases} x^{t+\frac{1}{2}} = \text{MT}(x^t - \alpha_1^t W(Ax^t - y), \theta_1^t, \bar{\theta}_1^t), \\ x^{t+1} = \text{MT}(x^t - \alpha_2^t W(Ax^{t+\frac{1}{2}} - y), \theta_2^t, \bar{\theta}_2^t). \end{cases} \quad (15)$$

Firstly, we consider the first step of (15). We have

$$\begin{aligned} x_i^{t+\frac{1}{2}} &= \text{MT}(x_i^t - \alpha_1^t (WA(x^t - x^*))_i, \theta_1^t, \bar{\theta}_1^t) \\ &= \text{MT}(-\alpha_1^t \sum_{j \in S} (WA)_{ij} (x_j^t - x_j^*), \theta_1^t, \bar{\theta}_1^t), \quad \forall i \notin S. \end{aligned} \quad (16)$$

Besides, we know $\theta_1^t = \mu(A) \sup_{x^*} \|x^t - x^*\|_1$. Since $W \in \mathcal{W}(A)$, we have $|(WA)_{ij}| \leq \mu(A)$, $\forall i \neq j$. Thus, we obtain

$$\begin{aligned} &\left| -\alpha_1^t \sum_{j \in S} (WA)_{ij} (x_j^t - x_j^*) \right| \\ &\leq \sum_{j \in S} \alpha_1^t \mu(A) |x_j^t - x_j^*| \\ &= \sum_{j \in \text{supp}(x^t)} \alpha_1^t \mu(A) |x_j^t - x_j^*| \\ &= \alpha_1^t \mu(A) \|x^t - x^*\|_1 \leq \theta_1^t \quad \forall i \notin S, \end{aligned}$$

which means that $x_i^{t+\frac{1}{2}} = 0, \forall i \notin S$, according to the definition of the multistage-thresholding function (MT) (6). Thus, we have $\text{supp}(x^{t+\frac{1}{2}}) \subset S$. Next, we consider the second step of (15). We have

$$\begin{aligned} x_i^{t+1} &= \text{MT}(x_i^t - \alpha_2^t (WA(x^{t+\frac{1}{2}} - x^*))_i, \theta_2^t, \bar{\theta}_2^t) \\ &= \text{MT}(-\alpha_2^t (WA(x^{t+\frac{1}{2}} - x^*))_i, \theta_2^t, \bar{\theta}_2^t) \quad \forall i \notin S. \end{aligned} \quad (17)$$

Besides, we have $\theta_2^t = \alpha_2^t \mu(A) \sup_{x^*} \|x^{t+\frac{1}{2}} - x^*\|_1$. Thus,

$$\begin{aligned} &\left| -\alpha_2^t \sum_{j \in S} (WA)_{ij} (x_j^{t+\frac{1}{2}} - x_j^*) \right| \\ &\leq \sum_{j \in S} \alpha_2^t \mu(A) |x_j^{t+\frac{1}{2}} - x_j^*| \\ &= \sum_{j \in \text{supp}(x^{t+\frac{1}{2}})} \alpha_2^t \mu(A) |x_j^{t+\frac{1}{2}} - x_j^*| \\ &= \alpha_2^t \mu(A) \|x^{t+\frac{1}{2}} - x^*\|_1 \leq \theta_2^t \quad \forall i \notin S, \end{aligned}$$

which means $x_i^{t+1} = 0, \forall i \notin S$, according to the definition of MT, i.e., $\text{supp}(x^{t+1}) \subset S$. Moreover, we know $\text{supp}(x^0) = \text{supp}(\mathbf{0}) \subset S$. Thus, we obtain $\forall t$, $\text{supp}(x^t) \subset S$ and $\text{supp}(x^{t+\frac{1}{2}}) \subset S$ when $\theta_1^t = \alpha_1^t \mu(A) \sup_{x^*} \|x^t - x^*\|_1$ and $\theta_2^t = \alpha_2^t \mu(A) \sup_{x^*} \|x^{t+\frac{1}{2}} - x^*\|_1$. Finally, we hence complete the proof of the no ‘‘false positive’’ property. \square

Then we present a key lemma about our MT, which is also an important lemma to prove the convergence result of ELISTA, i.e., Theorem 1.

Lemma 2. For MT, if we define $\tilde{z} = \text{ST}(x, \theta)$, when $0 \leq |\tilde{z}| < \theta$, and $\tilde{z} = x$, when $|\tilde{z}| \geq \bar{\theta}$, then we have

$$z = K \odot \tilde{z},$$

where

$$K_i = \begin{cases} \frac{\bar{\theta}}{\theta - \bar{\theta}} & 0 \leq |\tilde{z}| < \bar{\theta}, \\ 1 & |\tilde{z}| \geq \bar{\theta}. \end{cases}$$

Proof. We know that the definition of ST is

$$\text{ST}(x, \theta) = \begin{cases} 0, & 0 \leq |x| < \theta, \\ \text{sign}(x)(|x| - \theta), & |x| \geq \theta, \end{cases}$$

and the definition of our multistage-thresholding function (MT) is

$$z = \text{MT}(x, \theta, \bar{\theta}) = \begin{cases} 0, & 0 \leq |x| < \theta, \\ \frac{\bar{\theta}}{\theta - \bar{\theta}} \text{sign}(x)(|x| - \theta), & \theta \leq |x| < \bar{\theta}, \\ x, & |x| \geq \bar{\theta}. \end{cases} \quad (18)$$

Then we define $\tilde{z} = \text{ST}(x, \theta)$, when $0 \leq |\tilde{z}| < \bar{\theta}$. Thus, when $0 \leq |x| < \theta$, i.e., $|\tilde{z}| = 0$, $z = \text{MT}(x, \theta, \bar{\theta}) = 0 = \frac{\bar{\theta}}{\theta - \bar{\theta}} \cdot 0 = \frac{\bar{\theta}}{\theta - \bar{\theta}} \tilde{z}$, and when $\theta \leq |x| < \bar{\theta}$, i.e., $0 < |\tilde{z}| < \bar{\theta}$, $z = \frac{\bar{\theta}}{\theta - \bar{\theta}} \text{sign}(x)(|x| - \theta) = \frac{\bar{\theta}}{\theta - \bar{\theta}} \text{ST}(x, \theta) = \frac{\bar{\theta}}{\theta - \bar{\theta}} \tilde{z}$. Besides, when $|\tilde{z}| \geq \bar{\theta}$, i.e., $|x| \geq \bar{\theta}$, we define $\tilde{z} = x$. Obviously, we can get $z = 1 \cdot x = 1 \cdot \tilde{z}$. Therefore, we obtain

$$z = K \odot \tilde{z},$$

where

$$K_i = \begin{cases} \frac{\bar{\theta}}{\theta - \bar{\theta}} & 0 \leq |\tilde{z}| < \bar{\theta}, \\ 1 & |\tilde{z}| \geq \bar{\theta}. \end{cases}$$

□

Next we present a lemma, which mainly aims at the first step of each layer of ELISTA, gives the relationship between x^t and $x^{t+\frac{1}{2}}$, and is also a basis for proving Theorem 1 which shows the convergence result of ELISTA.

Lemma 3. *For the first step in the update rules of ELISTA, if*

$$\theta_1^t = \alpha_1^t \omega_{t+\frac{1}{2}}(k_{t+\frac{1}{2}}|\Theta)\mu(A) \sup_{x^*} \|x^t - x^*\|_1$$

is achieved, then we have

$$\sup_{x^*} \|x^{t+\frac{1}{2}} - x^*\|_1 \leq \beta_2^t \sup_{x^*} \|x^t - x^*\|_1,$$

where

$$\begin{aligned} \beta_2^t &= \sup_{x^*} ((|S| + |S^{(t+\frac{1}{2})}|) + (|S| - |\bar{S}^{(t+\frac{1}{2})}|) \\ &\quad + (|\bar{S}^{(t+\frac{1}{2})}| - |S^{(t+\frac{1}{2})}| - |U_1^{(t+\frac{1}{2})}| \\ &\quad + |U_2^{(t+\frac{1}{2})}|)(1 - \gamma^{t+\frac{1}{2}}))\omega_{t+\frac{1}{2}}(k_{t+\frac{1}{2}}|\Theta))\alpha_1^t \mu(A) \\ &\quad + |1 - \alpha_1^t|. \end{aligned}$$

The symbols in this lemma are all defined and explained in the process of the following proof.

Proof. The update rules of ELISTA is formulated as follows:

$$\begin{cases} x^{t+\frac{1}{2}} = \text{MT}(x^t - \alpha_1^t W(Ax^t - y), \theta_1^t, \bar{\theta}_1^t), \\ x^{t+1} = \text{MT}(x^t - \alpha_2^t W(Ax^{t+\frac{1}{2}} - y), \theta_2^t, \bar{\theta}_2^t). \end{cases} \quad (19)$$

For the first step in (19), we define $\tilde{x}_i^{t+\frac{1}{2}} = \text{ST}(x_i^t - (\alpha_1^t W(Ax^t - y))_i, \theta_1^t)$, when $0 \leq |\tilde{x}_i^{t+\frac{1}{2}}| < \bar{\theta}_1^t$, and $\tilde{x}_i^{t+\frac{1}{2}} =$

$x_i^t - (\alpha_1^t W(Ax^t - y))_i$, when $|\tilde{x}_i^{t+\frac{1}{2}}| \geq \bar{\theta}_1^t$. Then due to Lemma 2, we have

$$x^{t+\frac{1}{2}} = K^{t+\frac{1}{2}} \odot \tilde{x}^{t+\frac{1}{2}},$$

where

$$K_i^{t+\frac{1}{2}} = \begin{cases} \frac{\bar{\theta}_1^t}{\bar{\theta}_1^t - \theta_1^t} & 0 \leq |\tilde{x}_i^{t+\frac{1}{2}}| < \bar{\theta}_1^t, \\ 1 & |\tilde{x}_i^{t+\frac{1}{2}}| \geq \bar{\theta}_1^t. \end{cases}$$

Similarly, we can get $x^t = K^t \odot \tilde{x}^t$. Then we first consider the case of $0 \leq |\tilde{x}_i^{t+\frac{1}{2}}| < \bar{\theta}_1^t$, and have

$$\tilde{x}_i^{t+\frac{1}{2}} = ((I - \alpha_1^t WA)(K^t \odot \tilde{x}^t - x^*))_i + x_i^* - \theta_1^t \partial \ell_1(\tilde{x}_i^{t+\frac{1}{2}}).$$

Thus, we obtain

$$\begin{aligned} K_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}} &= (1 + B_i^{t+\frac{1}{2}}) \tilde{x}_i^{t+\frac{1}{2}} \\ &= ((I - \alpha_1^t WA)(K^t \odot \tilde{x}^t - x^*))_i + x_i^* \\ &\quad - \theta_1^t \partial \ell_1(\tilde{x}_i^{t+\frac{1}{2}}) + B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}}, \end{aligned}$$

where $B_i^{t+\frac{1}{2}} = \frac{\theta_1^t}{\bar{\theta}_1^t - \theta_1^t}$. Thus, we have

$$\begin{aligned} x_i^{t+\frac{1}{2}} - x_i^* &= ((I - \alpha_1^t WA)(x^t - x^*))_i \\ &\quad - \theta_1^t \partial \ell_1(\tilde{x}_i^{t+\frac{1}{2}}) + B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}}. \end{aligned} \quad (20)$$

To calculate all non-zero elements of $x^{t+\frac{1}{2}} - x^*$, we divide them into three parts: $i \in \bar{S}^{(t+\frac{1}{2})}$, $i \in S \setminus \bar{S}^{(t+\frac{1}{2})}$

and $i \in S^{(t+\frac{1}{2})}$, where $\bar{S}^{(t+\frac{1}{2})} \triangleq S \cap \text{supp}(x^{t+\frac{1}{2}})$ and

$S^{(t+\frac{1}{2})} \triangleq \{i | i \in \text{supp}(x^{t+\frac{1}{2}}), i \notin S\}$.

For $i \in \bar{S}^{(t+\frac{1}{2})}$, we know $\tilde{x}_i^{t+\frac{1}{2}} \neq 0$ and $x_i^* \neq 0$. Thus,

$\partial \ell_1(\tilde{x}_i^{t+\frac{1}{2}}) = \text{sign}(\tilde{x}_i^{t+\frac{1}{2}})$. Then we have

$$\begin{aligned} x_i^{t+\frac{1}{2}} - x_i^* &= ((I - \alpha_1^t WA)(x^t - x^*))_i - \theta_1^t \text{sign}(\tilde{x}_i^{t+\frac{1}{2}}) \\ &\quad + |\tilde{x}_i^{t+\frac{1}{2}}| B_i^{t+\frac{1}{2}} \text{sign}(\tilde{x}_i^{t+\frac{1}{2}}), \end{aligned}$$

which implies that

$$\begin{aligned} |x_i^{t+\frac{1}{2}} - x_i^*| &\leq |((I - \alpha_1^t WA)(x^t - x^*))_i| + |\tilde{x}_i^{t+\frac{1}{2}}| B_i^{t+\frac{1}{2}} - \theta_1^t. \end{aligned}$$

We assume $\bar{\theta}_1^t \geq \theta_1^t + |\tilde{x}_i^{t+\frac{1}{2}}|$. Then we can obtain

$$\begin{aligned} ||\tilde{x}_i^{t+\frac{1}{2}}| B_i^{t+\frac{1}{2}} - \theta_1^t| &= \theta_1^t - |\tilde{x}_i^{t+\frac{1}{2}}| B_i^{t+\frac{1}{2}} \\ &= (1 - \frac{|\tilde{x}_i^{t+\frac{1}{2}}|}{\theta_1^t} B_i^{t+\frac{1}{2}}) \theta_1^t \leq \theta_1^t. \end{aligned}$$

In addition, we define $\gamma^{t+\frac{1}{2}} \triangleq \frac{|\tilde{x}_i^{t+\frac{1}{2}}|}{\theta_1^t} B_i^{t+\frac{1}{2}}$. From the previous assumption, we know $\gamma^{t+\frac{1}{2}} \leq 1$. Thus, we have

$$|x_i^{t+\frac{1}{2}} - x_i^*| \leq |((I - \alpha_1^t WA)(x^t - x^*))_i| + (1 - \gamma^{t+\frac{1}{2}}) \theta_1^t.$$

For $i \in S \setminus \bar{S}^{(t+\frac{1}{2})}$, we know $\tilde{x}_i^{t+\frac{1}{2}} = 0$ and $x_i^* \neq 0$. Thus,

$$x_i^{t+\frac{1}{2}} - x_i^* = ((I - \alpha_1^t WA)(x^t - x^*))_i - \theta_1^t \partial \ell_1(\tilde{x}_i^{t+\frac{1}{2}}),$$

which implies that

$$|x_i^{t+\frac{1}{2}} - x_i^*| \leq |((I - \alpha_1^t WA)(x^t - x^*))_i| + \theta_1^t.$$

For $i \in S^{(t+\frac{1}{2})}$, we know $\tilde{x}_i^{t+\frac{1}{2}} \neq 0$ and $x_i^* = 0$. Then we have

$$\begin{aligned} x_i^{t+\frac{1}{2}} - x_i^* &= x_i^{t+\frac{1}{2}} \\ &= ((I - \alpha_1^t WA)(x^t - x^*))_i \\ &\quad - (\theta_1^t - |\tilde{x}_i^{t+\frac{1}{2}}|B_i^{t+\frac{1}{2}})\text{sign}(\tilde{x}_i^{t+\frac{1}{2}}). \end{aligned} \quad (21)$$

Due to $K_i^{t+\frac{1}{2}} \geq 1$, we know that $x_i^{t+\frac{1}{2}}$ has the same sign as $\tilde{x}_i^{t+\frac{1}{2}}$. Thus, we multiply both sides of (21) by $\text{sign}(\tilde{x}_i^{t+\frac{1}{2}})$, and thus we have

$$\begin{aligned} |x_i^{t+\frac{1}{2}} - x_i^*| &= |((I - \alpha_1^t WA)(x^t - x^*))_i| \\ &\quad \text{sign}(\tilde{x}_i^{t+\frac{1}{2}}) - (\theta_1^t - |\tilde{x}_i^{t+\frac{1}{2}}|B_i^{t+\frac{1}{2}}), \end{aligned}$$

i.e.,

$$\begin{aligned} |x_i^{t+\frac{1}{2}} - x_i^*| &+ (\theta_1^t - |\tilde{x}_i^{t+\frac{1}{2}}|B_i^{t+\frac{1}{2}}) \\ &= |((I - \alpha_1^t WA)(x^t - x^*))_i| \text{sign}(\tilde{x}_i^{t+\frac{1}{2}}). \end{aligned}$$

We know $|x_i^{t+\frac{1}{2}} - x_i^*| \geq 0$. Besides, we can obtain $\theta_1^t - |\tilde{x}_i^{t+\frac{1}{2}}|B_i^{t+\frac{1}{2}} \geq 0$ from the assumption $\bar{\theta}_1^t \geq \theta_1^t + |\tilde{x}_i^{t+\frac{1}{2}}|$. Thus, we know that the sign of $((I - \alpha_1^t WA)(x^t - x^*))_i$ is the same as that of $\text{sign}(\tilde{x}_i^{t+\frac{1}{2}})$. That is,

$$\begin{aligned} &|((I - \alpha_1^t WA)(x^t - x^*))_i| \\ &= |((I - \alpha_1^t WA)(x^t - x^*))_i| \text{sign}(\tilde{x}_i^{t+\frac{1}{2}}). \end{aligned} \quad (22)$$

Substituting (22) into (21), we get

$$\begin{aligned} x_i^{t+\frac{1}{2}} - x_i^* &= \text{sign}(\tilde{x}_i^{t+\frac{1}{2}}) |((I - \alpha_1^t WA)(x^t - x^*))_i| \\ &\quad - (1 - \gamma^{t+\frac{1}{2}})\theta_1^t. \end{aligned} \quad (23)$$

Multiplying both sides of (23) by $\text{sign}(\tilde{x}_i^{t+\frac{1}{2}})$, we obtain

$$|x_i^{t+\frac{1}{2}} - x_i^*| = |((I - \alpha_1^t WA)(x^t - x^*))_i| - (1 - \gamma^{t+\frac{1}{2}})\theta_1^t.$$

For the case of $|\tilde{x}_i^{t+\frac{1}{2}}| \geq \bar{\theta}_1^t$, we can get

$$\begin{aligned} x_i^{t+\frac{1}{2}} &= K_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}} = \tilde{x}_i^{t+\frac{1}{2}} \\ &= ((I - \alpha_1^t WA)(K^t \odot \tilde{x}^t - x^*))_i + x_i^*. \end{aligned}$$

Thus, we have

$$|x_i^{t+\frac{1}{2}} - x_i^*| = |((I - \alpha_1^t WA)(x^t - x^*))_i|.$$

We assume $\text{supp}(x^t) \subset \text{supp}(x^{t+\frac{1}{2}}) \subset \text{supp}(x^{t+1})$. Then by summing up all $|x_i^{t+\frac{1}{2}} - x_i^*|$ which satisfy $i \in$

$\text{supp}(\tilde{x}^{t+\frac{1}{2}}) \cup \text{supp}(x^*) = S + S^{(t+\frac{1}{2})} = \bar{S}^{(t+\frac{1}{2})} + S \setminus \bar{S}^{(t+\frac{1}{2})} + S^{(t+\frac{1}{2})}$ in the two cases, we can obtain

$$\begin{aligned} \|x^{t+\frac{1}{2}} - x^*\|_1 &= \sum_{i \in S + S^{(t+\frac{1}{2})}} |x_i^{t+\frac{1}{2}} - x_i^*| \\ &\leq \sum_{i \in S + S^{(t+\frac{1}{2})}} |((I - \alpha_1^t WA)(x^t - x^*))_i| \\ &\quad + (|S| - |\bar{S}^{(t+\frac{1}{2})}|)\theta_1^t + (|\bar{S}^{(t+\frac{1}{2})}| - |S^{(t+\frac{1}{2})}| \\ &\quad - |U_1^{(t+\frac{1}{2})}| + |U_2^{(t+\frac{1}{2})}|)(1 - \gamma^{t+\frac{1}{2}})\theta_1^t \\ &\leq \sum_{i \in S + S^{(t+\frac{1}{2})}} \left(\sum_{j \in \text{supp}(x^t) + S} \alpha_1^t |(I - WA)_{ij}(x_j^t - x_j^*)| \right. \\ &\quad \left. + |1 - \alpha_1^t| |x_i^t - x_i^*| \right) + (|\bar{S}^{(t+\frac{1}{2})}| - |S^{(t+\frac{1}{2})}| - |U_1^{(t+\frac{1}{2})}| \\ &\quad + |U_2^{(t+\frac{1}{2})}|)(1 - \gamma^{t+\frac{1}{2}})\theta_1^t + (|S| - |\bar{S}^{(t+\frac{1}{2})}|)\theta_1^t \\ &\leq (|S| + |S^{(t+\frac{1}{2})}|)\alpha_1^t \mu(A) \|x^t - x^*\|_1 \\ &\quad + |1 - \alpha_1^t| \|x^t - x^*\|_1 + (|S| - |\bar{S}^{(t+\frac{1}{2})}|)\theta_1^t \\ &\quad + (|\bar{S}^{(t+\frac{1}{2})}| - |S^{(t+\frac{1}{2})}| - |U_1^{(t+\frac{1}{2})}| + |U_2^{(t+\frac{1}{2})}|) \\ &\quad (1 - \gamma^{t+\frac{1}{2}})\theta_1^t, \end{aligned} \quad (24)$$

where

$$\begin{aligned} U_1^{(t+\frac{1}{2})} &\triangleq \{i | i \in \bar{S}^{(t+\frac{1}{2})}, |\tilde{x}_i^{t+\frac{1}{2}}| \geq \bar{\theta}_1^t\}, \\ U_2^{(t+\frac{1}{2})} &\triangleq \{i | i \in S^{(t+\frac{1}{2})}, |\tilde{x}_i^{t+\frac{1}{2}}| \geq \bar{\theta}_1^t\}. \end{aligned} \quad (25)$$

Then we set

$$\theta_1^t = \alpha_1^t \mu(A) \omega_{t+\frac{1}{2}}(k_{t+\frac{1}{2}}|\Theta) \sup_{x^*} \|x^t - x^*\|_1, \quad (26)$$

and take supremum of both sides of (24). Then we obtain

$$\begin{aligned} &\sup_{x^*} \|x^{t+\frac{1}{2}} - x^*\|_1 \\ &\leq \sup_{x^*} (|S| + |S^{(t+\frac{1}{2})}| + (|S| - |\bar{S}^{(t+\frac{1}{2})}|) \\ &\quad + (|\bar{S}^{(t+\frac{1}{2})}| - |S^{(t+\frac{1}{2})}| - |U_1^{(t+\frac{1}{2})}| \\ &\quad + |U_2^{(t+\frac{1}{2})}|)(1 - \gamma^{t+\frac{1}{2}})) \omega_{t+\frac{1}{2}}(k_{t+\frac{1}{2}}|\Theta) \alpha_1^t \mu(A) \\ &\quad + |1 - \alpha_1^t| \sup_{x^*} \|x^t - x^*\|_1. \end{aligned} \quad (27)$$

□ 838

Next, using Lemma 3, we prove another important lemma, which mainly considers the second step of each iteration of ELISTA, gives the relationship between x^t and x^{t+1} , and provides the basis for proving Theorem 1.

Lemma 4. For the second step in the update rules of ELISTA, if

$$\theta_1^t = \alpha_1^t \omega_{t+\frac{1}{2}}(k_{t+\frac{1}{2}}|\Theta) \mu(A) \sup_{x^*} \|x^t - x^*\|_1$$

and

$$\theta_2^t = \alpha_2^t \omega_{t+1}(k_{t+1}|\Theta) \mu(A) \sup_{x^*} \|x^{t+\frac{1}{2}} - x^*\|_1$$

839

840

841

842

are achieved, then we have

$$\sup_{x^*} \|x^{t+1} - x^*\|_1 \leq \exp(c'_{t+1}) \sup_{x^*} \|x^t - x^*\|_1,$$

where,

$$\begin{aligned} c'_{t+1} = & \sup_{x^*} \log(((|S| + |S^{(t+1)}|) - 1 + (1 + \gamma^{t+\frac{1}{2}}) \\ & (|S| + |S^{(t+1)}|)\omega_{t+\frac{1}{2}}(k_{t+\frac{1}{2}}|\Theta|)\alpha_1^t\mu(A) + \alpha_1^t) \\ & + ((|S| + |S^{(t+1)}| + (|S| - |\bar{S}^{(t+1)}|) \\ & + (|\bar{S}^{(t+1)}| - |S^{(t+1)}| - |U_1^{(t+1)}| + |U_2^{(t+1)}|) \\ & (1 - \gamma^{t+1})\omega_{t+1}(k_{t+1}|\Theta|)\alpha_2^t\mu(A) + |1 - \alpha_2^t|) \\ & ((|S| + |S^{(t+\frac{1}{2})}|) + (|S| - |\bar{S}^{(t+\frac{1}{2})}|) \\ & + (|\bar{S}^{(t+\frac{1}{2})}| - |S^{(t+\frac{1}{2})}| - |U_1^{(t+\frac{1}{2})}| + |U_2^{(t+\frac{1}{2})}|) \\ & (1 - \gamma^{t+\frac{1}{2}})\omega_{t+\frac{1}{2}}(k_{t+\frac{1}{2}}|\Theta|)\alpha_1^t\mu(A) + |1 - \alpha_1^t|)). \end{aligned}$$

843 All the symbols in this lemma are all defined and explained in the process of its proof.
844

Proof. We consider the second step of (19). Similar to $\tilde{x}^{t+\frac{1}{2}}$, we also define \tilde{x}^{t+1} . Firstly, for the case of $0 \leq |\tilde{x}_i^{t+1}| < \bar{\theta}_2^t$, we have

$$\begin{aligned} \tilde{x}_i^{t+1} = & \text{ST}((x_i^t - x_i^*) - (x_i^{t+\frac{1}{2}} - x_i^*)) \\ & + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i + x_i^*, \theta_2^t) \\ = & (x_i^t - x_i^*) - (x_i^{t+\frac{1}{2}} - x_i^*) \\ & + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i + x_i^* - \theta_2^t \partial \ell_1(\tilde{x}_i^{t+1}) \\ = & (x_i^t - x_i^*) - ((I - \alpha_1^t WA)(x^t - x^*))_i \\ & + \theta_1^t \partial \ell(x_i^{t+\frac{1}{2}}) - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}} \\ & + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i - \theta_2^t \partial \ell(\tilde{x}_i^{t+1}) + x_i^*. \end{aligned}$$

The last equality holds due to (20). Thus,

$$\begin{aligned} x_i^{t+1} - x_i^* = & K_i^{t+1} \tilde{x}_i^{t+1} - x_i^* \\ = & (\alpha_1^t WA(x^t - x^*))_i + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i \\ & + \theta_1^t \partial \ell(\tilde{x}_i^{t+\frac{1}{2}}) - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}} - \theta_2^t \partial \ell(\tilde{x}_i^{t+1}) + B_i^{t+1} \tilde{x}_i^{t+1}. \end{aligned}$$

845 We also divide all non-zero elements of $x^{t+1} - x^*$ into three
846 parts: $i \in \bar{S}^{t+1}$, $i \in S \setminus \bar{S}^{t+1}$ and $i \in S^{(t+1)}$, where
847 the definitions of $\bar{S}^{(t+1)}$ and $S^{(t+1)}$ are similar to those of
848 $\bar{S}^{(t+\frac{1}{2})}$ and $S^{(t+\frac{1}{2})}$, respectively.

For $i \in \bar{S}^{(t+1)}$, we know $\tilde{x}_i^{t+1} \neq 0$ and $x_i^* \neq 0$. Thus,
 $\partial \ell_1(\tilde{x}_i^{t+1}) = \text{sign}(\tilde{x}_i^{t+1})$. Then we have

$$\begin{aligned} x_i^{t+1} - x_i^* = & (\alpha_1^t WA(x^t - x^*))_i \\ & + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i \\ & + \theta_1^t \partial \ell(\tilde{x}_i^{t+\frac{1}{2}}) - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}} - \theta_2^t \text{sign}(\tilde{x}_i^{t+1}) \\ & + |\tilde{x}_i^{t+1}| B_i^{t+1} \text{sign}(\tilde{x}_i^{t+1}), \end{aligned}$$

which means that

$$\begin{aligned} |x_i^{t+1} - x_i^*| \leq & |(\alpha_1^t WA(x^t - x^*))_i + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i \\ & + \theta_1^t \partial \ell(\tilde{x}_i^{t+\frac{1}{2}}) - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}}| + |\tilde{x}_i^{t+1}| B_i^{t+1} - \theta_2^t|. \end{aligned}$$

We assume $\bar{\theta}_2^t \geq \theta_2^t + |\tilde{x}_i^{t+1}|$. Then we can obtain

$$\begin{aligned} ||\tilde{x}_i^{t+1}| B_i^{t+1} - \theta_2^t| &= \theta_2^t - |\tilde{x}_i^{t+1}| B_i^{t+1} \\ &= (1 - \frac{|\tilde{x}_i^{t+1}|}{\theta_2^t} B_i^{t+1}) \theta_2^t \leq \theta_2^t. \end{aligned}$$

In addition, we define $\gamma^{t+1} \triangleq \frac{|\tilde{x}_i^{t+1}|}{\theta_2^t} B_i^{t+1}$. From the previous assumption, we know $\gamma^{t+1} \leq 1$. Thus, we have

$$\begin{aligned} |x_i^{t+1} - x_i^*| \leq & |(\alpha_1^t WA(x^t - x^*))_i \\ & + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i \\ & + \theta_1^t \partial \ell(\tilde{x}_i^{t+\frac{1}{2}}) - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}}| + (1 - \gamma^{t+1}) \theta_2^t. \end{aligned}$$

For $i \in S \setminus \bar{S}^{(t+1)}$, we know $\tilde{x}_i^{t+1} = 0$ and $x_i^* \neq 0$. Thus,

$$\begin{aligned} x_i^{t+1} - x_i^* = & (\alpha_1^t WA(x^t - x^*))_i \\ & + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i \\ & + \theta_1^t \partial \ell(\tilde{x}_i^{t+\frac{1}{2}}) - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}} - \theta_2^t \partial \ell_1(\tilde{x}_i^{t+1}), \end{aligned}$$

which implies that

$$\begin{aligned} |x_i^{t+1} - x_i^*| \leq & |(\alpha_1^t WA(x^t - x^*))_i \\ & + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i \\ & + \theta_1^t \partial \ell(\tilde{x}_i^{t+\frac{1}{2}}) - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}}| + \theta_2^t. \end{aligned}$$

For $i \in S^{(t+1)}$, we know $\tilde{x}_i^{t+1} \neq 0$ and $x_i^* = 0$. Then we
849 have
850

$$\begin{aligned} x_i^{t+1} - x_i^* &= x_i^{t+1} \\ &= (\alpha_1^t WA(x^t - x^*))_i \\ &+ ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i + \theta_1^t \partial \ell(\tilde{x}_i^{t+\frac{1}{2}}) \\ &- B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}} - (\theta_2^t - |\tilde{x}_i^{t+1}| B_i^{t+1}) \text{sign}(\tilde{x}_i^{t+1}). \end{aligned} \quad (28)$$

Due to $K_i^{t+1} \geq 1$, we know that x_i^{t+1} has the same sign as \tilde{x}_i^{t+1} . Thus, we multiply both sides of (28) by $\text{sign}(\tilde{x}_i^{t+1})$, then we have

$$\begin{aligned} |x_i^{t+1} - x_i^*| &= |x_i^{t+1}| \\ &= ((\alpha_1^t WA(x^t - x^*))_i + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i \\ &+ \theta_1^t \partial \ell(\tilde{x}_i^{t+\frac{1}{2}}) - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}}) \text{sign}(\tilde{x}_i^{t+1}) \\ &- (\theta_2^t - |\tilde{x}_i^{t+1}| B_i^{t+1}), \end{aligned}$$

i.e.,

$$\begin{aligned} |x_i^{t+1} - x_i^*| &+ (\theta_2^t - |\tilde{x}_i^{t+1}| B_i^{t+1}) \\ &= ((\alpha_1^t WA(x^t - x^*))_i + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i \\ &+ \theta_1^t \partial \ell(\tilde{x}_i^{t+\frac{1}{2}}) - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}}) \text{sign}(\tilde{x}_i^{t+1}). \end{aligned}$$

We know $|x_i^{t+1} - x_i^*| \geq 0$. Besides, we can obtain $\theta_2^t - |\tilde{x}_i^{t+1}| B_i^{t+1} \geq 0$ from the assumption $\bar{\theta}_2^t \geq \theta_2^t + |\tilde{x}_i^{t+1}|$. Thus, we know that the sign of $(\alpha_1^t WA(x^t - x^*))_i + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i + \theta_1^t \partial \ell(\tilde{x}_i^{t+\frac{1}{2}}) - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}}$.

854 $\alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*)_i + \theta_1^t \partial\ell(\tilde{x}_i^{t+\frac{1}{2}}) - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}}$ is the
 855 same as that of $\text{sign}(\tilde{x}_i^{t+1})$. That is,

$$\begin{aligned} & (\alpha_1^t WA(x^t - x^*))_i + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i \\ & + \theta_1^t \partial\ell(\tilde{x}_i^{t+\frac{1}{2}}) - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}} \\ & = |(\alpha_1^t WA(x^t - x^*))_i + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i| \\ & + \theta_1^t \partial\ell(\tilde{x}_i^{t+\frac{1}{2}}) - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}} |\text{sign}(\tilde{x}_i^{t+1})|. \end{aligned} \quad (29)$$

856 Substituting (29) into (28), we obtain

$$\begin{aligned} & x_i^{t+1} - x_i^* \\ & = \text{sign}(\tilde{x}_i^{t+1}) (|(\alpha_1^t WA(x^t - x^*))_i \\ & + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i + \theta_1^t \partial\ell(\tilde{x}_i^{t+\frac{1}{2}}) \\ & - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}}| - (1 - \gamma^{t+1})\theta_2^t). \end{aligned} \quad (30)$$

Then by multiplying both sides of (30) by $\text{sign}(\tilde{x}_i^{t+1})$, we obtain

$$\begin{aligned} & |x_i^{t+1} - x_i^*| \\ & = |(\alpha_1^t WA(x^t - x^*))_i + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i \\ & + \theta_1^t \partial\ell(\tilde{x}_i^{t+\frac{1}{2}}) - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}}| - (1 - \gamma^{t+1})\theta_2^t. \end{aligned}$$

For the case of $|\tilde{x}_i^{t+1}| \geq \bar{\theta}_2^t$, we can get

$$\begin{aligned} x_i^{t+1} & = K_i^{t+1} \tilde{x}_i^{t+1} = \tilde{x}_i^{t+1} \\ & = (\alpha_1^t WA(x^t - x^*))_i \\ & + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i + \theta_1^t \partial\ell(\tilde{x}_i^{t+\frac{1}{2}}) \\ & - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}} + x_i^*. \end{aligned}$$

Thus, we have

$$\begin{aligned} |x_i^{t+1} - x_i^*| & = |(\alpha_1^t WA(x^t - x^*))_i \\ & + ((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i \\ & + \theta_1^t \partial\ell(\tilde{x}_i^{t+\frac{1}{2}}) - B_i^{t+\frac{1}{2}} \tilde{x}_i^{t+\frac{1}{2}}|. \end{aligned}$$

the two cases, we can obtain

859

$$\begin{aligned} \|x^{t+1} - x^*\|_1 & = \sum_{i \in S+S^{(t+1)}} |x_i^{t+1} - x_i^*| \\ & \leq \sum_{i \in S+S^{(t+1)}} |(\alpha_1^t WA(x^t - x^*))_i| \\ & + \sum_{i \in S+S^{(t+1)}} |((I - \alpha_2^t WA)(x^{t+\frac{1}{2}} - x^*))_i| \\ & + (|\bar{S}^{(t+1)}| - |S^{(t+1)}| - |U_1^{(t+1)}| + |U_2^{(t+1)}|) \\ & (1 - \gamma^{t+1})\theta_2^t + (|S| + |S^{(t+1)}|) \\ & (1 + \gamma^{t+\frac{1}{2}})\theta_1^t + (|S| + |S^{(t+1)}|)\theta_2^t \\ & \leq \sum_{i \in S+S^{(t+1)}} \left(\sum_{j \in \text{supp}(\tilde{x}^t), j \neq i} \alpha_1^t |(WA)_{ij}(x_j^t - x_j^*)| \right. \\ & + \alpha_1^t |x_i^t - x_i^*| \\ & + \sum_{j \in \text{supp}(x^{t+\frac{1}{2}})} \alpha_2^t |(I - WA)_{ij}(x_j^{t+\frac{1}{2}} - x_j^*)| \\ & \left. + |1 - \alpha_2^t| |x_i^{t+\frac{1}{2}} - x_i^*| \right) \\ & + (|\bar{S}^{(t+1)}| - |S^{(t+1)}| - |U_1^{(t+1)}| + |U_2^{(t+1)}|) \\ & (1 - \gamma^{t+1})\theta_2^t + (|S| + |S^{(t+1)}|)(1 + \gamma^{t+\frac{1}{2}})\theta_1^t \\ & + (|S| + |S^{(t+1)}|)\theta_2^t \\ & \leq (|S| + |S^{(t+1)}| - 1)\alpha_1^t \mu(A) \|x^t - x^*\|_1 \\ & + \alpha_1^t \|x^t - x^*\|_1 + |1 - \alpha_2^t| \|x^{t+\frac{1}{2}} - x^*\|_1 \\ & + (|S| + |S^{(t+1)}|)\alpha_2^t \mu(A) \|x^{t+\frac{1}{2}} - x^*\|_1 \\ & + (|S| + |S^{(t+1)}|)(1 + \gamma^{t+\frac{1}{2}})\theta_1^t \\ & + (|\bar{S}^{(t+1)}| - |S^{(t+1)}| - |U_1^{(t+1)}| + |U_2^{(t+1)}|) \\ & (1 - \gamma^{t+1})\theta_2^t + (|S| + |S^{(t+1)}|)\theta_2^t, \end{aligned} \quad (31)$$

where the definitions of $U_1^{(t+1)}$ and $U_2^{(t+1)}$ are similar to (25). Then we set

$$\theta_2^t = \alpha_2^t \mu(A) \omega_{t+1}(k_{t+1}|\Theta) \sup_{x^*} \|x^{t+\frac{1}{2}} - x^*\|_1.$$

857 Summing up all $|x_i^{t+1} - x_i^*|$ which satisfy $i \in \text{supp}(\tilde{x}^{t+1}) \cup$
 858 $\text{supp}(x^*) = S + S^{(t+1)} = \bar{S}^{(t+1)} + S \setminus \bar{S}^{(t+1)} + S^{(t+1)}$ in

The definition of θ_1^t is the same as (26). Then by taking 860

supremum of both sides of (31), we obtain

$$\begin{aligned}
& \sup_{x^*} \|x^{t+1} - x^*\|_1 \\
& \leq \sup_{x^*} ((|S| + |S^{(t+1)}| - 1 + (1 + \gamma^{t+\frac{1}{2}}) \\
& \quad (|S| + |S^{(t+1)}|) \omega_{t+\frac{1}{2}}(k_{t+\frac{1}{2}}|\Theta)) \alpha_1^t \mu(A) + \alpha_1^t) \\
& \sup_{x^*} \|x^t - x^*\|_1 + \sup_{x^*} ((|S| + |S^{(t+1)}| \\
& \quad + (|S| - |\bar{S}^{(t+1)}|) + (|\bar{S}^{(t+1)}| - |S^{(t+1)}| \\
& \quad - |U_1^{(t+1)}| + |U_2^{(t+1)}|)(1 - \gamma^{t+1}) \omega_{t+1}(k_{t+1}|\Theta)) \\
& \quad \alpha_2^t \mu(A) + |1 - \alpha_2^t|) \sup_{x^*} \|x^{t+\frac{1}{2}} - x^*\|_1 \\
& \leq \sup_{x^*} (((|S| + |S^{(t+1)}| - 1 + (1 + \gamma^{t+\frac{1}{2}})(|S| + |S^{(t+1)}|) \\
& \quad \omega_{t+\frac{1}{2}}(k_{t+\frac{1}{2}}|\Theta)) \alpha_1^t \mu(A) + \alpha_1^t) + ((|S| + |S^{(t+1)}| \\
& \quad + (|S| - |\bar{S}^{(t+1)}|) + (|\bar{S}^{(t+1)}| - |S^{(t+1)}| - |U_1^{(t+1)}| \\
& \quad + |U_2^{(t+1)}|)(1 - \gamma^{t+1}) \omega_{t+1}(k_{t+1}|\Theta)) \alpha_2^t \mu(A) \\
& \quad + |1 - \alpha_2^t|)((|S| + |S^{(t+\frac{1}{2})}| + (|S| - |\bar{S}^{(t+\frac{1}{2})}| \\
& \quad + (|\bar{S}^{(t+\frac{1}{2})}| - |S^{(t+\frac{1}{2})}| - |U_1^{(t+\frac{1}{2})}| \\
& \quad + |U_2^{(t+\frac{1}{2})}|)(1 - \gamma^{t+\frac{1}{2}}) \omega_{t+\frac{1}{2}}(k_{t+\frac{1}{2}}|\Theta)) \alpha_1^t \mu(A) \\
& \quad + |1 - \alpha_1^t|)) \sup_{x^*} \|x^t - x^*\|_1 \\
& = \exp(c'_{t+1}) \sup_{x^*} \|x^t - x^*\|_1,
\end{aligned} \tag{32}$$

where

$$\begin{aligned}
c'_{t+1} &= \sup_{x^*} \log(((|S| + |S^{(t+1)}| - 1 + (1 + \gamma^{t+\frac{1}{2}}) \\
&\quad (|S| + |S^{(t+1)}|) \omega_{t+\frac{1}{2}}(k_{t+\frac{1}{2}}|\Theta)) \alpha_1^t \mu(A) + \alpha_1^t) \\
&\quad + ((|S| + |S^{(t+1)}| + (|S| - |\bar{S}^{(t+1)}| \\
&\quad + (|\bar{S}^{(t+1)}| - |S^{(t+1)}| - |U_1^{(t+1)}| + |U_2^{(t+1)}|) \\
&\quad (1 - \gamma^{t+1}) \omega_{t+1}(k_{t+1}|\Theta)) \alpha_2^t \mu(A) + |1 - \alpha_2^t|) \\
&\quad ((|S| + |S^{(t+\frac{1}{2})}| + (|S| - |\bar{S}^{(t+\frac{1}{2})}| \\
&\quad + (|\bar{S}^{(t+\frac{1}{2})}| - |S^{(t+\frac{1}{2})}| - |U_1^{(t+\frac{1}{2})}| + |U_2^{(t+\frac{1}{2})}|) \\
&\quad (1 - \gamma^{t+\frac{1}{2}}) \omega_{t+\frac{1}{2}}(k_{t+\frac{1}{2}}|\Theta)) \alpha_1^t \mu(A) + |1 - \alpha_1^t|)).
\end{aligned}$$

Here the second inequality in (32) holds due to Lemma 3.

□

Finally, according to Lemmas 3 and 4, we analyze and discuss the coefficient c'_{t+1} obtained above, and finally give the proof of Theorem 1.

Proof of Theorem 1

Proof. According to Lemma 4, we can obtain

$$\begin{aligned}
\|x^t - x^*\|_2 &\leq \|x^t - x^*\|_1 \leq \sup_{x^*} \|x^t - x^*\|_1 \\
&\leq \exp(c'_t) \sup_{x^*} \|x^{t-1} - x^*\|_1 \\
&\leq \exp\left(\sum_{i=1}^t c'_i\right) sB.
\end{aligned}$$

Then we discuss the value of c'_i . Let $t_0 = \lceil \log(\frac{sB}{\sigma})/c \rceil$. When $i \leq t_0$, we assume $\gamma^{i-\frac{1}{2}} = \gamma^i = 0$. Then we have

$$\begin{aligned}
c'_i &= \sup_{x^*} \log(((|S| + |S^{(i)}| - 1 + (|S| + |S^{(i)}|) \\
&\quad \omega_{i-\frac{1}{2}}(k_{i-\frac{1}{2}}|\Theta)) \alpha_1^{i-1} \mu(A) + \alpha_1^{i-1}) \\
&\quad + ((|S| + |S^{(i)}| + (|S| - |S^{(i)}| - |U_1^{(i)}| + |U_2^{(i)}|) \\
&\quad \omega_i(k_i|\Theta)) \alpha_2^{i-1} \mu(A) + |1 - \alpha_2^{i-1}|) \\
&\quad ((|S| + |S^{(i-\frac{1}{2})}| + (|S| - |S^{(i-\frac{1}{2})}| \\
&\quad - |U_1^{(i-\frac{1}{2})}| + |U_2^{(i-\frac{1}{2})}|) \omega_{i-\frac{1}{2}}(k_{i-\frac{1}{2}}|\Theta)) \alpha_1^{i-1} \mu(A) \\
&\quad + |1 - \alpha_1^{i-1}|)) \\
&= \log((s + k_i - 1 + (s + k_i) \omega_{i-\frac{1}{2}}(k_{i-\frac{1}{2}}|\Theta)) \\
&\quad \alpha_1^{i-1} \mu(A) + \alpha_1^{i-1}) + ((s + k_i + (s - k_i) \\
&\quad - \inf_{x^*} |U_1^{(i)}| + \sup_{x^*} |U_2^{(i)}|) \omega_i(k_i|\Theta)) \alpha_2^{i-1} \mu(A) \\
&\quad + |1 - \alpha_2^{i-1}|)((s + k_{i-\frac{1}{2}} + (s - k_{i-\frac{1}{2}}) \\
&\quad - \inf_{x^*} |U_1^{(i-\frac{1}{2})}| + \sup_{x^*} |U_2^{(i-\frac{1}{2})}|) \omega_{i-\frac{1}{2}}(k_{i-\frac{1}{2}}|\Theta)) \\
&\quad \alpha_1^{i-1} \mu(A) + |1 - \alpha_1^{i-1}|).
\end{aligned}$$

(i). When $\inf_{x^*} |U_1^{(i-\frac{1}{2})}| \geq \sup_{x^*} |U_2^{(i-\frac{1}{2})}|$ and $\inf_{x^*} |U_1^{(i)}| \geq \sup_{x^*} |U_2^{(i)}|$, we have $c'_i \leq c^*_i < 0$.

(ii). When $\inf_{x^*} |U_1^{(i-\frac{1}{2})}| < \sup_{x^*} |U_2^{(i-\frac{1}{2})}|$ and $\inf_{x^*} |U_1^{(i)}| < \sup_{x^*} |U_2^{(i)}|$, we have

$$\begin{aligned}
c'_i &= \log((s + k_i - 1 + (s + k_i) \omega_{i-\frac{1}{2}}(k_{i-\frac{1}{2}}|\Theta)) \\
&\quad \alpha_1^{i-1} \mu(A) + \alpha_1^{i-1}) + ((s + k_i + (s - k_i + e_i) \\
&\quad \omega_i(k_i|\Theta)) \alpha_2^{i-1} \mu(A) + |1 - \alpha_2^{i-1}|) \\
&\quad ((s + k_{i-\frac{1}{2}} + (s - k_{i-\frac{1}{2}} + e_{i-\frac{1}{2}}) \\
&\quad \omega_{i-\frac{1}{2}}(k_{i-\frac{1}{2}}|\Theta)) \alpha_1^{i-1} \mu(A) + |1 - \alpha_1^{i-1}|),
\end{aligned} \tag{33}$$

where $e_{i-\frac{1}{2}} = \sup_{x^*} |U_2^{(i-\frac{1}{2})}| - \inf_{x^*} |U_1^{(i-\frac{1}{2})}| > 0$, $e_i = \sup_{x^*} |U_2^{(i)}| - \inf_{x^*} |U_1^{(i)}| > 0$.

According to the definition of $\omega_t(\cdot)$ and

$$\begin{aligned}
\theta_1^{i-1} &= \alpha_1^{i-1} \omega_{i-\frac{1}{2}}(k_{i-\frac{1}{2}}|\Theta) \mu(A) \sup_{x^*} \|x^{i-1} - x^*\|_1, \\
\theta_2^{i-1} &= \alpha_2^{i-1} \omega_i(k_i|\Theta) \mu(A) \sup_{x^*} \|x^{i-\frac{1}{2}} - x^*\|_1,
\end{aligned}$$

we know that, for $x^{i-\frac{1}{2}}$, the number of “false positive” is not larger than $k_{i-\frac{1}{2}}$, i.e., $|S^{(i-\frac{1}{2})}| \leq k_{i-\frac{1}{2}}$ and for x^i , the number of “false positive” is not larger than k_i , i.e., $|S^{(i)}| \leq$

k_i . From Lemma 1, we can know that the number of element in $x^{i-\frac{1}{2}}$ which is “false positive” satisfies $k_{i-\frac{1}{2}} = 0$, when $\theta_1^{i-1} = \mu(A) \sup_{x^*} \|x^{i-1} - x^*\|_1$. Similarly, the number of element in x^i which is “false positive” satisfies $k_i = 0$, when $\theta_2^{i-1} = \mu(A) \sup_{x^*} \|x^{i-\frac{1}{2}} - x^*\|_1$. These results mean that $\omega(0|\Theta) \leq 1$ and $\omega(k|\Theta) \leq 1$, when $k > 0$. Thus, we assume that

$$\begin{aligned} & \exists k_0^{i-\frac{1}{2}}, \text{s.t. } 0 < k_0^{i-\frac{1}{2}} < s + e_{i-\frac{1}{2}}, \\ & \omega_{i-\frac{1}{2}}(k_0^{i-\frac{1}{2}}|\Theta) < 1 - \frac{1}{s - k_0^{i-\frac{1}{2}} + e_{i-\frac{1}{2}}}, \\ & \exists k_0^i, \text{s.t. } 0 < k_0^i < s + e_i, \omega_i(k_0^i|\Theta) < 1 - \frac{1}{s - k_0^i + e_i}. \end{aligned}$$

Then we can get $k_{i-\frac{1}{2}} = k_0^{i-\frac{1}{2}}$, $k_i = k_0^i$ by selecting θ_1^{i-1} and θ_2^{i-1} appropriately. Substituting the above results into (33), we obtain

$$\begin{aligned} c'_i &= \log(((s + k_i - 1 + (s + k_i)(1 - \frac{1}{s - k_{i-\frac{1}{2}} + e_{i-\frac{1}{2}}})) \\ &\quad \alpha_1^{i-1}\mu(A) + \alpha_1^{i-1}) + ((s + k_i + (s - k_i + e_i) \\ &\quad (1 - \frac{1}{s - k_i + e_i}))\alpha_2^{i-1}\mu(A) + |1 - \alpha_2^{i-1}|) \\ &\quad ((s + k_{i-\frac{1}{2}} + (s - k_{i-\frac{1}{2}} + e_{i-\frac{1}{2}}) \\ &\quad (1 - \frac{1}{s - k_{i-\frac{1}{2}} + e_{i-\frac{1}{2}}}))\alpha_1^{i-1}\mu(A) + |1 - \alpha_1^{i-1}|)) \\ &= \log(((2s + 2k_i - \frac{1}{s - k_{i-\frac{1}{2}} + e_{i-\frac{1}{2}}} - 1) \\ &\quad \alpha_1^{i-1}\mu(A) + \alpha_1^{i-1}) \\ &\quad + ((2s - 1 + e_i)\alpha_2^{i-1}\mu(A) + |1 - \alpha_2^{i-1}|) \\ &\quad ((2s - 1 + e_{i-\frac{1}{2}})\alpha_1^{i-1}\mu(A) + |1 - \alpha_1^{i-1}|)). \end{aligned}$$

We assume $\alpha_1^{i-1}, \alpha_2^{i-1} > 0$ and

$$2s + 2k_i - \frac{s + k_i}{s - k_{i-\frac{1}{2}} + e_{i-\frac{1}{2}}} \leq 1 - \frac{1}{\mu(A)}.$$

Then we get

$$\begin{aligned} c'_i &\leq \log(((2s - 1 + e_i)\alpha_2^{i-1}\mu(A) + |1 - \alpha_2^{i-1}|) \\ &\quad ((2s - 1 + e_{i-\frac{1}{2}})\alpha_1^{i-1}\mu(A) + |1 - \alpha_1^{i-1}|)). \end{aligned}$$

Moreover, we assume $s < \min\{\frac{1}{2}(1 + e_{i-\frac{1}{2}} + \frac{1}{\mu(A)}), \frac{1}{2}(1 + e_i + \frac{1}{\mu(A)})\}$, which implies $(2s - 1 + e_{i-\frac{1}{2}})\mu(A) < 1$ and $(2s - 1 + e_i)\mu(A) < 1$. For α_1^{i-1} , we also assume $\alpha_1^{i-1} < \frac{2}{1+(2s-1+e_{i-\frac{1}{2}})\mu(A)}$, i.e., $\alpha_1^{i-1} \in (0, \frac{2}{1+(2s-1+e_{i-\frac{1}{2}})\mu(A)})$. Thus, when $0 < \alpha_1^{i-1} \leq 1$, we have

$$\begin{aligned} & (2s - 1 + e_{i-\frac{1}{2}})\alpha_1^{i-1}\mu(A) + |1 - \alpha_1^{i-1}| \\ &= (2s - 1 + e_{i-\frac{1}{2}})\alpha_1^{i-1}\mu(A) + 1 - \alpha_1^{i-1} < 1. \end{aligned}$$

If $1 < \alpha_1^{i-1} < \frac{2}{1+(2s-1+e_{i-\frac{1}{2}})\mu(A)}$, we have

$$\begin{aligned} & (2s - 1 + e_{i-\frac{1}{2}})\alpha_1^{i-1}\mu(A) + |1 - \alpha_1^{i-1}| \\ &= (2s - 1 + e_{i-\frac{1}{2}})\alpha_1^{i-1}\mu(A) + \alpha_1^{i-1} - 1 < 1. \end{aligned}$$

For α_2^{i-1} , we also assume $\alpha_2^{i-1} < \frac{2}{1+(2s-1+e_i)\mu(A)}$, i.e., $\alpha_2^{i-1} \in (0, \frac{2}{1+(2s-1+e_i)\mu(A)})$. Thus, when $0 < \alpha_2^{i-1} \leq 1$, we have

$$\begin{aligned} & (2s - 1 + e_i)\alpha_2^{i-1}\mu(A) + |1 - \alpha_2^{i-1}| \\ &= (2s - 1 + e_i)\alpha_2^{i-1}\mu(A) + 1 - \alpha_2^{i-1} < 1. \end{aligned}$$

If $1 < \alpha_2^{i-1} < \frac{2}{1+(2s-1+e_{i-\frac{1}{2}})\mu(A)}$, we have

$$\begin{aligned} & (2s - 1 + e_i)\alpha_2^{i-1}\mu(A) + |1 - \alpha_2^{i-1}| \\ &= (2s - 1 + e_i)\alpha_2^{i-1}\mu(A) + \alpha_2^{i-1} - 1 < 1. \end{aligned}$$

Summing up all the results above, we can obtain $c'_i < 0$. 874

(iii). By a similar derivation to (ii), we can get that $c'_i < 0$ when $\inf_{x^*} |U_1^{(i-\frac{1}{2})}| \geq \sup_{x^*} |U_2^{(i-\frac{1}{2})}|$ and $\inf_{x^*} |U_1^{(i)}| < \sup_{x^*} |U_2^{(i)}|$ as well as $\inf_{x^*} |U_1^{(i-\frac{1}{2})}| < \sup_{x^*} |U_2^{(i-\frac{1}{2})}|$ and $\inf_{x^*} |U_1^{(i)}| \geq \sup_{x^*} |U_2^{(i)}|$. 875
876
877
878

When $i > t_0$, $\sup_{x^*} \|x^i - x^*\|_1 < sB \exp(ci) \leq \sigma \leq \min_{i \in S} |x_i^*|$, which implies $\bar{S}^{(i)} = S$. Thus, we have $\bar{S}^{(i-\frac{1}{2})} = \bar{S}^{(i)} = S$, according to $\text{supp}(x^{i-\frac{1}{2}}) \subset \text{supp}(x^i)$. Then choose θ_1^{i-1} and θ_2^{i-1} so that $k_{i-\frac{1}{2}} = 0$, $\omega_{i-\frac{1}{2}}(k_{i-\frac{1}{2}}|\Theta) \leq 1$ and $k_i = 0$, $\omega_i(k_i|\Theta) \leq 1$. Thus, we obtain $|S^{(i-\frac{1}{2})}| = k_{i-\frac{1}{2}} = 0$ and $|S^{(i)}| = k_i = 0$, which mean $|U_2^{(i-\frac{1}{2})}| = 0$ and $|U_2^{(i)}| = 0$. So, we have

$$\begin{aligned} c'_i &= \sup_{x^*} \log(((|S| - 1 + (1 + \gamma^{i-\frac{1}{2}})|S|)\omega_{i-\frac{1}{2}}(k_{i-\frac{1}{2}}|\Theta)) \\ &\quad \alpha_1^{i-1}\mu(A) + \alpha_1^{i-1}) + ((|S| + (|S| - |U_1^{(i)}|))(1 - \gamma^i) \\ &\quad \omega_i(k_i|\Theta))\alpha_2^{i-1}\mu(A) + |1 - \alpha_2^{i-1}|)((|S| + (|S| \\ &\quad - |U_1^{(i-\frac{1}{2})}|)(1 - \gamma^{i-\frac{1}{2}})\omega_{i-\frac{1}{2}}(k_{i-\frac{1}{2}}|\Theta)) \\ &\quad \alpha_1^{i-1}\mu(A) + |1 - \alpha_1^{i-1}|)) \\ &\leq \log(((s - 1 + (1 + \gamma^{i-\frac{1}{2}})s\omega_{i-\frac{1}{2}}(k_{i-\frac{1}{2}}|\Theta))) \\ &\quad \alpha_1^{i-1}\mu(A) + \alpha_1^{i-1}) + ((s + (s - \inf_{x^*} |U_1^{(i)}|) \\ &\quad (1 - \gamma^i))\alpha_2^{i-1}\mu(A) + |1 - \alpha_2^{i-1}|) \\ &\quad ((s + (s - \inf_{x^*} |U_1^{(i-\frac{1}{2})}|)(1 - \gamma^{i-\frac{1}{2}})) \\ &\quad \alpha_1^{i-1}\mu(A) + |1 - \alpha_1^{i-1}|)). \end{aligned}$$

Because $1 - \omega_{i-\frac{1}{2}}(s|\Theta) < \gamma^{i-\frac{1}{2}} \leq 1$ and $1 - \omega_i(s|\Theta) <$ 879

880 $\gamma^i \leq 1$, we get

$$\begin{aligned}
c'_i &< \log(((s-1+2s\omega_{i-\frac{1}{2}}(k_{i-\frac{1}{2}}|\Theta))\alpha_1^{i-1}\mu(A)+\alpha_1^{i-1}) \\
&\quad + ((s+(s-\inf_{x^*}|U_1^{(i)}|)\omega_i(s|\Theta))\alpha_2^{i-1}\mu(A) \\
&\quad + |1-\alpha_2^{i-1}|t((s+(s-\inf_{x^*}|U_1^{(i-\frac{1}{2})}|) \\
&\quad \omega_{i-\frac{1}{2}}(s|\Theta))\alpha_1^{i-1}\mu(A)+|1-\alpha_1^{i-1}|)) \\
&< \log(((s+k_i-1+2(s+k_i)\omega_{i-\frac{1}{2}}(k_{i-\frac{1}{2}}|\Theta)) \\
&\quad \alpha_1^{i-1}\mu(A)+\alpha_1^{i-1})+((s+k_i+(s-k_i \\
&\quad -\inf_{x^*}|U_1^{(i)}|)\omega_i(k_i|\Theta))\alpha_2^{i-1}\mu(A)+|1-\alpha_2^{i-1}|) \\
&\quad ((s+k_{i-\frac{1}{2}}+(s-k_{i-\frac{1}{2}}-\inf_{x^*}|U_1^{(i-\frac{1}{2})}|) \\
&\quad \omega_{i-\frac{1}{2}}(k_{i-\frac{1}{2}}|\Theta))\alpha_1^{i-1}\mu(A)+|1-\alpha_1^{i-1}|)). \tag{34}
\end{aligned}$$

Similar to the proof in (ii), we assume that

$$\exists k_0^{i-\frac{1}{2}}, s.t. 0 < k_0^{i-\frac{1}{2}} < s, \omega_{i-\frac{1}{2}}(k_0^{i-\frac{1}{2}}|\Theta) < 1 - \frac{1}{s - k_0^{i-\frac{1}{2}}},$$

$$\exists k_0^i, s.t. 0 < k_0^i < s, \omega_i(k_0^i|\Theta) < 1 - \frac{1}{s - k_0^i}.$$

Then we can get $k_{i-\frac{1}{2}} = k_0^{i-\frac{1}{2}}$, $k_i = k_0^i$ by selecting θ_1^{i-1} and θ_2^{i-1} appropriately. Substituting the above results into (34), we obtain

$$\begin{aligned}
c'_i &< \log(((s+k_i-1+2(s+k_i)(1-\frac{1}{s-k_{i-\frac{1}{2}}})) \\
&\quad \alpha_1^{i-1}\mu(A)+\alpha_1^{i-1})+((s+k_i+(s-k_i \\
&\quad -\inf_{x^*}|U_1^{(i)}|)(1-\frac{1}{s-k_i})\alpha_2^{i-1}\mu(A)+|1-\alpha_2^{i-1}|) \\
&\quad ((s+k_{i-\frac{1}{2}}+(s-k_{i-\frac{1}{2}}-\inf_{x^*}|U_1^{(i-\frac{1}{2})}|) \\
&\quad (1-\frac{1}{s-k_{i-\frac{1}{2}}}))\alpha_1^{i-1}\mu(A)+|1-\alpha_1^{i-1}|)).
\end{aligned}$$

Besides, we have $\alpha_1^{i-1} > 0$. We assume

$$3s+3k_{t+1}-\frac{2s+2k_i}{s-k_{i-\frac{1}{2}}} \leq 1 - \frac{1}{\mu(A)}.$$

Then we have

$$\begin{aligned}
c'_i &< \log(((s+k_i+(s-k_i-\inf_{x^*}|U_1^{(i)}|)(1-\frac{1}{s-k_i})) \\
&\quad \alpha_2^{i-1}\mu(A)+|1-\alpha_2^{i-1}|)((s+k_{i-\frac{1}{2}}+(s-k_{i-\frac{1}{2}} \\
&\quad -\inf_{x^*}|U_1^{(i-\frac{1}{2})}|)(1-\frac{1}{s-k_{i-\frac{1}{2}}}))\alpha_1^{i-1}\mu(A) \\
&\quad +|1-\alpha_1^{i-1}|)) \\
&\leq \log(((s+k_i+(s-k_i)(1-\frac{1}{s-k_i}))\alpha_2^{i-1}\mu(A) \\
&\quad +|1-\alpha_2^{i-1}|)((s+k_{i-\frac{1}{2}}+(s-k_{i-\frac{1}{2}}) \\
&\quad (1-\frac{1}{s-k_{i-\frac{1}{2}}}))\alpha_1^{i-1}\mu(A)+|1-\alpha_1^{i-1}|)) \\
&= \log(((2s-1)\alpha_2^{i-1}\mu(A)+|1-\alpha_2^{i-1}|) \\
&\quad ((2s-1)\alpha_1^{i-1}\mu(A)+|1-\alpha_1^{i-1}|)).
\end{aligned}$$

In addition, we assume $s < \frac{1}{2}(1+\frac{1}{\mu(A)})$, from which we can obtain $(2s-1)\mu(A) < 1$. Besides, we have $\alpha_1^{i-1}, \alpha_2^{i-1} \in (0, \frac{2}{1+(2s-1)\mu(A)})$. If $0 < \alpha_1^{i-1}, \alpha_2^{i-1} \leq 1$, we have $(2s-1)\alpha_1^{i-1}\mu(A)+|1-\alpha_1^{i-1}| < 1$ and $(2s-1)\alpha_2^{i-1}\mu(A)+|1-\alpha_2^{i-1}| < 1$. Thus, we obtain $c'_i < 0$. If $1 < \alpha_1^{i-1}, \alpha_2^{i-1} < \frac{2}{1+(2s-1)\mu(A)}$, we have

$$(2s-1)\alpha_1^{i-1}\mu(A)+|1-\alpha_1^{i-1}| = (2s-1)\alpha_1^{i-1}\mu(A)+\alpha_1^{i-1}-1 < 1,$$

$$(2s-1)\alpha_2^{i-1}\mu(A)+|1-\alpha_2^{i-1}| = (2s-1)\alpha_2^{i-1}\mu(A)+\alpha_2^{i-1}-1 < 1,$$

which imply that $c'_i < 0$.

Combining all the above theoretical derivation, we get $c'_i < 0, i = 1, 2, \dots, t$. We define $c' = \max\{c'_i, i = 1, 2, \dots, t\}$. Then

$$\|x^t - x^*\|_2 \leq sB \exp\left(\sum_{i=1}^t c'_i\right) < sB \exp(c't). \tag{35}$$

Note that $\forall i, c'_i < 0$. Thus, we have $c' < 0$, which implies the proposed algorithm achieves linear convergence. □

S3. Comparison between MT($\cdot, \theta, \bar{\theta}$) and HELU $_{\sigma}(\cdot)$

• Difference of Motivation

The motivation of HELU $_{\sigma}(\cdot)$ is to propose a deep ℓ_0 -norm network, but the original HT is a non-continuous and non-differentiable function, which will cause the network to be difficult to train and bad performance. Therefore, in order to address this issue, Wang, Ling, and Huang (2016) **improved HT** and proposed HELU $_{\sigma}(\cdot)$. As for our MT, we find that there are some problems in the gain gate proposed in (Wu et al. 2020), so **in order to further improve the ability of ST and make the algorithm get a more sparse solution**, we propose MT. Therefore, their motivations are totally different.

• Differences of Some Specific Details

The definition of $\text{HELU}_\sigma(\cdot)$ is formulated as follows:

$$[\text{HELU}_\sigma(\mathbf{u})]_i = \begin{cases} 0, & \text{if } |\mathbf{u}_i| \leq 1 - \sigma, \\ \frac{(\mathbf{u}_i - 1 + \sigma)}{\sigma}, & \text{if } 1 - \sigma < \mathbf{u}_i < 1, \\ \frac{(\mathbf{u}_i + 1 - \sigma)}{\sigma}, & \text{if } -1 < \mathbf{u}_i < \sigma - 1, \\ \mathbf{u}_i, & \text{if } |\mathbf{u}_i| \geq 1, \end{cases}$$

where σ is a hyper parameter. Therefore, we know that $\text{HELU}_\sigma(\cdot)$ is a function with a hyper parameter, while our MT has no hyper parameter, and all thresholds in MT are obtained through the learning of the network. In addition, for $\text{HELU}_\sigma(\cdot)$, Wang, Ling, and Huang (2016) artificially make the hyper parameter σ tend to 0, i.e., make the first threshold gradually move to the second one and cause $K^t \rightarrow \infty$, but they did not provide the corresponding theoretical analysis. However, through theoretical analysis, we can obtain that the first threshold in MT satisfies $\theta^t \rightarrow 0$ as $t \rightarrow \infty$, i.e., $K^t \rightarrow 1$, which can also be verified by our experimental results in Figure 21. Thus, through comparison, we know that the adjustment direction of the hyper parameter σ in $\text{HELU}_\sigma(\cdot)$ is opposite to the reality, and thus $\text{HELU}_\sigma(\cdot)$ will somehow reduce the performance of the network, which is also mentioned in Section 7.3.2 in (Xin et al. 2016a). As for our MT, according to the above experimental results (Section 4.1), we know that MT can effectively improve the performance of the algorithms.

• Theoretical Analysis and Experimental Verification

Proposition 1. For

$$\begin{cases} x^{t+\frac{1}{2}} = \text{MT}(x^t - \alpha_1^t W(Ax^t - y), \theta_1^t, \bar{\theta}_1^t), \\ x^{t+1} = \text{MT}(x^t - \alpha_2^t W(Ax^{t+\frac{1}{2}} - y), \theta_2^t, \bar{\theta}_2^t), \end{cases}$$

if the following holds uniformly for any $x^* \in \mathcal{X}(B, s)$:

$$x^t, x^{t+\frac{1}{2}} \rightarrow x^*, \text{ as } t \rightarrow \infty,$$

and $W \in \mathcal{W}(A)$, $\forall t = 0, 1, 2, \dots$, then θ_1^t and θ_2^t satisfy $\theta_1^t, \theta_2^t \rightarrow 0$ as $t \rightarrow \infty$.

Proof. According the proof of Lemma 3, when $0 \leq |\tilde{x}_i^{t+\frac{1}{2}}| < \bar{\theta}_1^t$, from (20), we have

$$\begin{aligned} x_i^{t+\frac{1}{2}} &= \text{MT}(x_i^t - (W(Ax^t - y))_i, \theta_1^t, \bar{\theta}_1^t) \\ &= ((I - WA)x^t)_i + (WAx^*)_i \\ &\quad - \theta_1^t \partial \ell_1\left(\frac{1}{K_i^{t+\frac{1}{2}}} x_i^{t+\frac{1}{2}}\right) + \frac{B_i^{t+\frac{1}{2}}}{K_i^{t+\frac{1}{2}}} x_i^{t+\frac{1}{2}} \\ &= ((I - WA)x^t)_i + (WAx^*)_i \\ &\quad - \theta_1^t \partial \ell_1(x_i^{t+\frac{1}{2}}) + \frac{\theta_1^t}{\bar{\theta}_1^t} x_i^{t+\frac{1}{2}}. \end{aligned}$$

Then we consider the situation of $t \rightarrow \infty$. When $t \rightarrow \infty$, we have

$$x_i^* = ((I - WA)x^*)_i + (WAx^*)_i - \theta_1^t \partial \ell_1(x_i^*) + \frac{\theta_1^t}{\bar{\theta}_1^t} x_i^*,$$

i.e.,

$$\theta_1^t (\partial \ell_1(x_i^*) - \frac{1}{\bar{\theta}_1^t} x_i^*) = 0. \quad (36)$$

For $\forall i \in S$, we let $x_i^* \rightarrow 0$ but $x_i^* \neq 0$, thus $\partial \ell_1(x_i^*) = \text{sign}(x_i^*)$. Since $x_i^* \rightarrow 0$ and $\partial \ell_1(x_i^*) = \text{sign}(x_i^*) \neq 0$, in order to make the left side of (36) tend to zero, we obtain

$$\theta_1^t \rightarrow 0, \text{ as } t \rightarrow \infty.$$

Similarly, we can also obtain $\theta_2^t \rightarrow 0$, as $t \rightarrow \infty$. Finally, $\theta_1^t, \theta_2^t \rightarrow 0$, as $t \rightarrow \infty$ is proved \square

Besides, we also provide the experimental verification of Proposition 1. The experimental results are shown in Figure 21.

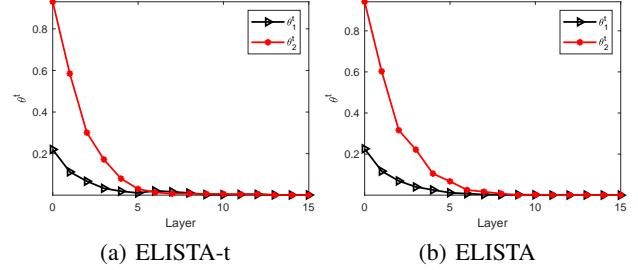


Figure 21: Justification of ELISTA-t and ELISTA (noiseless case): θ_1^t and θ_2^t satisfy $\theta_1^t, \theta_2^t \rightarrow 0$ as $t \rightarrow \infty$.

929