

Stable and Interpretable Unrolled Dictionary Learning

Anonymous authors

Paper under double-blind review

Abstract

The dictionary learning problem, representing data as a combination of a few atoms, has long stood as a popular method for learning representations in statistics and signal processing. The most popular dictionary learning algorithm alternates between sparse coding and dictionary update steps, and a rich literature has studied its theoretical convergence. The success of dictionary learning relies on access to a “good” initial estimate of the dictionary and the ability of the sparse coding step to provide an unbiased estimate of the code. The growing popularity of unrolled sparse coding networks has led to the empirical finding that backpropagation through such networks performs dictionary learning. We offer the theoretical analysis of these empirical results through PUDLE, a Provable Unrolled Dictionary LEarning method. We provide conditions on the network initialization and data distribution sufficient to recover and preserve the support of the latent code. Additionally, we address two challenges; first, the vanilla unrolled sparse coding computes a biased code estimate, and second, gradients during backpropagated learning can become unstable. We show approaches to reduce the bias of the code estimate in the forward pass, and that of the dictionary estimate in the backward pass. We propose strategies to resolve the learning instability by tuning network parameters and modifying the loss function. Overall, we highlight the impact of loss, unrolling, and backpropagation on convergence. We complement our findings through synthetic and image denoising experiments. Finally, we demonstrate PUDLE’s interpretability, a driving factor in designing deep networks based on iterative optimizations, by building a mathematical relation between network weights, its output, and the training set.

1 Introduction

This paper considers the dictionary learning problem, namely representing data $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^m$ as linear combinations of a few atoms from a dictionary $\mathbf{D} \in \mathcal{D} \subset \mathbb{R}^{m \times p}$. Given \mathbf{x} and \mathbf{D} , the problem of recovering the sparse (few non-zero elements) coefficients $\mathbf{z} \in \mathbb{R}^p$ is referred to as sparse coding, and can be solved through the lasso (Tibshirani, 1996) (also known as basis pursuit (Chen et al., 2001)):

$$\ell_{\mathbf{x}}(\mathbf{D}) := \min_{\mathbf{z} \in \mathbb{R}^p} \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D}) + h(\mathbf{z}) \quad (1)$$

where $\mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D}) = \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|_2^2$, and $h(\mathbf{z}) = \lambda \|\mathbf{z}\|_1$. Specifically, the problem aims to recover a dictionary \mathbf{D}^* that generates the data, i.e.,

$$\mathbf{x} = \mathbf{D}^* \mathbf{z}^* \quad (2)$$

where \mathbf{z}^* is sparse. Olshausen and Field (Olshausen & Field, 1997) introduced (2) in computational neuroscience as a model for how early layers of the visual cortex process natural images. Sparse coding has been widely studied and utilized in the statistics (Hastie et al., 2015) and signal processing communities (Elad, 2010). A few practical examples are denoising (Elad & Aharon, 2006), super-resolution (Yang et al., 2010), text processing (Jenatton et al., 2011), and classification (Mairal et al., 2009b), where it enables the extraction of sparse high-dimensional features representing data. Moreover, sparse modelling is ubiquitous in many other fields such as seismic signal processing (Nose-Filho et al., 2018), radar sensing for target detections (Bajwa et al., 2011), and astrophysics for image reconstruction from interferometric data (Akiyama et al., 2017). Furthermore, Cleary et al. (2017; 2021) use this model to learn a dictionary consisting of gene modules for efficient imaging transcriptomics.

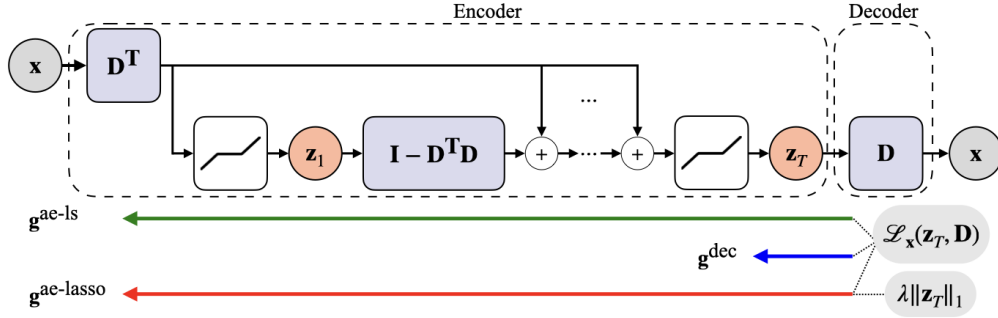


Figure 1: Provable unrolled dictionary learning (PUDLE): Unrolled network architecture with dictionary D .

Sparse coding has been utilized to construct neural architectures through approaches such as sparse energy-based models (Ranzato et al., 2007; 2008) or recurrent sparsifying encoders (Gregor & LeCun, 2010). The latter has initiated a growing literature on constructing interpretable deep networks based on an approach referred to as algorithm unrolling (Hershey et al., 2014; Monga et al., 2019). Deep unrolled neural networks have gained popularity as inference maps in recent years due to their computational efficiency and their performance in various domains such as image denoising (Simon & Elad, 2019; Tolooshams et al., 2021a; 2020), super-resolution (Wang et al., 2015), medical imaging (Solomon et al., 2020), deblurring (Schuler et al., 2016; Li et al., 2020), radar sensing (Tolooshams et al., 2021b), and speech processing (Hershey et al., 2014).

Prior to the advent of unrolled networks, gradient-based dictionary learning relied on analytic gradients computed from the lasso given the sparse code. With unrolled networks, automatic differentiation (Baydin et al., 2018), referred to as backpropagation (LeCun et al., 2012) in the reverse-mode, gained attention for parameter estimation (Tolooshams et al., 2018). The automatic gradient is obtained by backpropagation through the algorithm used to estimate the code. Automatic differentiation in reverse and forward-mode (Franceschi et al., 2017) is used in other areas, e.g., hyperparameter selection (Feurer & Hutter, 2019), and in a more relevant context, in the seminal work of LISTA (Gregor & LeCun, 2010). Other works demonstrated empirically the convergence of ℓ_1 -based dictionary learning by backpropagation through unrolled networks (Tolooshams et al., 2021a). Given finite computational power, Tolooshams et al. (2021a) convert sparse coding into an encoder by unrolling T iterations of ISTA (Daubechies et al., 2004; Blumensath & Davies, 2008), and attach to it a linear decoder for reconstructing. Unrolled networks obtained in this manner suffer from two important limitations.

First, the sparse coding step in the forward pass computes a biased estimate of the code. This results, in turn, in a biased estimate of the backward gradient and, hence, a degradation of dictionary recovery performance. Second, as studied recently (Malézieux et al., 2022), inaccuracies in the early iterations of the unrolled network make backpropagation unstable. We address both of these shortcomings in this paper. Moreover, while Malézieux et al. (2022) analyze the gradient computed by backpropagation through unrolled sparse-coding networks, there is no known theoretical analysis of how weight updates using this gradient impact the recovery of a ground-truth code z^* , nor of their convergence to a ground-truth dictionary D^* .

This paper proposes a Provable Unrolled Dictionary LEarning (PUDLE) (Figure 1). We aim to recover D^* by training the network using backpropagation with a learning rate of η . Three different choices affect the gradient: the number of unrolled iterations, the loss, and whether one backpropagates through the decoder only or through both the encoder and decoder. We highlight the impact of such choices on the convergence of the training algorithm. Backpropagation through the decoder results in the analytic gradient g_t^{dec} using the code estimate z_t . The gradients $g_t^{\text{ae-lasso}}$ and $g_t^{\text{ae-ls}}$ are computed by backpropagation through the autoencoder using the lasso and least-squares objectives, respectively (Algorithm 2). We compare the gradients with the classical gradient-based alternating-minimization algorithm for dictionary learning (Chatterji & Bartlett, 2017) (i.e., cycling between sparse coding and dictionary update steps using the analytic gradient \hat{g} (Algorithm 1)), and provide a theoretical analysis of gradient-based recovery of the dictionary D^* . We provide sufficient conditions under which the gradient computation, hence the learning, is stable. We show how using the reconstruction loss with backpropagation ameliorates the propagation of the forward pass bias into the gradient estimate from the backward pass. Finally, we demonstrate the interpretability of the unrolled network. Our contributions are:

- Unrolled sparse coding** Unlike prior work (Malézieux et al., 2022) that studies the ability of sparse coding to recover the *solution of the lasso* (1) given the current estimate of the dictionary (we call this local estimation), we study unrolled sparse coding for recovery of the *true generating code in* (2) (we call this global estimation). We provide sufficient conditions on the network and data distributions such that the forward pass recovers (Theorem 4.1) and preserves (Theorem 4.2) the correct code support. Assuming support identification, we show the linear convergence of the code estimated through the unrolled iterations to the solution of the lasso (Theorem 4.3). Moreover, in a more general scenario, we show that the error in the code estimate is upper bounded by two terms, i.e., one associated with the dictionary error and the other to the bias of the estimate of code amplitude, due to ℓ_1 -based optimization (Theorem 4.4). The latter highlights that vanilla lasso (ℓ_1 -based) sparse coding computes a biased estimate of codes, and below we discuss strategies to either alleviate this bias in the forward pass or mitigate its propagation into the backward pass for dictionary learning.
- Mitigation of coding bias propagation into dictionary learning** We study gradient estimation for dictionary learning in PUDLE. We decompose the upper bound on the gradient errors compared to the gradient direction to recover \mathbf{D}^* into terms involving the current dictionary error, the bias of the code estimate, and the lasso loss used to compute the gradient. We show that using only the reconstruction loss while backpropagating (i.e., $\mathbf{g}_t^{\text{ae-ls}}$) results in the vanishing of the upper bound due to the usage of lasso loss. This means that given fixed λ , $\mathbf{g}_t^{\text{ae-ls}}$ ameliorates the propagation of the forward pass bias into the backward pass. Specifically, we show that $\mathbf{g}_t^{\text{ae-ls}}$ is a better estimator of the direction to recover \mathbf{D}^* than $\mathbf{g}_t^{\text{dec}}$ and $\mathbf{g}_t^{\text{ae-lasso}}$. Hence, weight updates using $\mathbf{g}_t^{\text{ae-ls}}$ converges to a closer neighbourhood of \mathbf{D}^* (Theorem 4.8). In a supervised image denoising task, we show that the advantage of $\mathbf{g}_t^{\text{ae-ls}}$ goes beyond dictionary learning; $\mathbf{g}_t^{\text{ae-ls}}$ results in better image denoising compared to $\mathbf{g}_t^{\text{dec}}$. Furthermore, our network outperforms the sparse coding scheme in NOODL, a state-of-the-art online dictionary learning algorithm (Rambhatla et al., 2018) (Table 1). Moreover, we show that the bias in the estimate of \mathbf{D}^* vanishes as $\lambda_t = \lambda \nu^t$ (with $0 < \nu < 1$) decays within the forward unrolled layers (Figure 10). This strategy results in an unbiased estimate of the code \mathbf{z}^* and, hence, of \mathbf{D}^* .
- Stability of unrolled learning** We show that under proper dictionary initialization, the instability of the gradient $\mathbf{g}_t^{\text{ae-lasso}}$ computation, studied by Malézieux et al. (2022), as T increases is resolved. We give a condition under which the code support is identified and recovered after one iteration and, hence, gradient computation stays stable. Second, in the absence of support identification in early iterations, we propose to use the gradient $\mathbf{g}_t^{\text{ae-ls}}$ which resolves the stability issue introduced by lasso loss in the backward pass. We highlight this stability through image denoising training without gradient explosion.
- Interpretable sparse codes and dictionary** Prior work has discussed algorithm unrolling for designing interpretable deep architectures based on optimization models (Monga et al., 2019), or interpretability of sparse representations in dictionary learning models (Kim et al., 2010). However, there is no known work to mathematically characterize the interpretability of unrolled network architectures. In this regard, first, we construct a mathematical relation between learned weights (dictionary) at gradient convergence and the training data (Theorem 5.1). Second, we relate the inferred representation/reconstruction of test examples to the training data. We highlight several interpretable features of the unrolled dictionary learning network. Specifically, we perform analysis that provide insights into questions such as *why am I learning a particular feature in the dictionary?* or *from what part of the training set or an image I am learning that feature?* (Figure 7). Moreover, we provide an explanation of the relation between the new test image denoised/reconstructed through the network and the training dataset. The model provides insights on *how training images are used to reconstruct a new test image* (Figure 8) or *how the test image picks up training images that have a similar representation to itself to reconstruct* (Figure 9).

2 Related Works

There is vast literature on the theoretical convergence of dictionary learning. Spielman et al. (2012) proposed a factorization method to recover the dictionary in the undercomplete setting (i.e., $p \leq m$). Barak et al. (2015) proposed to solve dictionary learning via sum-of-squares semidefinite program. K-SVD (Aharon et al., 2006) and MOD (Engan et al., 1999) are popular greedy approaches. Alternating-minimization-based methods have been used extensively in theory and practice (Jain et al., 2013; Agarwal et al., 2014; Arora et al., 2014).

Recent work has incorporated gradient-based updates into alternating minimization (Chatterji & Bartlett, 2017; Arora et al., 2015; Rambhatla et al., 2018). Chatterji & Bartlett (2017) provided a finite sample analysis and convergence guarantees when updating the dictionary using the analytic gradient. Arora et al. (2015) proposed neurally plausible sparse coding approaches with analytic gradients. Another work focused on online dictionary learning (Mairal et al., 2009a) with an unbiased gradient updates (Rambhatla et al., 2018). Arora et al. (2015) discussed methods to reduce the bias of dictionary estimate, and Rambhatla et al. (2018) showed how to reduce bias in code and dictionary estimates. A common feature in the above-mentioned work is the use of analytic gradients, i.e., explicitly designing gradient updates independent of the sparse coding step and not utilizing automatic gradients with deep learning optimizers. A theoretical analysis of backpropagation for dictionary learning exists only for shallow autoencoders (Rangamani et al., 2018; Nguyen et al., 2019).

The theoretical analysis of unrolled neural networks has mainly analyzed the convergence speed of variants of LISTA (Gregor & LeCun, 2010), where the focus is on sparse coding (i.e., the encoder) not dictionary learning (Sprechmann et al., 2012; Xin et al., 2016; Moreau & Bruna, 2017; Giryes et al., 2018; Chen et al., 2018; Liu & Chen, 2019; Ablin et al., 2019). Moreau & Bruna (2017) showed that upon successful factorization of the Gram matrix of the dictionary within layers, the network achieves accelerated convergence. Giryes et al. (2018) examined the tradeoffs between reconstruction accuracy and convergence speed of LISTA. Moreover, Chen et al. (2018) studied the learning dynamics of the weights and biases of unrolled-ISTA and proved that it achieves linear convergence. Follow-up works investigated the dynamics of step size in a recursive sparse coding encoder (Liu & Chen, 2019; Ablin et al., 2019). Ablin et al. (2019) minimized the lasso through backpropagation but still assumed the knowledge of the dictionary at the decoder.

Ablin et al. (2020) compared analytic and automatic gradient estimators of min-min optimizations with smooth and differentiable functions. Moreover, Malézieux et al. (2022) studied the stability of gradient approximation in the early regime of unrolling for dictionary learning. Unlike our work, where we evaluate the gradients for model recovery, Ablin et al. (2020) and Malézieux et al. (2022) studied the asymptotic gradient errors locally in each step of an alternating minimization and did not provide errors concerning \mathbf{z}^* or \mathbf{D}^* .

3 Preliminaries

Given n independent samples, dictionary learning aims to minimize the empirical risk, i.e.,

$$\min_{\mathbf{D} \in \mathcal{D}} \mathcal{R}_n(\mathbf{D}) \quad \text{with} \quad \mathcal{R}_n(\mathbf{D}) \triangleq \frac{1}{n} \sum_{i=1}^n \ell_{\mathbf{x}^i}(\mathbf{D}) \quad (3)$$

where $\lim_{n \rightarrow \infty} \mathcal{R}_n(\mathbf{D}) = \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [\ell_{\mathbf{x}}(\mathbf{D})]$ a.s. To prevent scaling ambiguity between the code \mathbf{z} and dictionary \mathbf{D} , it is common to constrain the norm of the dictionary columns. Hence, we define the set of feasible solutions for the dictionary as $\mathcal{D} \triangleq \{\mathbf{D} \in \mathbb{R}^{m \times p} \text{ s.t. } \forall j \in \{1, 2, \dots, p\}, \|\mathbf{D}_j\|_2^2 \leq 1\}$. We can project estimates of \mathbf{D} onto the feasible set by performing $\mathbf{D}_j \leftarrow 1/\max(\|\mathbf{D}_j\|_2, 1) \mathbf{D}_j$, either at every update or at the end of training. We assume certain properties on the data, specifically its domain (Assumption 3.1), energy (Assumption 3.2), code distribution (Assumption 3.3), and generating dictionary (Assumption 3.4).

Assumption 3.1 (Domain signals). \mathcal{X} and \mathcal{D} are both compact convex sets.

Assumption 3.2 (Bounded signals). $\exists M > 0$ s.t. $\|\mathbf{x}\|_2 < M \forall \mathbf{x} \in \mathcal{X}$.

Assumption 3.3 (Code distribution). The code \mathbf{z}^* is at most s -sparse with the support $S^* = \text{supp}(\mathbf{z}^*)$. Each element in S^* is chosen from the set $[1, p]$, uniformly at random without replacement. $p_i = P(i \in S^*) = \Theta(s/p)$, and $p_{ij} = P(i, j \in S^*) = \Theta(s^2/p^2)$. Given the support, \mathbf{z}_S^* is i.i.d., has symmetric probability distribution density function, $\mathbb{E}[\mathbf{z}_i^* \mid i \in S^*] = 0$ and $\mathbb{E}[\mathbf{z}_S^* \mathbf{z}_S^{*T} \mid S^*] = \mathbf{I}$. Moreover, the non-zero entries of the code are sub-Gaussian and lower bounded, i.e., for $i \in S^*$, $|\mathbf{z}_i^*| \geq C_{\min}$ where $C_{\min} \leq 1$.

Assumption 3.4 (Generating dictionary). \mathbf{D}^* is μ -incoherent (see Definition A.1) where $\mu = \mathcal{O}(\log(m))$. \mathbf{D}^* is unit-norm columns matrix ($\|\mathbf{D}_i^*\|_2 = 1$), $\|\mathbf{D}^*\|_2 = \mathcal{O}(\sqrt{p/m})$, and $p = \mathcal{O}(m)$.

To achieve model recovery using gradient descent, we assume an appropriate dictionary initialization, i.e.,

Assumption 3.5 (Dictionary closeness). The initial dictionary $\mathbf{D}^{(0)}$ is $(\delta_0, 2)$ -close to \mathbf{D}^* (see Definition A.2). The dictionary closeness at every update is denoted by $\|\mathbf{D}_j^{(l)} - \mathbf{D}_j^*\|_2 \leq \delta_l \forall j$. Furthermore, $\delta_l = \mathcal{O}^*(1/\log p)$.

Given the μ -incoherence of \mathbf{D}^* (Assumption 3.4) and δ_l -closeness of the dictionary, $\mathbf{D}^{(l)}$ is μ_l -incoherent, i.e., **Lemma 3.1** (μ_l -incoherent). $\mathbf{D}^{(l)}$ is μ_l -incoherent where $\mu_l = \mu/\sqrt{m} + 2\delta_l$.

The recurrent encoder and decoder, which perform the computations shown in Algorithm 2, use the loss \mathcal{L} and proximal operator $\mathcal{P}_{\alpha h}(v) \triangleq \text{sign}(v) \max(|v| - \alpha\lambda, 0)$ for the ℓ_1 norm $h: \mathbb{R}^p \rightarrow \mathbb{R}$. The encoder implements ISTA (Daubechies et al., 2004; Blumensath & Davies, 2008) with step size α , assumed to be less than $1/\sigma_{\max}^2(\mathbf{D})$. With infinite encoder unrolling, the encoder’s output is the solution to the lasso (1), following the optimality condition (Lemma A.3) where we denote $f_{\mathbf{x}}(\mathbf{z}, \mathbf{D}) \triangleq \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D}) + h(\mathbf{z})$. One immediate observation is that $\lambda \geq \|\mathbf{D}^T \mathbf{x}\|_{\infty} \Leftrightarrow \{\mathbf{0}\} \in \arg \min f_{\mathbf{x}}(\mathbf{z}, \mathbf{D})$. We assume $\lambda < \|\mathbf{D}^T \mathbf{x}\|_{\infty}$. We specify in Theorem 4.1 and Theorem 4.2 the conditions on λ at every encoder iteration to ensure support recovery and its preservation through the encoder. In case of a constant λ across encoder iterations while using \mathbf{D}^* as the dictionary (i.e., sparse coding using ℓ_1 norm), the network recovers a biased code $\hat{\mathbf{z}}^*$. We denote this amplitude error in the code by $\hat{\delta}^* \triangleq \|\hat{\mathbf{z}}^* - \mathbf{z}^*\|_2$ which is small and goes to zero with λ decaying through the encoder.

In addition, we assume the solution to (1) is unique; sufficient conditions for uniqueness in the overcomplete case (i.e., $p > m$) are extensively studied in the literature (Wainwright, 2009; Candès & Plan, 2009; Tibshirani, 2013). Tibshirani (2013) discussed that the solution is unique with probability one if entries of \mathbf{D} are drawn from a continuous probability distribution (Tibshirani, 2013) (Assumption 3.6). This assumption implies that $\mathbf{D}_S^T \mathbf{D}_S$ is full-rank. We argue that as long as the data $\mathbf{x} \in \mathcal{X}$ are sampled from a continuous distribution, this assumption holds for the entire learning process. The preservation of this property is guaranteed at all iterations of the alternating minimization proposed in (Agarwal et al., 2014). Moreover, this assumption has been previously considered in analyses of unrolled sparse coding networks (Ablin et al., 2019; Malézieux et al., 2022) and can be extended to ℓ_1 -based optimization problems (Tibshirani, 2013; Rosset et al., 2004).

Assumption 3.6 (Lasso uniqueness). *The entries of the dictionary \mathbf{D} are continuously distributed. Hence, the minimizer of (1) is unique, i.e., $\{\hat{\mathbf{z}}\} \in \arg \min f_{\mathbf{x}}(\mathbf{z}, \mathbf{D})$ with probability one.*

Lemma 3.2 states the fixed-point property of the encoder recursion (Parikh & Boyd, 2014). Given the definitions for Lipschitz and Lipschitz differentiable functions, (Definitions A.3 and A.4), the loss \mathcal{L} and function h satisfy following Lipschitz properties.

Lemma 3.2 (Fixed-point property of lasso). *Given Assumption 3.6, we have $\mathbf{0} \in \nabla_1 \mathcal{L}(\hat{\mathbf{z}}, \mathbf{D}) + \partial h(\hat{\mathbf{z}})$. The minimizer is a fixed-point of the mapping, i.e., $\hat{\mathbf{z}} = \mathcal{P}_{\alpha h}(\hat{\mathbf{z}} - \alpha \nabla_1 \mathcal{L}(\hat{\mathbf{z}}, \mathbf{D})) = \Phi(\hat{\mathbf{z}})$ (Parikh & Boyd, 2014).*

Lemma 3.3 (Lipschitz differentiable least squares). *Given $\mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D}) = \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|_2^2$, \mathcal{D} , and Assumption 3.2, the loss is Lipschitz differentiable. Let L_1 and L_2 denote the Lipschitz constants of the first derivatives $\nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D})$ and $\nabla_2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D})$, L_{11} and L_{21} the Lipschitz constants of the second derivatives $\nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D})$ and $\nabla_{21}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D})$, all w.r.t \mathbf{z} . Let $\nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D})$ be L_D -Lipschitz w.r.t \mathbf{D} .*

Lemma 3.4 (Lipschitz proximal). *Given $h(\mathbf{z}) = \lambda \|\mathbf{z}\|_1$, its proximal operator has bounded sub-derivative, i.e., $\|\partial \mathcal{P}_h(\mathbf{z})\|_2 \leq c_{\text{prox}}$.*

4 Unrolled Dictionary Learning

The gradients defined in PUDLE (Algorithm 2) can be compared against the local direction at each update of classical alternating-minimization (Algorithm 1). Assuming there are infinite samples, i.e.,

$$\text{Best local direction : } \hat{\mathbf{g}} \triangleq \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \nabla_2 \mathcal{L}_{\mathbf{x}^i}(\hat{\mathbf{z}}^i, \mathbf{D}) = \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [\nabla_2 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D})] \quad (4)$$

where $\hat{\mathbf{z}} = \arg \min_{\mathbf{z} \in \mathbb{R}^p} \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D}) + h(\mathbf{z})$. Additionally, to assess the estimators for model recovery, hence dictionary learning, we compare them against gradient pointing towards \mathbf{D}^* , namely

$$\text{Desired global gradient for } \mathbf{D}^* : \mathbf{g}^* \triangleq \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \nabla_2 \mathcal{L}_{\mathbf{x}^i}(\mathbf{z}^{i*}, \mathbf{D}) = \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [\nabla_2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D})]. \quad (5)$$

To see why the above is the desired direction, $(\mathbf{z}^*, \mathbf{D}^*)$ is a critical point of the loss \mathcal{L} which reaches zero for data following the model (2). Hence, to reach $\mathbf{D}^* \in \arg \min_{\mathbf{D} \in \mathcal{D}} \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [\mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D})]$, we move towards the direction minimizing the loss in expectation. Specifically, using the gradient $\nabla_2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D}) = -(\mathbf{x} - \mathbf{D}\mathbf{z}^*)\mathbf{z}^{*\top} = (\mathbf{D} - \mathbf{D}^*)\mathbf{z}^*\mathbf{z}^{*\top}$ as a descent direction, we move from \mathbf{D} toward \mathbf{D}^* modulo the code presence matrix $\mathbf{z}^*\mathbf{z}^{*\top}$.

Algorithm 1: Classical alternating-minimization-based dictionary learning using lasso (1).

Initialize: Samples $\{\mathbf{x}^i\}_{i=1}^n \in \mathcal{X}$, initial dictionary $\mathbf{D}^{(0)}$

Repeat: $l = 0, 1, \dots$, number of epochs

Sparse coding step: $\mathbf{z}^{i(l)} = \arg \min_{\mathbf{z}} \mathcal{L}_{\mathbf{x}^i}(\mathbf{z}, \mathbf{D}^{(l)}) + h(\mathbf{z})$, (for $i \in [1, n]$)

Dictionary update: $\mathbf{D}^{(l+1)} = \mathbf{D}^{(l)} - \eta \hat{\mathbf{g}}^{(l)}$ where $\hat{\mathbf{g}}^{(l)} \triangleq \frac{1}{n} \sum_{i=1}^n \nabla_2 \mathcal{L}_{\mathbf{x}^i}(\mathbf{z}^{i(l)}, \mathbf{D}^{(l)})$

Algorithm 2: PUDLE: Provable unrolled dictionary learning framework.

Initialize: Samples $\{\mathbf{x}^i\}_{i=1}^n \in \mathcal{X}$, initial dictionary $\mathbf{D}^{(0)}$, and $\mathbf{z}_0 = \mathbf{0}$.

Repeat: $l = 0, 1, \dots$, number of epochs

Forward pass: (for $i \in [1, n]$)

$$\begin{aligned} \text{Encoder: } \mathbf{z}_{t+1}^{i(l)} &= \Phi(\mathbf{z}_t^{i(l)}, \mathbf{D}^{(l)}) = \mathcal{P}_{\alpha h}(\mathbf{z}_t^{i(l)} - \alpha \nabla_1 \mathcal{L}_{\mathbf{x}^i}(\mathbf{z}_t^{i(l)}, \mathbf{D}^{(l)})) \text{ (repeat for } T) \\ \text{Decoder: } \hat{\mathbf{x}}^{i(l)} &= \mathbf{D}^{(l)} \mathbf{z}_T^{i(l)} \end{aligned} \quad (6)$$

Backward pass: $\mathbf{D}^{(l+1)} = \mathbf{D}^{(l)} - \eta \mathbf{g}_T^{(l)}$ where $\mathbf{g}_T^{(l)}$ is either of

$$\begin{aligned} \mathbf{g}_t^{(l) \text{ dec}} &\triangleq \frac{1}{n} \sum_{i=1}^n \nabla_2 \mathcal{L}_{\mathbf{x}^i}(\mathbf{z}_t^{i(l)}, \mathbf{D}^{(l)}) \\ \mathbf{g}_t^{(l) \text{ ae-lasso}} &\triangleq \frac{1}{n} \sum_{i=1}^n \nabla_2 \mathcal{L}_{\mathbf{x}^i}(\mathbf{z}_t^{i(l)}, \mathbf{D}^{(l)}) + \frac{\partial \mathbf{z}_t^{i(l)}}{\partial \mathbf{D}^{(l)}} \left(\nabla_1 \mathcal{L}_{\mathbf{x}^i}(\mathbf{z}_t^{i(l)}, \mathbf{D}^{(l)}) + \partial h(\mathbf{z}_t^{i(l)}) \right) \\ \mathbf{g}_t^{(l) \text{ ae-ls}} &\triangleq \frac{1}{n} \sum_{i=1}^n \nabla_2 \mathcal{L}_{\mathbf{x}^i}(\mathbf{z}_t^{i(l)}, \mathbf{D}^{(l)}) + \frac{\partial \mathbf{z}_t^{i(l)}}{\partial \mathbf{D}^{(l)}} \nabla_1 \mathcal{L}_{\mathbf{x}^i}(\mathbf{z}_t^{i(l)}, \mathbf{D}^{(l)}) \end{aligned} \quad (7)$$

Given these directions, we analyze the error of the gradients $\mathbf{g}_t^{\text{dec}}$, $\mathbf{g}_t^{\text{ae-lasso}}$, and $\mathbf{g}_t^{\text{ae-ls}}$ assuming infinite samples. In local analysis, we compare the code and gradient estimates to the lasso optimization in each update of the alternating minimization. In global analysis, we evaluate the performance in recovery of the ground-truth code \mathbf{z}^* and the dictionary \mathbf{D}^* . In this regard, we first study the forward pass.

4.1 Forward pass

We show convergence results in the forward pass for \mathbf{z} and the Jacobian, i.e.,

Definition 4.1 (Code Jacobian). *Given \mathbf{D} , the Jacobian of \mathbf{z}_t is defined as $\mathbf{J}_t \triangleq \frac{\partial \mathbf{z}_t}{\partial \mathbf{D}}$.*

The forward pass analyses give upper bounds on the error between \mathbf{z}_t and $\hat{\mathbf{z}}$ and the error between \mathbf{J}_t and $\hat{\mathbf{J}}$ as a function of unrolled iterations t . We will require these errors in [Section 4.2](#), where we analyze the gradient estimation errors. Similar to ([Chatterji & Bartlett, 2017](#)), the error associated with $\mathbf{g}_t^{\text{dec}}$ depends on the code convergence. Unlike $\mathbf{g}_t^{\text{dec}}$, the convergence of backpropagation with gradient estimates $\mathbf{g}_t^{\text{ae-lasso}}$ and $\mathbf{g}_t^{\text{ae-ls}}$ relies on the convergence properties of the code and the Jacobian ([Ablin et al., 2020](#)). Forward-pass theories are based on studies by [Gilbert \(1992\)](#) on the convergence of variables and their derivatives in an iterative process governed by a smooth operator ([Gilbert, 1992](#)). Moreover, [Hale et al. \(2007\)](#) studied the convergence analysis of fixed point iterations for ℓ_1 regularized optimization problems ([Hale et al., 2007](#)).

Support recovery and preservation We re-state a result from ([Hale et al., 2007](#)) on support selection.

Proposition 4.1 (Finite-iteration support selection). *Given [Assumption 3.6](#), let $\hat{\mathbf{z}} = \arg \min_{\mathbf{z}} f_{\mathbf{x}}(\mathbf{z}, \mathbf{D})$ with $S \triangleq \text{supp}(\hat{\mathbf{z}})$. There exists a $B > 0$ such that $\text{supp}(\mathbf{z}_B) = S, \forall t > B$.*

This means the unrolled encoder identifies the support in finite iterations. Support recovery in finite iterations has been studied in the literature for LISTA ([Chen et al., 2018](#)), Step-LISTA ([Ablin et al., 2019](#)), and shallow autoencoders ([Arora et al., 2015](#); [Rangamani et al., 2018](#); [Nguyen et al., 2019](#); [Tolooshams et al., 2020](#)). We show that under proper initialization of the dictionary, the encoder achieves linear convergence. [Arora et al. \(2015\)](#) discussed some appropriate initialization which is used by [Rambhatla et al. \(2018\)](#). Given initial closeness δ_0 , the encoder selects and recovers the correct signed support of the code with high probability in one iteration $B = 1$ ([Theorem 4.1](#)), and the iterations preserve the correct support ([Theorem 4.2](#)).

Theorem 4.1 (Forward pass support recovery). *Given the [Assumption 3.3](#), [Assumption 3.4](#), suppose $\mathbf{D}^{(l)}$ is $\delta_l = \mathcal{O}^*(1/\sqrt{\log p})$ close to \mathbf{D}^* . If $s = \mathcal{O}^*(\sqrt{m}/\mu \log m)$, and $\mu = \mathcal{O}(\log m)$, then with high probability of at least $1 - \epsilon_{\text{supp-rec}}^{(l)}$, the choice of $\lambda_0 = C_{\min}/4$ recovers the support of the code \mathbf{z}^* in one encoder iteration, i.e., $\text{sign}(\text{ReLU}(\alpha(\mathbf{D}^{(l)T} \mathbf{x} - \lambda_0))) = \text{sign}(\mathbf{z}^*)$, where $\epsilon_{\text{supp-rec}}^{(l)} = 2p \exp(-\frac{C_{\min}^2}{\mathcal{O}^*(\delta_l^2)})$.*

Theorem 4.2 (Forward pass support preservation). *Given the [Assumption 3.3](#), [Assumption 3.4](#), suppose $\mathbf{D}^{(l)}$ is $\delta_l = \mathcal{O}^*(1/\log p)$ close to \mathbf{D}^* . If $s = \mathcal{O}^*(\sqrt{m}/\mu \log m)$, $\mu = \mathcal{O}(\log m)$, and the regularizer and step size are chosen such that $\lambda_t^{(l)} = \frac{\mu_l}{\sqrt{m}} \|\mathbf{z}^* - \mathbf{z}_t\|_1 + a_\gamma = \Omega(\frac{s \log m}{\sqrt{m}})$ and $\alpha^{(l)} \leq 1 - \frac{2\lambda_t - (1 - \frac{\delta_l^2}{2})C_{\min}}{\lambda_{t-1}}$, then with high probability of at least $1 - \epsilon_{\text{supp-pres}}^{(l)}$, the support, recovered at the first iteration, is preserved through the encoder iterations. We have $a_\gamma = \mathcal{O}(\sqrt{s\delta_l})$ and $\epsilon_{\text{supp-pres}}^{(l)} := \epsilon_{\text{supp-rec}}^{(l)} + \epsilon_\gamma^{(l)} = 2p \exp(-\frac{C_{\min}^2}{\mathcal{O}^*(\delta_l^2)}) + 2s \exp(-\frac{1}{\mathcal{O}(\delta_l)})$.*

The support preservation conditions on λ_t and α introduce two insights. First, with an increase of t , the code error decrease, hence the lower bound on λ_t . Second, the decay of λ_t as the encoder unrolls increases the upper bound on α . Hence, we suggest a decaying strategy in values of λ_t as t increases.

Code convergence and error Given the support recovery and its preservation, the encoder convergence studied in ([Malézieux et al., 2022](#)) can achieve linear convergences after its first iteration. We re-state this result on the rate of convergence of the encoder in [Theorem 4.3](#). We drop the superscript (l) to simplify the notation.

Theorem 4.3 (Local forward pass code convergence). *Given the encoder $\mathbf{z}_{t+1} = \Phi(\mathbf{z}_t, \mathbf{D})$, [Assumption 3.6](#), [Lemmas 3.2, A.1 and A.2](#), then $\exists \rho < 1, B > 0$ s.t. $\|\mathbf{z}_t - \hat{\mathbf{z}}\|_2 \leq \mathcal{O}(\rho^t) \forall t > B$, where $\hat{\mathbf{z}}$ is the unique minimizer of lasso (1). Furthermore, given [Theorem 4.1](#) and [Theorem 4.2](#), $B = 1$.*

[Theorem 4.3](#) shows that in PUDLE, \mathbf{z}_t converges to $\hat{\mathbf{z}}$ at a linear rate eventually after a certain number of unrolling ([Figure 2](#)). The local linear convergence of ISTA and FISTA ([Beck & Teboulle, 2009](#)) (with global rates of $\mathcal{O}(1/t)$ and $\mathcal{O}(1/t^2)$) in the neighbourhood of a fixed-point is studied in ([Tao et al., 2016](#)). The speed of convergence depends on when support selection happens ([Proposition 4.1](#)) ([Bredies & Lorenz, 2008](#); [Zhang et al., 2017](#)). We showed in [Theorem 4.1](#) and [Theorem 4.2](#) that under mild assumptions, the support is selected and recovered after one encoder iteration.

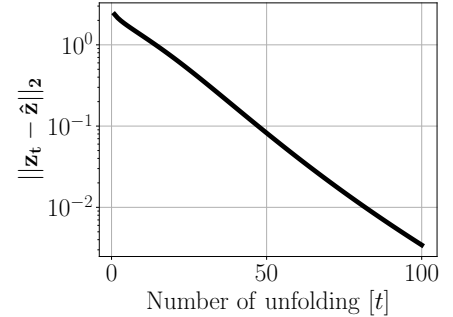


Figure 2: Code convergence ([Theorem 4.3](#)). As the network unrolls, \mathbf{z}_t converges to $\hat{\mathbf{z}}$, the solution of lasso.

Next, we focus on recovery of \mathbf{z}^* and prove the upper bound on the error between the converged code and \mathbf{z}^* can be decomposed into two terms: the dictionary error and the biased amplitude estimate of the code.

Theorem 4.4 (Global forward pass code error). *Let $\hat{\mathbf{z}}$ be the fixed-point of the encoder with iterations $\mathbf{z}_{t+1} = \Phi(\mathbf{z}_t, \mathbf{D})$. Given [Assumption 3.6](#), [Lemmas 3.2, A.1 and A.2](#), we have $\|\hat{\mathbf{z}} - \mathbf{z}^*\|_2 \leq \mathcal{O}(\|\mathbf{D} - \mathbf{D}^*\|_2 + \hat{\delta}^*)$, where $\hat{\delta}^* = \|\hat{\mathbf{z}}^* - \mathbf{z}^*\|_2$, $\hat{\mathbf{z}}^*$ is the unique minimizer of lasso (1) given the dictionary \mathbf{D} , $\hat{\mathbf{z}}^*$ is the unique minimizer of lasso (1) given the dictionary \mathbf{D}^* , and \mathbf{z}^* is the ground-truth code.*

This general decomposition is to emphasize that aside from the current estimate of the dictionary, the code error is a function of the forward pass algorithm used to solve the sparse coding problem. Specifically, the upper bound states that at the best scenario where there is access to the generating dictionary \mathbf{D}^* , the forward pass solving lasso with fixed λ still gives a biased amplitude estimate of \mathbf{z}^* . Overall, the assumptions to get this bound are mild; the bound is valid independent of successful support recovery or data distribution. With incorporation of data distribution and conditions stated in [Theorem 4.1](#) and [Theorem 4.2](#), the upper bound $\hat{\delta}^*$ can be replaced with terms involving λ , and reaches at zero as λ decays across forward iterations.

Jacobian convergence and error Following properties similar to those used in [Theorem 4.3](#), and assuming \mathbf{J}_t is bounded ([Assumption 4.1](#)), we show in [Theorem 4.5](#) that, as the PUDLE unrolls, the code Jacobian \mathbf{J}_t converges to $\hat{\mathbf{J}}$, the Jacobian of the solution of the lasso. The convergence of the Jacobian of proximal

gradient descent is also studied in (Bertrand et al., 2021) for hyperparameter selection through implicit differentiation (Bengio, 2000), where the Jacobian is taken w.r.t to the hyperparameter λ as opposed to \mathbf{D} .

Assumption 4.1 (Bounded Jacobian). *The Jacobian is bounded, i.e., $\exists M_J > 0$, s.t. $\|\mathbf{J}_t\|_2 \leq M_J \forall t$.*

Theorem 4.5 (Local forward pass Jacobian convergence). *Given the recursion $\mathbf{z}_{t+1} = \Phi(\mathbf{z}_t, \mathbf{D})$, and $\hat{\mathbf{z}}$ the unique minimizer of lasso with Jacobian $\hat{\mathbf{J}}$, then $\exists \rho < 1, B > 0$ s.t. $\|\mathbf{J}_t - \hat{\mathbf{J}}\|_2 \leq \mathcal{O}(t\rho^t) \forall t > B$. Furthermore, given Theorem 4.1 and Theorem 4.2, $B = 1$.*

The forward pass code and Jacobian convergences *after* support selection is similar to the results from (Malézieux et al., 2022). The highlights of our finding are that the order of upper bound convergences can be achieved from the first iteration of the encoder. In other words, we specify, in Theorem 4.1 and Theorem 4.2, the dictionary and data conditions such that the support can be recovered with $B = 1$. This resolves the instability issue discussed by Malézieux et al. (2022) in computation of the gradient $\mathbf{g}_t^{\text{ae-lasso}}$ outside of the support. Finally, we show that the global Jacobian error is in the order of dictionary error.

Theorem 4.6 (Global forward pass Jacobian error). *Let $\hat{\mathbf{z}}$ be the fixed-point of the encoder with iterations $\mathbf{z}_{t+1} = \Phi(\mathbf{z}_t, \mathbf{D})$. Given Assumption 3.6, Lemmas 3.2, A.1 and A.2, we have $\|\hat{\mathbf{J}} - \mathbf{J}^*\|_2 \leq \mathcal{O}(\|\mathbf{D} - \mathbf{D}^*\|_2 + \hat{\delta}^*)$, where $\hat{\delta}^* = \|\hat{\mathbf{z}}^* - \mathbf{z}^*\|_2$, $\hat{\mathbf{z}}$ is the unique minimizer of lasso (1) given \mathbf{D} , $\hat{\mathbf{z}}^*$ is the unique minimizer of lasso (1) given \mathbf{D}^* , and \mathbf{z}^* is the ground-truth code. Moreover, $\hat{\mathbf{J}}$ and \mathbf{J}^* are Jacobians of $\hat{\mathbf{z}}$ and \mathbf{z}^* , respectively.*

4.2 Backward pass

We show two results for local gradient $\hat{\mathbf{g}}$ and global gradient \mathbf{g}^* convergence. The goal is not to provide a finite sample analysis but to emphasize the relative differences between the gradients in Algorithm 2. The impact of gradient error for parameter estimation in the convex setting has been studied by Devolder et al. (2013) indicating that the convergence to the parameter’s neighbourhood is dictated by the gradient error (Devolder et al., 2013; 2014). As dictionary learning is a bi-convex problem, findings of Devolder et al. (2013) hold as well for better estimation of the local dictionary at every step of alternating minimization. Moreover, Arora et al. (2015), provided a detailed analysis of sparse coding and various gradient estimations for dictionary learning, showing that by computing a more accurate gradient at every step of the alternating minimization scheme, the dictionary estimates converge to a closer neighbourhood of \mathbf{D}^* . Overall, the intuition is that the size of the gradient error dictates the size of the neighbourhood of the dictionary within which one can guarantee convergence. We argue that the method with lower gradient error recovers the dictionary better.

Local gradient estimations We highlight the effect of finite unrolling on the gradient for parameter estimation (Abblin et al., 2020). Theorem 4.7 shows the convergence rate of gradients to $\hat{\mathbf{g}}$, determining the similarity of PUDLE and Algorithm 1.

Theorem 4.7 (Local convergence of gradients). *Given the convergence results from the forward pass (Theorems 4.3 and 4.5), $\exists \rho < 1, B > 0$ such that $\forall t > B$, the errors of gradients defined in Algorithm 2 w.r.t $\hat{\mathbf{g}}$ (4) satisfy*

$$\begin{aligned} \|\mathbf{g}_t^{\text{dec}} - \hat{\mathbf{g}}\|_2 &\leq \mathcal{O}(\rho^t) \\ \|\mathbf{g}_t^{\text{ae-lasso}} - \hat{\mathbf{g}}\|_2 &\leq \mathcal{O}(t\rho^{2t}) \\ \|\mathbf{g}_t^{\text{ae-ls}} - \hat{\mathbf{g}}\|_2 &\leq \mathcal{O}(t\rho^{2t} + M_J\lambda\sqrt{s}). \end{aligned} \quad (8)$$

Moreover, the order of upper bounds is tight (see Lemma A.4).

First, upper bounds on the errors related to $\mathbf{g}_t^{\text{dec}}$ and $\mathbf{g}_t^{\text{ae-lasso}}$ go to zero as t increases. Hence, both gradients converge to $\hat{\mathbf{g}}$. This means that asymptotically as t increases, training PUDLE with $\mathbf{g}_t^{\text{dec}}$ and $\mathbf{g}_t^{\text{ae-lasso}}$ is equivalent to classical alternating-minimization (Algorithm 1). Second, as t increases, $\mathbf{g}_t^{\text{ae-lasso}}$ has faster convergence than $\mathbf{g}_t^{\text{dec}}$. Lastly, $\mathbf{g}_t^{\text{ae-ls}}$ is a biased estimator of $\hat{\mathbf{g}}$ (Figure 3). The convergence results on the error $\|\mathbf{g}_t^{\text{ae-lasso}} - \hat{\mathbf{g}}\|_2$ is previously studied by Malézieux et al. (2022).

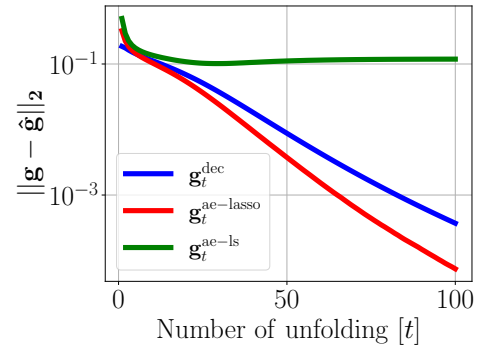


Figure 3: Convergence rate of gradients (Theorem 4.7).

Given the above convergence results, one may conclude that $\mathbf{g}_t^{\text{ae-lasso}}$ should be used for dictionary recovery. However, we show next that for dictionary recovery, the gradient $\mathbf{g}_t^{\text{ae-lasso}}$, used by Malézieux et al. (2022), is indeed a biased estimator of the global gradient \mathbf{g}^* for recovery of \mathbf{D}^* . We decrease this bias by replacing $\mathbf{g}_t^{\text{ae-lasso}}$ with $\mathbf{g}_t^{\text{ae-ls}}$ and show that $\mathbf{g}_t^{\text{ae-ls}}$ results in a better recovery of \mathbf{D}^* than $\mathbf{g}_t^{\text{ae-lasso}}$.

Global gradient estimations Theorem 4.8 shows the global gradient errors w.r.t \mathbf{g}^* from (5). We omit the gradient $\mathbf{g}_t^{\text{dec}}$, as it is asymptotically equivalent to $\mathbf{g}_t^{\text{ae-lasso}}$. We study the errors in the limit to unrolling, i.e., as $t \rightarrow \infty$. This determines which PUDLE gradients recover \mathbf{D}^* better (Devolder et al., 2013; 2014).

Theorem 4.8 (Global error of gradients). *Given the convergence results from the forward pass, (Theorems 4.4 and 4.6), the errors of gradients defined in Algorithm 2 w.r.t global direction \mathbf{g}^* (defined in (5)) satisfy*

$$\begin{aligned} \|\mathbf{g}_\infty^{\text{ae-lasso}} - \mathbf{g}^*\|_2 &\leq \mathcal{O}(\|\mathbf{D} - \mathbf{D}^*\|_2^2 + \|\mathbf{D} - \mathbf{D}^*\|_2 + \|\mathbf{D} - \mathbf{D}^*\|_2 \hat{\delta}^* + \hat{\delta}^* + \hat{\delta}^{*2} + M_J \lambda \sqrt{s}) \\ \|\mathbf{g}_\infty^{\text{ae-ls}} - \mathbf{g}^*\|_2 &\leq \mathcal{O}(\|\mathbf{D} - \mathbf{D}^*\|_2^2 + \|\mathbf{D} - \mathbf{D}^*\|_2 + \|\mathbf{D} - \mathbf{D}^*\|_2 \hat{\delta}^* + \hat{\delta}^* + \hat{\delta}^{*2}). \end{aligned} \quad (9)$$

Several factors affect the order of upper bounds: the current estimate of the dictionary, code amplitude-bias error due to ℓ_1 norm, and the usage of ℓ_1 norm in the loss used for backpropagation. To study the bias in the gradient computation, let consider the scenario where $\mathbf{D} = \mathbf{D}^*$. We denote those gradients by superscript \mathbf{D}^* . If the gradients are not biased, then the upper bounds should goes to zero. The gradient errors are

$$\|\mathbf{g}_\infty^{\text{ae-lasso}, \mathbf{D}^*} - \mathbf{g}^*\|_2 \leq \mathcal{O}(\hat{\delta}^* + \hat{\delta}^{*2} + M_J \lambda \sqrt{s}) \quad \text{and} \quad \|\mathbf{g}_\infty^{\text{ae-ls}, \mathbf{D}^*} - \mathbf{g}^*\|_2 \leq \mathcal{O}(\hat{\delta}^* + \hat{\delta}^{*2}). \quad (10)$$

For $\mathbf{g}_\infty^{\text{ae-ls}, \mathbf{D}^*}$, the radius of the error ball is only a function of the amplitude error of the code estimated through lasso compare to the ground-truth code \mathbf{z}^* . However, the error ball for the gradient $\mathbf{g}_\infty^{\text{ae-lasso}, \mathbf{D}^*}$ includes an additional term concerning the usage of lasso loss containing the regularization term λ . This implies that the \mathbf{D}^* neighbourhood at which the gradient $\mathbf{g}_\infty^{\text{ae-ls}, \mathbf{D}^*}$ is guaranteed to converge to is smaller than of the $\mathbf{g}_\infty^{\text{ae-lasso}, \mathbf{D}^*}$ (Figure 4a). Implications of such gradient estimation are seen in dictionary learning where $\mathbf{g}_\infty^{\text{ae-ls}}$ recovers \mathbf{D}^* better (Figures 4b and 4c). In Figure 4b, the encoder unrolls for $T = 25$, hence the phenomenon of implicit acceleration is seen in faster and better dictionary learning performance of $\mathbf{g}_\infty^{\text{ae-lasso}}$ than $\mathbf{g}_\infty^{\text{dec}}$. In Figure 4c where $T = 100$, similar performance of $\mathbf{g}_\infty^{\text{dec}}$ and $\mathbf{g}_\infty^{\text{ae-lasso}}$ illustrates their asymptotic equivalence as $t \rightarrow \infty$.

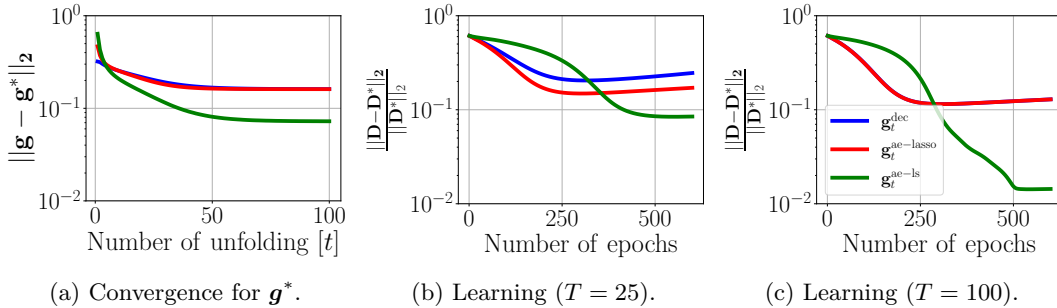


Figure 4: Results for PUDLE’s global convergence (Theorem 4.8) and dictionary learning.

Towards unbiased estimation As long as λ is fixed within PUDLE, all defined gradients remain biased estimators of \mathbf{g}^* , due to the biased estimate of the code \mathbf{z}^* through ℓ_1 norm. This bias exists while dictionary learning is performed strictly using lasso through Algorithm 1. Given the conditions on the regularizer in Theorem 4.2 which we discussed in Section 4.1 and the derived upper bounds in Theorem 4.8, we suggest the decaying of λ across the encoder to reduce the gradient biases and improve dictionary learning. We show in Section 4.3 that by decaying λ at each unrolled layer, this bias vanishes, and the training converges to \mathbf{D}^* .

4.3 Experiments

Dictionary learning We focus on the performance of the best-performing gradient estimator $\mathbf{g}_t^{\text{ae-ls}}$, and compare it with NOODL (Rambhatla et al., 2018), a state-of-the-art online dictionary learning algorithm, and SPORCO (Wohlberg, 2017), an alternating-minimization dictionary learning algorithm that

uses lasso. NOODL, which uses iterative hard-thresholding (HT) for sparse coding and a gradient update employing the code’s sign, has linear convergence upon proper initialization (Rambhatla et al., 2018). We note that the results from $\mathbf{g}_t^{\text{ae-lasso}}$ are not shown, as the gradient computation was unstable (Malézieux et al., 2022). We emphasize that our proposed gradient $\mathbf{g}_t^{\text{ae-ls}}$ does not suffer such instability. We train:

- $\mathbf{g}_t^{\text{ae-ls}}$: λ is fixed across iterations.
- $\mathbf{g}_t^{\text{ae-ls, decay}}$: λ decays (i.e., $\lambda_t = \lambda \nu^t$, with $0 < \nu < 1$) where ν decreases as training progresses.
- $\mathbf{g}_t^{\text{ae-ls, HT}}$: $\mathcal{P}_{\alpha h}(v)$ is replaced with $\text{HT}_b(v) \triangleq v \mathbf{1}_{|v| \geq b}$.

With HT, the sparse coding step reduces to that from NOODL. In this case, we highlight the difference between the gradient update of our method (backpropagation) with NOODL. We focus on convergence, as η across methods is not comparable.

Figure 5 shows the convergence of $\mathbf{D} \in \mathbb{R}^{1000 \times 1500}$ to \mathbf{D}^* when the code is 20-sparse (for other sparsity levels and details see Appendix C). A biased estimate of the code amplitudes results in convergence only to a neighbourhood of the dictionary (Rambhatla et al., 2018). This is observed in the convergence of $\mathbf{g}_t^{\text{ae-ls}}$ and SPORCO (final error is shown). The convergence of $\mathbf{g}_t^{\text{ae-ls}}$ to a closer neighbourhood than SPORCO supports Theorem 4.8. Moreover, with decaying λ , the code bias vanishes, hence $\mathbf{g}_t^{\text{ae-ls, decay}}$ and $\mathbf{g}_t^{\text{ae-ls, HT}}$ converges to \mathbf{D}^* similar to NOODL.

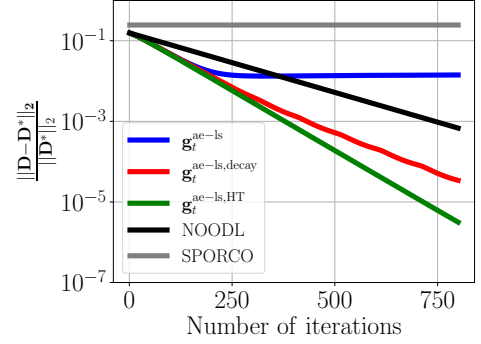


Figure 5: Dictionary convergences.

Image denoising To further highlight the advantage of $\mathbf{g}_t^{\text{ae-ls}}$ over the other gradients, we compare them in a supervised task of image denoising. We focus on $\mathbf{g}_t^{\text{ae-ls}}$ and $\mathbf{g}_t^{\text{dec}}$ to highlight the impact of backpropagation, and also consider $\mathbf{g}_t^{\text{ae-ls, HT}}$ where the proximal operator is replaced with HT. This is to compare with sparse coding scheme of NOODL. We do not compare against NOODL’s dictionary update, as this computation for two-dimensional convolutions is not straightforward. $\mathbf{g}_t^{\text{ae-ls}}$ uses full backpropagation and stays stable. This highlights stability of the gradient computation as opposed to $\mathbf{g}_t^{\text{ae-lasso}}$ (Malézieux et al. (2022)).

Prior works have shown that variants of PUDLE either rival or outperform state-of-the-art architectures (Simon & Elad, 2019; Tolooshams et al., 2020). Thus, we focus on a comparative analysis of the gradients. We trained on 432 and tested on 68 images from BSD (Martin et al., 2001). We used a convolutional dictionary and corrupted images with zero-mean Gaussian noise of standard deviation of 25 (see Appendix C for details). Table 1 shows the denoising performance of soft-thresholding using λ and HT with b in peak signal-to-noise-ratio (PSNR). The result shows that the advantage of $\mathbf{g}_t^{\text{ae-ls}}$ over $\mathbf{g}_t^{\text{dec}}$ is not limited to dictionary learning and is seen in denoising. Additionally, the superior performance of $\mathbf{g}_t^{\text{ae-ls}}$ compared to $\mathbf{g}_t^{\text{ae-ls, HT}}$ highlights the benefits of PUDLE (i.e., ℓ_1 -based unrolling) against HT used in NOODL.

Table 1: Denoising of BSD68.

METHOD	PSNR [dB]				
	λ	0.08	0.12	0.16	0.2
$\mathbf{g}_t^{\text{dec}}$		24.31	24.73	25.24	24.89
$\mathbf{g}_t^{\text{ae-ls}}$		24.80	25.48	25.67	25.51
$\mathbf{g}_t^{\text{ae-ls, HT}}$	b	0.02	0.05	0.08	0.1
		22.87	25.40	24.68	23.77

5 Interpretable Sparse Codes and Dictionary

One motivation behind using algorithm unrolling to design deep architectures is interpretability (Monga et al., 2019); they argue that the designed networks are interpretable as they capture domain knowledge via an optimization model. For example, Tolooshams et al. (2021a) takes advantage of the interpretability of learned weights in an unrolled dictionary learning network to solve spike sorting, an unsupervised source separation problem in computational neuroscience. Moreover, Kim et al. (2010) uses sparse coding to learn interpretable representations of human motions. However, none of the existing methods in the literature provide interpretability results that open the black-box network through building a mathematical relation

between the learned dictionary, training data, and test representation/reconstruction. This section analyzes the interpretability of the unrolled sparse coding method in this context. We provide the following theorem.

Theorem 5.1 (Interpretable unrolled network). *Consider the dictionary learning optimization of the form $\min_{\mathbf{Z}, \mathbf{D}} \frac{1}{2} \|\mathbf{X} - \mathbf{D}\mathbf{Z}\|_2^2 + \lambda \|\mathbf{Z}\|_1 + \omega/2 \|\mathbf{D}\|_F^2$, where $\mathbf{X} = [\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^n] \in \mathbb{R}^{m \times n}$ and $\mathbf{Z} = [\mathbf{z}^1, \mathbf{z}^2, \dots, \mathbf{z}^n] \in \mathbb{R}^{p \times n}$. Let $\tilde{\mathbf{Z}}$ be the given converged sparse codes, then stationary points of the problem w.r.t the network weights (dictionary) follows $\tilde{\mathbf{D}} = \mathbf{X}\mathbf{G}^{-1}\tilde{\mathbf{Z}}^T$, where we denote $\mathbf{G} := (\tilde{\mathbf{Z}}^T\tilde{\mathbf{Z}} + \omega\mathbf{I})$.*

The dictionary interpolates the training data Given Theorem 5.1, each learned atom interpolates the training data, i.e.,

$$\tilde{\mathbf{D}}_j = \mathbf{X}(\mathbf{G}^{-1}\mathbf{w}_j) = \sum_{k=1}^n (\mathbf{G}^{-1}\mathbf{w}_j)_k \mathbf{x}^k \quad (11)$$

where $\mathbf{w}_j = [\tilde{z}_j^1, \tilde{z}_j^2, \dots, \tilde{z}_j^n]^T \in \mathbb{R}^n$ is a vector containing the training code activity for dictionary atom j . Specifically, the importance of training image \mathbf{x}^k in learning dictionary atom j is captured by the term $(\mathbf{G}^{-1}\mathbf{w}_j)_k$. This proves the dictionary spans the training set. Given the small number of atoms compared to the size of training set, (11) shows that the dictionary summarizes the training examples. We trained the network on digits of $\{0, 1, 2, 3, 4\}$ MNIST (Figure 6 shows a fraction of the most used learned atoms). Figure 7 visualizes dictionary atoms along with training images with the highest contribution (green color) and the lowest contribution (red color). In addition, we used (11) on the partial training data to reconstruct the learned training image (shown as Estimate). Next, we interpret the relation between the new test data to the training data using representer point selection, similar to (Yeh et al., 2018).

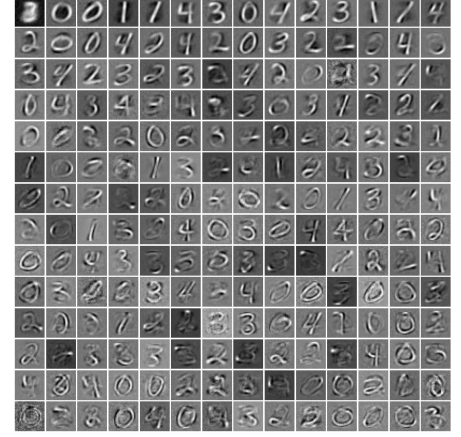


Figure 6: Fraction of dictionary atoms learned from $\{0, 1, 2, 3, 4\}$ MNIST.

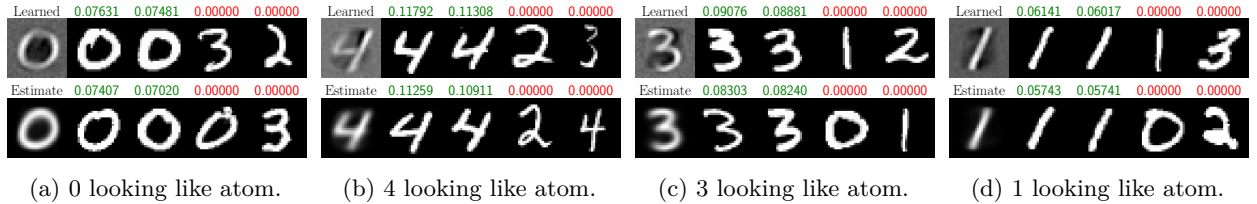


Figure 7: Training image contributions to learning the dictionary.

Relation between new test image and training data For representation of a new data, we observe that the reconstruction of a new example \mathbf{x}^j is a linear combination of all the training examples, i.e.,

$$\mathbf{x}^j = \tilde{\mathbf{D}}\mathbf{z}^j = \mathbf{X}\boldsymbol{\beta}^j = \sum_{k=1}^n \beta_k^j \mathbf{x}^k \quad (12)$$

where $\boldsymbol{\beta}^j = \mathbf{G}^{-1}\tilde{\mathbf{Z}}^T\mathbf{z}^j \in \mathbb{R}^n$ and $\beta_k^j = \sum_{a=1}^n \mathbf{G}_{ka}^{-1} \langle \mathbf{z}^a, \mathbf{z}_{\text{test}}^j \rangle$. We observe that the contribution of image k into the reconstruction of the test image is a function of β_k^j , and the energy of β_k^j itself depends on the whole training set, and \mathbf{G}^{-1} . (12) shows how each image is reconstructed as interpolation of the training images. Figure 8 shows this results, where images with high (green color) β_k^j contribution are similar to the test image and those with low (red color) β_k^j contribution are different. From another perspective, we can write the new image as

$$\mathbf{x}^j = \tilde{\mathbf{D}}\mathbf{z}^j = \sum_{k=1}^n (\mathbf{X}\mathbf{G}^{-1})_k \langle \tilde{\mathbf{z}}^k, \mathbf{z}^j \rangle \quad (13)$$

i.e., the contribution of each training image for reconstruction is a function of their code similarity to the new image and properties of the Gram matrix of training set code similarities. Specifically, the relation rules the contribution of transformed image k (i.e., $(\mathbf{X}\mathbf{G}^{-1})_k$) into reconstruction of the test image as a function of its

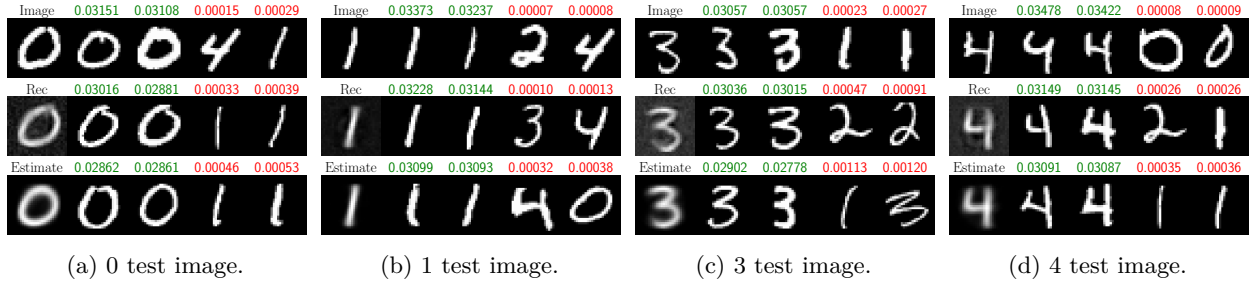


Figure 8: Interpolation of training data to reconstruct a new image. Contribution of training images are shown in green color (high contribution) and red color (low contribution).

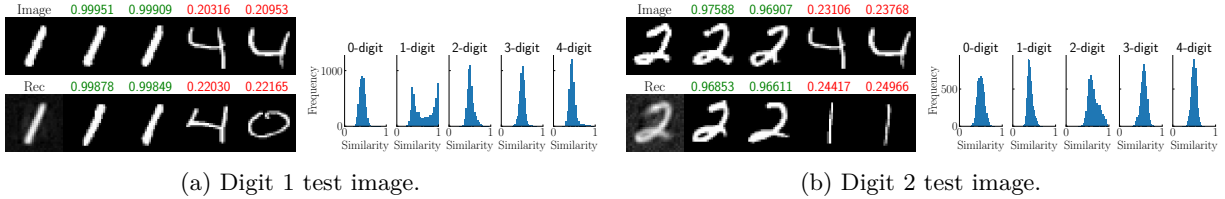


Figure 9: Contribution of images with code similarity into reconstruction of a new test image along with the histograms of the similarity of the test code to training codes from each class.

code similarity $\langle \mathbf{z}^{*k}, \mathbf{z}^j \rangle$. In other words, (13) shows that training images with the highest code similarity to the representation of the new image have the highest contribution to its reconstruction. This interpretation is demonstrated in Figure 9. The training images with the highest code similarity (green color) and the lowest similarity (red) are shown. In addition, the figure demonstrates the histogram of the code similarity between the test image and the training set, grouped by their class digit. For example, for digit 1 test image, its code similarity to train images from class 1 are bimodal. This corresponds to 1 digits that are tilted to the left (low similarity) and right (high similarity). Moreover, for digit 2 test image, we observe that the histogram of images corresponding to digit 2 are shifted the most to the right (highest similarity) than the other classes.

6 Conclusions

This paper studied dictionary learning and analyzed the dynamics of unrolled sparse coding networks through a provable unrolled dictionary learning (PUDLE) framework. First, we provided a theoretical analysis of the forward pass for code recovery. We discussed the bias introduced by ℓ_1 -based sparse coding in the forward pass, and how this affects the dictionary estimate in the backward pass. Second, we showed strategies to mitigate the propagation of this code bias into the backward pass; this is achieved by modification of the loss function used in the backward pass. We demonstrated that this bias could be further reduced and eliminated by decaying the regularization parameter within the unrolled layers. Additionally, we provided sufficient conditions on the data distribution and network initialization to guarantee stability of backpropagated gradient computations. In the absence of such conditions, we proposed a modification to the loss function that resolves the gradient explosion and allows stable learning. In an image denoising task, we showed PUDLE outperforms the NOODL sparse coding scheme (Rambhatla et al., 2018). Motivated by interpretability as a popular feature unrolled deep neural networks, we derived a mathematical relation between the network weights (dictionary) and the training dataset. We proved that the network weights span the training dataset, and constructed a relation between predictions of new test examples and the training set. The latter allows the user to extract images from the training set that are similar/dissimilar to the test image in representation/reconstruction.

References

Pierre Ablin, Thomas Moreau, Mathurin Massias, and Alexandre Gramfort. Learning step sizes for unfolded sparse coding. In *Proceedings of Advances in Neural Information Processing Systems*, volume 32, pp. 1–11, 2019.

- Pierre Ablin, Gabriel Peyré, and Thomas Moreau. Super-efficiency of automatic differentiation for functions defined as a minimum. In *Proceedings of International Conference on Machine Learning*, pp. 32–41. PMLR, 2020.
- Alekh Agarwal, Animashree Anandkumar, Prateek Jain, Praneeth Netrapalli, and Rashish Tandon. Learning sparsely used overcomplete dictionaries. In Maria Florina Balcan, Vitaly Feldman, and Csaba Szepesvári (eds.), *Proc the 27th Conference on Learning Theory*, volume 35 of *Proceedings of Machine Learning Research*, pp. 123–137, Barcelona, Spain, 13–15 Jun 2014. PMLR.
- M. Aharon, M. Elad, and A. Bruckstein. K-svd: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on Signal Processing*, 54(11):4311–4322, 2006.
- Kazunori Akiyama, Kazuki Kuramochi, Shiro Ikeda, Vincent L Fish, Fumie Tazaki, Mareki Honma, Sheperd S Doleman, Avery E Broderick, Jason Dexter, Monika Mościbrodzka, et al. Imaging the schwarzschild-radius-scale structure of m87 with the event horizon telescope using sparse modeling. *The Astrophysical Journal*, 838(1):1, 2017.
- Sanjeev Arora, Rong Ge, and Ankur Moitra. New algorithms for learning incoherent and overcomplete dictionaries. In Maria Florina Balcan, Vitaly Feldman, and Csaba Szepesvári (eds.), *Proceedings of the 27th Conference on Learning Theory*, volume 35 of *Proceedings of Machine Learning Research*, pp. 779–806, Barcelona, Spain, 13–15 Jun 2014. PMLR.
- Sanjeev Arora, Rong Ge, Tengyu Ma, and Ankur Moitra. Simple, efficient, and neural algorithms for sparse coding. In Peter Grünwald, Elad Hazan, and Satyen Kale (eds.), *Proceedings of Conference on Learning Theory*, volume 40 of *Proceedings of Machine Learning Research*, pp. 113–149, Paris, France, 03–06 Jul 2015. PMLR.
- Hedy Attouch and Jérôme Bolte. On the convergence of the proximal algorithm for nonsmooth functions involving analytic features. *Mathematical Programming*, 116(1):5–16, 2009.
- Waheed U. Bajwa, Kfir Gedalyahu, and Yonina C. Eldar. Identification of parametric underspread linear systems and super-resolution radar. *IEEE Transactions on Signal Processing*, 59(6):2548–2561, 2011.
- Boaz Barak, Jonathan A. Kelner, and David Steurer. Dictionary learning and tensor decomposition via the sum-of-squares method. In *Proceedings of Annual ACM Symposium on Theory of Computing*, STOC ’15, pp. 143–151, 2015.
- Atilim Gunes Baydin, Barak A Pearlmutter, Alexey Andreyevich Radul, and Jeffrey Mark Siskind. Automatic differentiation in machine learning: a survey. *Journal of machine learning research*, 18, 2018.
- Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
- Yoshua Bengio. Gradient-based optimization of hyperparameters. *Neural computation*, 12(8):1889–1900, 2000.
- Quentin Bertrand, Quentin Klopfenstein, Mathurin Massias, Mathieu Blondel, Samuel Vaiter, Alexandre Gramfort, and Joseph Salmon. Implicit differentiation for fast hyperparameter selection in non-smooth convex learning. *arXiv:2105.01637*, 2021.
- Thomas Blumensath and Mike E Davies. Iterative thresholding for sparse approximations. *Journal of Fourier analysis and Applications*, 14(5-6):629–654, 2008.
- Kristian Bredies and Dirk A Lorenz. Linear convergence of iterative soft-thresholding. *Journal of Fourier Analysis and Applications*, 14(5-6):813–837, 2008.
- Emmanuel J. Candès and Yaniv Plan. Near-ideal model selection by ℓ_1 minimization. *The Annals of Statistics*, 37(5A):2145 – 2177, 2009.

- Niladri S Chatterji and Peter L Bartlett. Alternating minimization for dictionary learning: Local convergence guarantees. *arXiv:1711.03634*, pp. 1–26, 2017.
- Scott Shaobing Chen, David L Donoho, and Michael A Saunders. Atomic decomposition by basis pursuit. *SIAM review*, 43(1):129–159, 2001.
- Xiaohan Chen, Jialin Liu, Zhangyang Wang, and Wotao Yin. Theoretical linear convergence of unfolded ista and its practical weights and thresholds. In *Proceedings of Advances in Neural Information Processing Systems*, volume 31, pp. 1–11, 2018.
- Brian Cleary, Le Cong, Anthea Cheung, Eric S. Lander, and Aviv Regev. Efficient generation of transcriptomic profiles by random composite measurements. *Cell*, 171(6):1424–1436.e18, 2017. ISSN 0092-8674.
- Brian Cleary, Brooke Simonton, Jon Bezney, Evan Murray, Shahul Alam, Anubhav Sinha, Ehsan Habibi, Jamie Marshall, Eric S Lander, Fei Chen, et al. Compressed sensing for highly efficient imaging transcriptomics. *Nature Biotechnology*, pp. 1–7, 2021.
- I. Daubechies, M. Defrise, and C. De Mol. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Mathematics*, 57(11):1413–1457, 2004.
- Olivier Devolder, François Glineur, Yurii Nesterov, et al. First-order methods with inexact oracle: the strongly convex case. Technical report, Université catholique de Louvain, Center for Operations Research and . . . , 2013.
- Olivier Devolder, François Glineur, and Yurii Nesterov. First-order methods of smooth convex optimization with inexact oracle. *Mathematical Programming*, 146(1):37–75, 2014.
- Michael Elad. *Sparse and redundant representations: from theory to applications in signal and image processing*. Springer Science & Business Media, 2010.
- Michael Elad and Michal Aharon. Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Transactions on Image Processing*, 15(12):3736–3745, 2006.
- K. Engan, S.O. Aase, and J. Hakon Husoy. Method of optimal directions for frame design. In *Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 5, pp. 2443–2446 vol.5, 1999.
- Matthias Feurer and Frank Hutter. Hyperparameter optimization. In *Automated Machine Learning*, pp. 3–33. Springer, Cham, 2019.
- Luca Franceschi, Michele Donini, Paolo Frasconi, and Massimiliano Pontil. Forward and reverse gradient-based hyperparameter optimization. In *Proceedings of International Conference on Machine Learning*, pp. 1165–1173, 2017.
- Jean Charles Gilbert. Automatic differentiation and iterative processes. *Optimization methods and software*, 1(1):13–21, 1992.
- Raja Giryes, Yonina C. Eldar, Alex M. Bronstein, and Guillermo Sapiro. Tradeoffs between convergence speed and reconstruction accuracy in inverse problems. *IEEE Transactions on Signal Processing*, 66(7):1676–1690, 2018.
- Karol Gregor and Yann LeCun. Learning fast approximations of sparse coding. In *Proceedings of international conference on international conference on machine learning*, pp. 399–406, 2010.
- Elaine T Hale, Wotao Yin, and Yin Zhang. A fixed-point continuation method for l1-regularized minimization with applications to compressed sensing. *CAAM TR07-07, Rice University*, 43:44, 2007.
- Trevor Hastie, Robert Tibshirani, and Martin Wainwright. *Statistical learning with sparsity: the lasso and generalizations*. CRC press, 2015.

- John R. Hershey, Jonathan Le Roux, and Felix Weninger. Deep unfolding: Model-based inspiration of novel deep architectures. *arXiv:1409.2574*, pp. 1–27, 2014.
- Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi. Low-rank matrix completion using alternating minimization. In *Proceedings of Annual ACM Symposium on Theory of Computing*, pp. 665–674, 2013. ISBN 9781450320290.
- Rodolphe Jenatton, Julien Mairal, Guillaume Obozinski, and Francis Bach. Proximal methods for hierarchical sparse coding. *The Journal of Machine Learning Research*, 12:2297–2334, 2011.
- Taehwan Kim, Gregory Shakhnarovich, and Raquel Urtasun. Sparse coding for learning interpretable spatio-temporal primitives. In J. Lafferty, C. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta (eds.), *Proceedings of Advances in Neural Information Processing Systems*, volume 23, 2010.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2014.
- Yann A LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient backprop. In *Neural networks: Tricks of the trade*, pp. 9–48. Springer, 2012.
- Yuelong Li, Mohammad Tofghi, Junyi Geng, Vishal Monga, and Yonina C. Eldar. Efficient and interpretable deep blind image deblurring via algorithm unrolling. *IEEE Transactions on Computational Imaging*, 6: 666–681, 2020.
- Jialin Liu and Xiaohan Chen. Alista: Analytic weights are as good as learned weights in lista. In *Proceedings of International Conference on Learning Representations*, 2019.
- Julien Mairal, Francis Bach, Jean Ponce, and Guillermo Sapiro. Online dictionary learning for sparse coding. In *Proceedings of Annual International Conference on Machine Learning*, pp. 689–696, 2009a.
- Julien Mairal, Jean Ponce, Guillermo Sapiro, Andrew Zisserman, and Francis Bach. Supervised dictionary learning. In *Proceedings of Advances in Neural Information Processing Systems*, volume 21, pp. 1–8, 2009b.
- Benoît Malézieux, Thomas Moreau, and Matthieu Kowalski. Understanding approximate and unrolled dictionary learning for pattern recovery. In *Proceedings of International Conference on Learning Representations*, 2022.
- D. Martin, C. Fowlkes, D. Tal, and J. Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proceedings of IEEE International Conference on Computer Vision*, volume 2, pp. 416–423, 2001.
- Vishal Monga, Yuelong Li, and Yonina C Eldar. Algorithm unrolling: Interpretable, efficient deep learning for signal and image processing. *arXiv:1912.10557*, pp. 1–27, 2019.
- Thomas Moreau and Joan Bruna. Understanding trainable sparse coding via matrix factorization. In *Proceedings of 5th International Conference on Learning Representations*, pp. 1–13, 2017.
- Thanh V Nguyen, Raymond KW Wong, and Chinmay Hegde. On the dynamics of gradient descent for autoencoders. In *Proceedings of International Conference on Artificial Intelligence and Statistics*, pp. 2858–2867. PMLR, 2019.
- Kenji Nose-Filho, Andre Kazuo Takahata, Renato Lopes, and Joao Marcos Travassos Romano. Improving sparse multichannel blind deconvolution with correlated seismic data: Foundations and further results. *IEEE Signal Processing Magazine*, 35(2):41–50, 2018.
- Bruno A Olshausen and David J Field. Sparse coding with an overcomplete basis set: A strategy employed by v1? *Vision research*, 37(23):3311–3325, 1997.
- Neal Parikh and Stephen Boyd. Proximal algorithms. *Foundations and Trends in optimization*, 1(3):127–239, 2014.

- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- Sirisha Rambhatla, Xingguo Li, and Jarvis Haupt. Noodl: Provable online dictionary learning and sparse coding. In *Proceedings of International Conference on Learning Representations*, pp. 1–11, 2018.
- Akshay Rangamani, Anirbit Mukherjee, Amitabh Basu, Ashish Arora, Tejaswini Ganapathi, Sang Chin, and Trac D. Tran. Sparse coding and autoencoders. In *Proceedings of IEEE International Symposium on Information Theory (ISIT)*, pp. 36–40, 2018.
- Marc aurelio Ranzato, Christopher Poultney, Sumit Chopra, and Yann Cun. Efficient learning of sparse representations with an energy-based model. In *Advances in Neural Information Processing Systems*, volume 19. MIT Press, 2007.
- Marc aurelio Ranzato, Y-lan Boureau, and Yann Cun. Sparse feature learning for deep belief networks. In J. Platt, D. Koller, Y. Singer, and S. Roweis (eds.), *Proceedings of Advances in Neural Information Processing Systems*, volume 20, 2008.
- Saharon Rosset, Ji Zhu, and Trevor Hastie. Boosting as a regularized path to a maximum margin classifier. *The Journal of Machine Learning Research*, 5:941–973, 2004.
- Christian J. Schuler, Michael Hirsch, Stefan Harmeling, and Bernhard Schölkopf. Learning to deblur. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(7):1439–1451, 2016.
- Dror Simon and Michael Elad. Rethinking the csc model for natural images. In *Proceedings of Advances in Neural Information Processing Systems*, volume 32, pp. 1–11, 2019.
- Oren Solomon, Regev Cohen, Yi Zhang, Yi Yang, Qiong He, Jianwen Luo, Ruud J. G. van Sloun, and Yonina C. Eldar. Deep unfolded robust pca with application to clutter suppression in ultrasound. *IEEE Transactions on Medical Imaging*, 39(4):1051–1063, 2020.
- Daniel A. Spielman, Huan Wang, and John Wright. Exact recovery of sparsely-used dictionaries. In *Proceedings of Annual Conference on Learning Theory*, volume 23 of *PMRL*, pp. 37.1–37.18, 2012.
- Pablo Sprechmann, Alex Bronstein, and Guillermo Sapiro. Learning efficient structured sparse models. In *Proceedings of International Conference on Machine Learning*, pp. 219–226, 2012.
- Shaozhe Tao, Daniel Boley, and Shuzhong Zhang. Local linear convergence of ista and fista on the lasso problem. *SIAM Journal on Optimization*, 26(1):313–336, 2016.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):267–288, 1996. ISSN 00359246.
- Ryan J. Tibshirani. The lasso problem and uniqueness. *Electronic Journal of Statistics*, 7(none):1456–1490, 2013.
- Bahareh Tolooshams, Sourav Dey, and Demba Ba. Scalable convolutional dictionary learning with constrained recurrent sparse auto-encoders. In *2018 IEEE 28th International Workshop on Machine Learning for Signal Processing (MLSP)*, pp. 1–6, 2018.
- Bahareh Tolooshams, Andrew Song, Simona Temereanca, and Demba Ba. Convolutional dictionary learning based auto-encoders for natural exponential-family distributions. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 9493–9503. PMLR, 13–18 Jul 2020.
- Bahareh Tolooshams, Sourav Dey, and Demba Ba. Deep residual autoencoders for expectation maximization-inspired dictionary learning. *IEEE Transactions on Neural Networks and Learning Systems*, 32(6):2415–2429, 2021a.

- Bahareh Tolooshams, Satish Mulleti, Demba Ba, and Yonina C Eldar. Unfolding neural networks for compressive multichannel blind deconvolution. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 2890–2894. IEEE, 2021b.
- Martin J Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery using ℓ_1 -constrained quadratic programming (lasso). *IEEE transactions on information theory*, 55(5):2183–2202, 2009.
- Zhaowen Wang, Ding Liu, Jianchao Yang, Wei Han, and Thomas Huang. Deep networks for image super-resolution with sparse prior. In *Proceedings of IEEE International Conference on Computer Vision*, pp. 370–378, 2015.
- Brendt Wohlberg. Sporco: A python package for standard and convolutional sparse representations. In *Proceedings of the 15th Python in Science Conference, Austin, TX, USA*, pp. 1–8, 2017.
- Bo Xin, Yizhou Wang, Wen Gao, David Wipf, and Baoyuan Wang. Maximal sparsity with deep networks? In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 29, pp. 1–9, 2016.
- Jianchao Yang, John Wright, Thomas S. Huang, and Yi Ma. Image super-resolution via sparse representation. *IEEE Transactions on Image Processing*, 19(11):2861–2873, 2010.
- Chih-Kuan Yeh, Joon Sik Kim, Ian EH Yen, and Pradeep Ravikumar. Representer point selection for explaining deep neural networks. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pp. 9311–9321, 2018.
- Lufang Zhang, Yaohua Hu, Chong Li, and Jen-Chih Yao. A new linear convergence result for the iterative soft thresholding algorithm. *Optimization*, 66(7):1177–1189, 2017.

A Appendix - proofs

A.1 Notation

Bold-lower-case and upper-case letters refer to vectors \mathbf{d} and matrices \mathbf{D} . We use \mathbf{d}_j to denote the j^{th} element of the vector \mathbf{d} , and \mathbf{D}_j is the j^{th} column of the matrix \mathbf{D} . $\lambda > 0$ is the regularization (sparsity-enforcing) parameter. $\sigma_{\max}(\mathbf{D})$ is the maximum singular value of \mathbf{D} . When taking the derivatives or norms w.r.t the matrix \mathbf{D} , we assume that \mathbf{D} is vectorized. $\nabla_1 \mathcal{L}(\mathbf{z}, \mathbf{D})$ and $\nabla_2 \mathcal{L}(\mathbf{z}, \mathbf{D})$ are the first derivatives of the loss w.r.t \mathbf{z} and \mathbf{D} , respectively. $\nabla_{11}^2 \mathcal{L}(\mathbf{z}, \mathbf{D})$ is the second derivative of the loss w.r.t \mathbf{z} . $\nabla_{21}^2 \mathcal{L}(\mathbf{z}, \mathbf{D})$ is the derivative of $\nabla_1 \mathcal{L}(\mathbf{z}, \mathbf{D})$ w.r.t \mathbf{D} . The support of \mathbf{z} is $\text{supp}(\mathbf{z}) \triangleq \{j: z_j \neq 0\}$.

A.2 Basic definitions and Lemmas

We list four definitions used throughout the paper below.

Definition A.1 (μ -incoherence). \mathbf{D} is μ -incoherent, i.e., for every pair (i, j) of columns, $|\langle \mathbf{D}_i, \mathbf{D}_j \rangle| \leq \mu/\sqrt{m}$.

Definition A.2 $((\delta, \kappa)$ -closeness). Dictionary \mathbf{D} is δ -close to \mathbf{D}^* , i.e., there is a permutation π and sign flip operator u such that $\forall i \|\mathbf{D}_{\pi(i)} - \mathbf{D}_i^*\|_2 \leq \delta$. Additionally, $\|\mathbf{D} - \mathbf{D}^*\|_2 \leq \kappa \|\mathbf{D}^*\|_2$.

Definition A.3 (Lipschitz function). A function $f: \mathbb{R}^m \rightarrow \mathbb{R}^p$ is L -Lipschitz w.r.t a norm $\|\cdot\|$ if $\exists L > 0$ s.t. $\|f(a) - f(b)\| \leq L\|a - b\| \forall a, b \in \mathbb{R}^m$.

Definition A.4 (Lipschitz differentiable function). A twice differentiable function $f: \mathbb{R}^m \rightarrow \mathbb{R}^p$ is L -Lipschitz differentiable w.r.t a norm $\|\cdot\|$ iff $\exists L > 0$ s.t. $\|\nabla^2 f(a)\| \leq L \forall a \in \mathbb{R}^m$.

Definition A.5 (Strong convexity). A twice differentiable function $f: \mathbb{R}^m \rightarrow \mathbb{R}^p$ is strongly convex if $\exists \mu > 0$ s.t. $\nabla^2 f(a) \succeq \mu \mathbf{I}$.

Definition A.6 (Norm of subgradient). For norms involving subgradients, we define $\|\partial h(\mathbf{z})\| := \max_{\mathbf{v} \in \partial h(\mathbf{z})} \|\mathbf{v}\|$.

In the proof of the theorems, we use the strong convexity of the reconstruction loss after support selection and the bounded property of the Lipschitz mapping stated below.

Lemma A.1 (Strong convexity of reconstruction loss). Given the support selection ([Proposition 4.1](#)), $\mathbf{D}_S^T \mathbf{D}_S$ is full-rank. Thus, $\forall t > B$, $\mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D}) = \mathcal{L}_{\mathbf{x}}(\mathbf{z}_{t,S}, \mathbf{D}_S)$ is strongly convex ([Definition A.5](#)) in \mathbf{z} .

Lemma A.2 (Lipschitz mapping). Given the recursion $\mathbf{z}_{t+1} = \Phi(\mathbf{z}_t) = \mathcal{P}_{\alpha h}(\mathbf{z}_t - \alpha \nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D}))$, from [Lemma A.1](#), there exist $B > 0$ such that loss $\mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D})$ is μ -strongly convex $\forall t > B$. Hence, using [Lemma 3.4](#),

$$\|\nabla_1 \Phi(\mathbf{z}_t, \mathbf{D})\|_2 = \|(I - \alpha \nabla_{11}^2 \mathcal{L}(\mathbf{z}_t, \mathbf{D})) \partial \mathcal{P}_{\alpha h}(\mathbf{z}_t)\|_2 \leq \rho \quad (14)$$

where $\rho \triangleq c_{\text{prox}}(1 - \alpha\mu) < 1$.

One key term, used in the proofs, is that $\mathbf{0} \in \nabla_1 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D}) + \partial h(\hat{\mathbf{z}})$ which is followed by the lasso optimality, i.e.,

Lemma A.3 (Lasso optimality). Lasso Karush-Kuhn-Tucker (KKT) optimality conditions are

$$\hat{\mathbf{z}} \in \arg \min_{\mathbf{z} \in \mathbb{R}^p} f_{\mathbf{x}}(\mathbf{z}, \mathbf{D}) \Leftrightarrow \mathbf{D}^T(\mathbf{x} - \mathbf{D}\hat{\mathbf{z}}) \in \lambda \partial \|\hat{\mathbf{z}}\|_1, \text{ and } |\hat{z}_j| = \begin{cases} \{\text{sign}(\hat{z}_j)\} & \text{if } \hat{z}_j \neq 0 \\ [-1, 1] & \text{if } \hat{z}_j = 0 \end{cases}, \forall j \in \{1, 2, \dots, p\}. \quad (15)$$

A.3 Forward pass proof details

Given the μ -incoherence of \mathbf{D}^* , and current dictionary closeness of δ_l , we re-state [Lemma 3.1](#) and proof it below. It shows that the current dictionary is μ_l -close to \mathbf{D}^* .

Lemma 3.1 (μ_l -incoherent). $\mathbf{D}^{(l)}$ is μ_l -incoherent where $\mu_l = \mu/\sqrt{m} + 2\delta_l$.

Proof.

$$\begin{aligned} \langle \mathbf{D}_i^{(l)}, \mathbf{D}_j^{(l)} \rangle &= \langle \mathbf{D}_i^*, \mathbf{D}_j^* \rangle - \langle \mathbf{D}_i^* - \mathbf{D}_i^{(l)}, \mathbf{D}_j^* \rangle - \langle \mathbf{D}_i^{(l)}, \mathbf{D}_j^* - \mathbf{D}_j^{(l)} \rangle \\ |\langle \mathbf{D}_i^{(l)}, \mathbf{D}_j^{(l)} \rangle| &\leq \mu/\sqrt{m} + \|\mathbf{D}_i^* - \mathbf{D}_i^{(l)}\|_2 \|\mathbf{D}_j^*\|_2 + \|\mathbf{D}_i^{(l)}\|_2 \|\mathbf{D}_j^* - \mathbf{D}_j^{(l)}\|_2 \leq \mu/\sqrt{m} + 2\delta_l \end{aligned} \quad (16)$$

■

We re-state and proof the forward pass support recovery (Theorem 4.1). This shows that given proper initialization and under mild conditions, the support of the true code \mathbf{z}^* is recovered with high probability in one iteration of the encoder.

Theorem 4.1 (Forward pass support recovery). *Given the Assumption 3.3, Assumption 3.4, suppose $\mathbf{D}^{(l)}$ is $\delta_l = \mathcal{O}^*(1/\sqrt{\log p})$ close to \mathbf{D}^* . If $s = \mathcal{O}^*(\sqrt{m}/\mu \log m)$, and $\mu = \mathcal{O}(\log m)$, then with high probability of at least $1 - \epsilon_{\text{supp-rec}}^{(l)}$, the choice of $\lambda_0 = C_{\min}/4$ recovers the support of the code \mathbf{z}^* in one encoder iteration, i.e., $\text{sign}(\text{ReLU}(\alpha(\mathbf{D}^{(l)\top} \mathbf{x} - \lambda_0))) = \text{sign}(\mathbf{z}^*)$, where $\epsilon_{\text{supp-rec}}^{(l)} = 2p \exp(-\frac{C_{\min}^2}{\mathcal{O}^*(\delta_l^2)})$.*

Proof. The code estimate after one iteration is $\mathbf{z}_1 = \mathcal{P}_{\alpha h}(\alpha \mathbf{D}^{(l)\top} \mathbf{x}) = \text{sign}(\mathbf{D}^{(l)\top} \mathbf{x}) \text{ReLU}(\alpha(|\mathbf{D}^{(l)\top} \mathbf{D}^* \mathbf{z}^*| - \lambda_0))$. We focus on the positive entries. The analysis for negative entries is similar. Writting the relation for i -th entry,

$$\mathbf{z}_1^i = \text{sign}(\mathbf{D}^{(l)\top} \mathbf{x}) \text{ReLU}(\alpha(\sum_{j \in S^*} \langle \mathbf{D}_i^{(l)}, \mathbf{D}_j^* \rangle \mathbf{z}_j^* - \lambda_0)) = \text{ReLU}(\alpha(\langle \mathbf{D}_i^{(l)}, \mathbf{D}_i^* \rangle \mathbf{z}_i^* + \sum_{j \in S^* \setminus \{i\}} \langle \mathbf{D}_i^{(l)}, \mathbf{D}_j^* \rangle \mathbf{z}_j^* - \lambda_0)) \quad (17)$$

We focus on the term inside ReLU and discard α , shared by all terms. We shows that under proper choice of λ_0 , $\langle \mathbf{D}_i^{(l)}, \mathbf{D}_i^* \rangle \mathbf{z}_i^*$ is greater than λ_0 and $\mathbf{v}_i = \sum_{j \in S^* \setminus \{i\}} \langle \mathbf{D}_i^{(l)}, \mathbf{D}_j^* \rangle \mathbf{z}_j^*$ is small with respect to λ_0 , hence getting cancelled by ReLU. The small value of \mathbf{v}_i , compared to $\langle \mathbf{D}_i^{(l)}, \mathbf{D}_i^* \rangle \mathbf{z}_i^*$, results in $\text{sign}(\mathbf{D}^{(l)\top} \mathbf{x})$ be equal to the $\text{sign}(\mathbf{D}^{(l)\top} \mathbf{D}^* \mathbf{z}^*)$ which is equal to the sign of \mathbf{z}^* .

Given the current dictionary distance $\|\mathbf{D}_i^{(l)} - \mathbf{D}_i^*\|_2 \leq \delta_l$, we can find a lower bound on $\langle \mathbf{D}_i^{(l)}, \mathbf{D}_i^* \rangle \mathbf{z}_i^*$ as follows

$$\begin{aligned} \langle \mathbf{D}_i^{(l)}, \mathbf{D}_i^* \rangle &= \frac{1}{2}(\|\mathbf{D}_i^*\|_2^2 + \|\mathbf{D}_i^{(l)}\|_2^2 - \|\mathbf{D}_i^{(l)} - \mathbf{D}_i^*\|_2^2) = 1 - \frac{1}{2}\|\mathbf{D}_i^{(l)} - \mathbf{D}_i^*\|_2^2 \\ |\langle \mathbf{D}_i^{(l)}, \mathbf{D}_i^* \rangle| &\geq 1 - \delta_l^2/2 \end{aligned} \quad (18)$$

Hence, for $i \in S^*$

$$|\langle \mathbf{D}_i^{(l)}, \mathbf{D}_i^* \rangle \mathbf{z}_i^*| \geq (1 - \delta_l^2/2) C_{\min} \quad (19)$$

otherwise, it is 0. Given, $\text{var}(\mathbf{z}_i^*) = 1$ for $i \in S^*$, we find an upper bound on the variance \mathbf{v}_i of as follows

$$\begin{aligned} \text{var}(\mathbf{v}_i) &= \sum_{j \in S^* \setminus \{i\}} \langle \mathbf{D}_i^{(l)}, \mathbf{D}_j^* \rangle^2 = \sum_{j \in S^* \setminus \{i\}} (\langle \mathbf{D}_i^*, \mathbf{D}_j^* \rangle + \langle \mathbf{D}_i^{(l)} - \mathbf{D}_i^*, \mathbf{D}_j^* \rangle)^2 \leq \sum_{j \in S^* \setminus \{i\}} 2(\langle \mathbf{D}_i^*, \mathbf{D}_j^* \rangle^2 + \langle \mathbf{D}_i^{(l)} - \mathbf{D}_i^*, \mathbf{D}_j^* \rangle^2) \\ &\leq \sum_{j \in S^* \setminus \{i\}} (2\mu^2/m) + 2\|(\mathbf{D}_i^{(l)} - \mathbf{D}_i^*)^\top \mathbf{D}_{S^* \setminus \{i\}}^*\|_2^2 \leq (2s\mu^2/m) + 2\|(\mathbf{D}_i^{(l)} - \mathbf{D}_i^*)\|_2^2 \|\mathbf{D}_{S^* \setminus \{i\}}^*\|_2^2 \\ &\leq 2(s\mu^2/m + 4\delta_l^2) = \mathcal{O}^*(\delta_l^2) \end{aligned} \quad (20)$$

where we used the Gershgorin Circle Theorem for the bound $\|\mathbf{D}_{S^* \setminus \{i\}}^*\|_2 \leq 2$. With the sub-Gaussian assumption on the coefficients \mathbf{z}^* , we get the following using Chernoff bound concerning \mathbf{v}_i .

$$P(|\mathbf{v}_i| \geq \frac{C_{\min}}{4}) \leq 2 \exp(-\frac{C_{\min}^2}{4s\mu^2/m + 16\delta_l^2}) = 2 \exp(-\frac{C_{\min}^2}{\mathcal{O}^*(\delta_l^2)}) \quad (21)$$

Taking a union bound over all indices $i \in [1, p]$ will result in

$$P(\max_i |\mathbf{v}_i| \geq \frac{C_{\min}}{4}) \leq 2p \exp(-\frac{C_{\min}^2}{\mathcal{O}^*(\delta_l^2)}) := \epsilon_{\text{supp-rec}}^{(l)} \quad (22)$$

Hence, we can set $\lambda_0 = C_{\min}/2$. ■

We re-state and prove the forward pass support preservation (Theorem 4.2).

Theorem 4.2 (Forward pass support preservation). *Given the [Assumption 3.3](#), [Assumption 3.4](#), suppose $\mathbf{D}^{(l)}$ is $\delta_l = \mathcal{O}^*(1/\log p)$ close to \mathbf{D}^* . If $s = \mathcal{O}^*(\sqrt{m}/\mu \log m)$, $\mu = \mathcal{O}(\log m)$, and the regularizer and step size are chosen such that $\lambda_t^{(l)} = \frac{\mu_l}{\sqrt{m}} \|\mathbf{z}^* - \mathbf{z}_t\|_1 + a_\gamma = \Omega(\frac{s \log m}{\sqrt{m}})$ and $\alpha^{(l)} \leq 1 - \frac{2\lambda_t - (1 - \frac{\delta_l^2}{2})C_{\min}}{\lambda_{t-1}}$, then with high probability of at least $1 - \epsilon_{\text{supp-pres}}^{(l)}$, the support, recovered at the first iteration, is preserved through the encoder iterations. We have $a_\gamma = \mathcal{O}(\sqrt{s\delta_l})$ and $\epsilon_{\text{supp-pres}}^{(l)} := \epsilon_{\text{supp-rec}}^{(l)} + \epsilon_\gamma^{(l)} = 2p \exp(-\frac{C_{\min}^2}{\mathcal{O}^*(\delta_l^2)}) + 2s \exp(-\frac{1}{\mathcal{O}(\delta_l)})$.*

Proof. Given current dictionary $\mathbf{D}^{(l)}$, in each iteration of the forward pass, we have $\mathbf{z}_{t+1} = \mathcal{P}_{\alpha h}(\mathbf{z}_t + \alpha \mathbf{D}^T(\mathbf{D}^* \mathbf{z}^* - \mathbf{D} \mathbf{z}_t))$. We focus on the entries that are non-negative. Then procedure for negative code entries is similar. We follow similar steps as in ([Rambhatla et al., 2018](#)). We get

$$\begin{aligned}
\mathbf{z}_{t+1}^j &= \text{ReLU}((\mathbf{I} - \alpha \mathbf{D}^{(l)T} \mathbf{D}^{(l)})_{(j,:)} \mathbf{z}_t + \alpha (\mathbf{D}^{(l)T} \mathbf{D}^*)_{(j,:)} \mathbf{z}^* - \alpha \lambda_t^j) \\
&= \text{ReLU}((\mathbf{I} - \alpha \mathbf{D}^{(l)T} \mathbf{D}^{(l)})_{(j,:)} \mathbf{z}_t + \alpha ((\mathbf{D}^{(l)} - \mathbf{D}^*)^T \mathbf{D}^*)_{(j,:)} \mathbf{z}^* + \alpha (\mathbf{D}^{*T} \mathbf{D}^*)_{(j,:)} \mathbf{z}^* - \alpha \lambda_t^j) \\
&= \text{ReLU}((1 - \alpha) \mathbf{z}_t^j - \alpha \sum_{i \neq j} \langle \mathbf{D}_j^{(l)}, \mathbf{D}_i^{(l)} \rangle \mathbf{z}_t^i + \alpha \langle (\mathbf{D}_j^{(l)} - \mathbf{D}_j^*), \mathbf{D}_j^* \rangle \mathbf{z}^{*j} \\
&\quad + \alpha \sum_{i \neq j} \langle \mathbf{D}_j^{(l)} - \mathbf{D}_j^*, \mathbf{D}_i^* \rangle \mathbf{z}^{*i} + \alpha \mathbf{z}^{*j} + \alpha \sum_{i \neq j} \langle \mathbf{D}_j^*, \mathbf{D}_i^* \rangle \mathbf{z}^{*i} - \alpha \lambda_t^j) \\
&= \text{ReLU}((1 - \alpha) \mathbf{z}_t^j + \alpha (1 - \beta_j^{(l)}) \mathbf{z}^{*j} + \alpha \eta_t^j - \alpha \lambda_t^j)
\end{aligned} \tag{23}$$

where $\beta_j^{(l)} = \langle \mathbf{D}_j^* - \mathbf{D}_j^{(l)}, \mathbf{D}_j^* \rangle$, and $\eta_{t,j}^{(l)} = -\sum_{i \neq j} \langle \mathbf{D}_j^{(l)}, \mathbf{D}_i^{(l)} \rangle \mathbf{z}_t^i + (\langle \mathbf{D}_j^{(l)} - \mathbf{D}_j^*, \mathbf{D}_i^* \rangle + \langle \mathbf{D}_j^*, \mathbf{D}_i^* \rangle) \mathbf{z}^{*i}$. With $\|\mathbf{D}_j^{(l)} - \mathbf{D}_j^*\|_2 \leq \delta_l$, $\beta_j^{(l)}$ can be bounded as follows

$$\beta_j^{(l)} = \langle \mathbf{D}_j^* - \mathbf{D}_j^{(l)}, \mathbf{D}_j^* \rangle \leq \delta_l^2/2 \tag{24}$$

where we used the relation $\|\mathbf{D}_j^{(l)} - \mathbf{D}_j^*\|_2^2 = 2(1 - \langle \mathbf{D}_j^{(l)}, \mathbf{D}_j^* \rangle)$. We re-write $\eta_{t,j}^{(l)}$ below

$$\begin{aligned}
\eta_{t,j}^{(l)} &= -\sum_{i \neq j} \langle \mathbf{D}_j^{(l)}, \mathbf{D}_i^{(l)} \rangle \mathbf{z}_t^i + (\langle \mathbf{D}_j^{(l)} - \mathbf{D}_j^*, \mathbf{D}_i^* \rangle + \langle \mathbf{D}_j^*, \mathbf{D}_i^* \rangle) \mathbf{z}^{*i} \\
&= -\sum_{i \neq j} \langle \mathbf{D}_j^{(l)}, \mathbf{D}_i^{(l)} \rangle \mathbf{z}_t^i + \sum_{i \neq j} (\langle \mathbf{D}_j^{(l)} - \mathbf{D}_j^*, \mathbf{D}_i^* \rangle + \langle \mathbf{D}_j^*, \mathbf{D}_i^* \rangle) \mathbf{z}^{*i} + \sum_{i \neq j} \langle \mathbf{D}_j^{(l)}, \mathbf{D}_i^{(l)} \rangle \mathbf{z}^{*i} - \sum_{i \neq j} \langle \mathbf{D}_j^{(l)}, \mathbf{D}_i^{(l)} \rangle \mathbf{z}^{*i} \\
&= \sum_{i \neq j} \langle \mathbf{D}_j^{(l)}, \mathbf{D}_i^{(l)} \rangle (\mathbf{z}^{*i} - \mathbf{z}_t^i) + \sum_{i \neq j} (\langle \mathbf{D}_j^{(l)} - \mathbf{D}_j^*, \mathbf{D}_i^* \rangle + \langle \mathbf{D}_j^*, \mathbf{D}_i^* \rangle - \langle \mathbf{D}_j^{(l)}, \mathbf{D}_i^{(l)} \rangle) \mathbf{z}^{*i} \\
&= \sum_{i \neq j} \langle \mathbf{D}_j^{(l)}, \mathbf{D}_i^{(l)} \rangle (\mathbf{z}^{*i} - \mathbf{z}_t^i) + \sum_{i \neq j} \langle \mathbf{D}_j^{(l)}, \mathbf{D}_i^* - \mathbf{D}_i^{(l)} \rangle \mathbf{z}^{*i} \\
&= \sum_{i \neq j} \langle \mathbf{D}_j^{(l)}, \mathbf{D}_i^{(l)} \rangle (\mathbf{z}^{*i} - \mathbf{z}_t^i) + \gamma_j^{(l)}
\end{aligned} \tag{25}$$

where $\gamma_j^{(l)} = \sum_{i \neq j} \langle \mathbf{D}_j^{(l)}, \mathbf{D}_i^* - \mathbf{D}_i^{(l)} \rangle \mathbf{z}^{*i}$. Given the sub-Gaussian entries of the code \mathbf{z}^* , we provide a bound on the variance of $\gamma_j^{(l)}$ below:

$$\text{var}(\gamma_j^{(l)}) = \sum_{i \neq j} \langle \mathbf{D}_j^{(l)}, \mathbf{D}_i^* - \mathbf{D}_i^{(l)} \rangle^2 \leq s \delta_l^2 \tag{26}$$

Now, using Chernoff bound on the sub-Gaussian code entries, we get

$$P(|\gamma_j^{(l)}| > a) \leq 2 \exp(-\frac{a^2}{2s\delta_l^2}) \tag{27}$$

To bound all the terms in the support, for $j \in S^*$, we have

$$P(\max |\gamma_j^{(l)}| > a_\gamma) \leq \epsilon_\gamma^{(l)} \tag{28}$$

where $\epsilon_\gamma^{(l)} = 2s \exp(\frac{-a_\gamma^2}{2s\delta_l^2})$. Let $a_\gamma = \mathcal{O}(\sqrt{s\delta_l})$, then $\epsilon_\gamma^{(l)} = 2s \exp(\frac{-1}{\mathcal{O}(\delta_l)})$. The above analysis states that with high probability of at least $1 - \epsilon_\gamma^{(l)}$, $|\gamma_j^{(l)}| \leq a_\gamma = \mathcal{O}(\sqrt{s\delta_l})$. Next, we write the recursion for when the support is identified (see [Theorem 4.1](#)). For the code at iteration T , we have

$$\begin{aligned} \mathbf{z}_T^j &= (1 - \alpha)^T \mathbf{z}_0^j + \mathbf{z}^{*j} \sum_{t=1}^T \alpha(1 - \beta_j^{(l)})(1 - \alpha)^{T-t} + \sum_{t=1}^T \alpha(\eta_{t,j}^{(l)} - \lambda_t^j)(1 - \alpha)^{T-t} \\ &= (1 - \alpha)^T \mathbf{z}_0^j + \mathbf{z}^{*j}(1 - \beta_j^{(l)})(1 - (1 - \alpha)^T) + \sum_{t=1}^T \alpha(\eta_{t,j}^{(l)} - \lambda_t^j)(1 - \alpha)^{T-t} \\ &= \mathbf{z}^{*j}(1 - \beta_j^{(l)}) + (1 - \alpha)^T(\mathbf{z}_0^j - \mathbf{z}^{*j}(1 - \beta_j^{(l)})) + \sum_{t=1}^T \alpha(\eta_{t,j}^{(l)} - \lambda_t^j)(1 - \alpha)^{T-t} \\ &= \mathbf{z}^{*j}(1 - \beta_j^{(l)}) + \zeta^j \end{aligned} \quad (29)$$

where $\zeta_{T,j}^{(l)} = (1 - \alpha)^T(\mathbf{z}_0^j - \mathbf{z}^{*j}(1 - \beta_j^{(l)})) + \sum_{t=1}^T \alpha(\eta_{t,j}^{(l)} - \lambda_t^j)(1 - \alpha)^{T-t}$. With the support correctly identified at iteration $t - 1$, we show that the support is preserved at iteration t . With $\|\mathbf{z}^{*i} - \mathbf{z}_t^i\|_1 = \mathcal{O}(s)$, for each $j \in S^*$, we have

$$\eta_{t,j}^{(l)} = \sum_{i \neq j} \langle \mathbf{D}_j^{(l)}, \mathbf{D}_i^{(l)} \rangle (\mathbf{z}^{*i} - \mathbf{z}_t^i) + \gamma_j^{(l)} \leq \frac{\mu_t}{\sqrt{m}} \|\mathbf{z}^* - \mathbf{z}_t\|_1 + a_\gamma = \mathcal{O}(\frac{s \log m}{\sqrt{m}}) \quad (30)$$

We make sure the regularizer is chosen such that

$$\lambda_t \geq \frac{\mu_t}{\sqrt{m}} \|\mathbf{z}^* - \mathbf{z}_t\|_1 + a_\gamma \quad (31)$$

We see that the larger the code error and coherence between the columns of the current dictionary, the larger λ_t should be. This is to suppress the noise component in the code recursion and make sure no false support is introduced. Furthermore, we want λ_t to be lower than half of the signal component, i.e.,

$$\begin{aligned} \alpha \lambda_t &\leq \frac{1 - \alpha}{2} \mathbf{z}_t^j + \frac{\alpha}{2} (1 - \beta_j^{(l)}) \mathbf{z}^{*j}, \forall j \in S^* \\ \alpha \lambda_t &\leq \frac{1 - \alpha}{2} \mathbf{z}_t^{\min} + \frac{\alpha}{2} (1 - \frac{\delta_l^2}{2}) C_{\min} \end{aligned} \quad (32)$$

where $\mathbf{z}_t^{\min} = \min_j \mathbf{z}_t^j$. We further shrink the upper bound, given the code from previous iteration $t - 1$ (i.e., $\alpha \lambda_{t-1} \leq \mathbf{z}_t^{\min}$). Hence, we want the regularizer to follow

$$\begin{aligned} \alpha \lambda_t &\leq \frac{1 - \alpha}{2} \alpha \lambda_{t-1} + \frac{\alpha}{2} (1 - \frac{\delta_l^2}{2}) C_{\min} \\ \lambda_t &\leq \frac{1 - \alpha}{2} \lambda_{t-1} + \frac{1}{2} (1 - \frac{\delta_l^2}{2}) C_{\min} \end{aligned} \quad (33)$$

This condition is to make sure the identified supports are not killed in the recursion. We use the condition to set the step size α . We get

$$\alpha \leq 1 - \frac{2\lambda_t - (1 - \frac{\delta_l^2}{2}) C_{\min}}{\lambda_{t-1}} \quad (34)$$

Hence, $\lambda_t = \Omega(\frac{s \log m}{\sqrt{m}})$ and α should be chosen sufficiently small such that the condition above is met. We denote $\epsilon_{\text{supp-pres}}^{(l)} := \epsilon_{\text{supp-rec}}^{(l)} + \epsilon_\gamma^{(l)} = 2p \exp(\frac{-C_{\min}^2}{\mathcal{O}^*(\delta_l^2)}) + 2s \exp(\frac{-1}{\mathcal{O}(\delta_l)})$. Hence, with high probability of at least $1 - \epsilon_{\text{supp-pres}}^{(l)}$, the support, recovered at the first iteration, is preserved through the encoder iterations. ■

With support recovery at first iteration and its preservation, we now re-state the forward pass code ([Theorem 4.3](#)) and Jacobian ([Theorem 4.5](#)) convergences. [Theorem 4.1](#) and [Theorem 4.2](#) allow to achieve linear convergence in the forward pass right after the first encoder iteration, i.e., $B = 1$.

Theorem 4.3 (Local forward pass code convergence). *Given the encoder $\mathbf{z}_{t+1} = \Phi(\mathbf{z}_t, \mathbf{D})$, [Assumption 3.6](#), [Lemmas 3.2](#), [A.1](#) and [A.2](#), then $\exists \rho < 1, B > 0$ s.t. $\|\mathbf{z}_t - \hat{\mathbf{z}}\|_2 \leq \mathcal{O}(\rho^t) \forall t > B$, where $\hat{\mathbf{z}}$ is the unique minimizer of lasso (1). Furthermore, given [Theorem 4.1](#) and [Theorem 4.2](#), $B = 1$.*

Proof. Given the support selection at iteration B , from [Lemma A.1](#), we have $\nabla_{11}^2 \mathcal{L}(\mathbf{z}_t, \mathbf{D}) \succeq \mu \mathbf{I}$ restricted to the support for $t > B$. Then, from [Lemma A.2](#), we get

$$\|\nabla_1 \Phi(\mathbf{z}_t, \mathbf{D})\|_2 = \|(\mathbf{I} - \alpha \nabla_{11}^2 \mathcal{L}(\mathbf{z}_t, \mathbf{D})) \partial \mathcal{P}_{\alpha h}(\mathbf{z}_t)\|_2 \leq \rho$$

where $\rho \triangleq c_{\text{prox}}(1 - \alpha\mu) < 1$. Hence, using fixed-point property ([Lemma 3.2](#))

$$\exists B > 0, \text{ s.t. } \|\mathbf{z}_{t+1} - \hat{\mathbf{z}}\|_2 = \|\Phi(\mathbf{z}_t) - \Phi(\hat{\mathbf{z}})\|_2 \leq \rho \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2 \forall t > B$$

where $\hat{\mathbf{z}} = \arg \min_{\mathbf{z}} f_{\mathbf{x}}(\mathbf{z}, \mathbf{D})$. Unrolling the recursion,

$$\|\mathbf{z}_t - \hat{\mathbf{z}}\|_2 \leq \rho^{t-B} \|\mathbf{z}_B - \hat{\mathbf{z}}\|_2.$$

■

Theorem 4.5 (Local forward pass Jacobian convergence). *Given the recursion $\mathbf{z}_{t+1} = \Phi(\mathbf{z}_t, \mathbf{D})$, and $\hat{\mathbf{z}}$ the unique minimizer of lasso with Jacobian $\hat{\mathbf{J}}$, then $\exists \rho < 1, B > 0$ s.t. $\|\mathbf{J}_t - \hat{\mathbf{J}}\|_2 \leq \mathcal{O}(t\rho^t) \forall t > B$. Furthermore, given [Theorem 4.1](#) and [Theorem 4.2](#), $B = 1$.*

Proof. Differentiating the recursion,

$$\mathbf{J}_{t+1} = \nabla_1 \Phi(\mathbf{z}_t, \mathbf{D}) \mathbf{J}_t + \nabla_2 \Phi(\mathbf{z}_t, \mathbf{D}).$$

Similarly,

$$\hat{\mathbf{J}} = \nabla_1 \Phi(\hat{\mathbf{z}}, \mathbf{D}) \hat{\mathbf{J}} + \nabla_2 \Phi(\hat{\mathbf{z}}, \mathbf{D})$$

where $\hat{\mathbf{z}}$ is a minimizer of lasso and fixed-point of the mapping (see [Lemma 3.2](#)). Subtract the terms

$$\mathbf{J}_{t+1} - \hat{\mathbf{J}} = \nabla_1 \Phi(\mathbf{z}_t, \mathbf{D})(\mathbf{J}_t - \hat{\mathbf{J}}) + (\nabla_1 \Phi(\mathbf{z}_t, \mathbf{D}) - \nabla_1 \Phi(\hat{\mathbf{z}}, \mathbf{D}))\hat{\mathbf{J}} + (\nabla_2 \Phi(\mathbf{z}_t, \mathbf{D}) - \nabla_2 \Phi(\hat{\mathbf{z}}, \mathbf{D}))$$

Given the Lipschitz properties of \mathcal{L} and h , we can further get the upper bounds on $\|\nabla_1 \Phi(\mathbf{a}, \mathbf{D}) - \nabla_1 \Phi(\mathbf{b}, \mathbf{D})\|_2 \leq L_{\Phi_1}$ and $\|\nabla_2 \Phi(\mathbf{a}, \mathbf{D}) - \nabla_2 \Phi(\mathbf{b}, \mathbf{D})\|_2 \leq L_{\Phi_2}$. Hence, with upper bound on the norm of Jacobian ([Assumption 4.1](#)), there exists $B > 0$ such that $\forall t > B$

$$\begin{aligned} \|\mathbf{J}_{t+1} - \hat{\mathbf{J}}\|_2 &\leq \|\nabla_1 \Phi(\mathbf{z}_t, \mathbf{D})\|_2 \|\mathbf{J}_t - \hat{\mathbf{J}}\|_2 + \|\nabla_1 \Phi(\mathbf{z}_t, \mathbf{D}) - \nabla_1 \Phi(\hat{\mathbf{z}}, \mathbf{D})\|_2 \|\hat{\mathbf{J}}\|_2 \\ &\quad + \|\nabla_2 \Phi(\mathbf{z}_t, \mathbf{D}) - \nabla_2 \Phi(\hat{\mathbf{z}}, \mathbf{D})\|_2 \\ &\leq \rho \|\mathbf{J}_t - \hat{\mathbf{J}}\|_2 + c \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2 \end{aligned}$$

where $c \triangleq M_J L_{\Phi_1} + L_{\Phi_2}$. Hence,

$$\|\mathbf{J}_{t+1} - \hat{\mathbf{J}}\|_2 \leq \rho \|\mathbf{J}_t - \hat{\mathbf{J}}\|_2 + \mathcal{O}(\rho^t).$$

Unrolling the recursion,

$$\|\mathbf{J}_{t+1} - \hat{\mathbf{J}}\|_2 \leq \mathcal{O}((t+1)\rho^t).$$

■

Theorem 4.4 (Global forward pass code error). *Let $\hat{\mathbf{z}}$ be the fixed-point of the encoder with iterations $\mathbf{z}_{t+1} = \Phi(\mathbf{z}_t, \mathbf{D})$. Given [Assumption 3.6](#), [Lemmas 3.2](#), [A.1](#) and [A.2](#), we have $\|\hat{\mathbf{z}} - \mathbf{z}^*\|_2 \leq \mathcal{O}(\|\mathbf{D} - \mathbf{D}^*\|_2 + \hat{\delta}^*)$, where $\hat{\delta}^* = \|\hat{\mathbf{z}}^* - \mathbf{z}^*\|_2$, $\hat{\mathbf{z}}$ is the unique minimizer of lasso (1) given the dictionary \mathbf{D} , $\hat{\mathbf{z}}^*$ is the unique minimizer of lasso (1) given the dictionary \mathbf{D}^* , and \mathbf{z}^* is the ground-truth code.*

Proof. We first find the error between $\hat{\mathbf{z}}$ and $\hat{\mathbf{z}}^*$ which is the unique minimizer of lasso (1) given the true dictionary \mathbf{D}^* . Using fixed-point property (Lemma 3.2), we get

$$\|\hat{\mathbf{z}} - \hat{\mathbf{z}}^*\|_2 = \|\Phi(\hat{\mathbf{z}}, \mathbf{D}) - \Phi(\hat{\mathbf{z}}^*, \mathbf{D}^*)\|_2 \leq \|\Phi(\hat{\mathbf{z}}, \mathbf{D}) - \Phi(\hat{\mathbf{z}}^*, \mathbf{D})\|_2 + \|\Phi(\hat{\mathbf{z}}^*, \mathbf{D}) - \Phi(\hat{\mathbf{z}}^*, \mathbf{D}^*)\|_2 \quad (35)$$

Using the μ -strongly convexity of $\mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D})$ on the support, and L_{21} Lipschitz constants of $\nabla_{21}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D})$, we upper bound the term as follows:

$$\|\hat{\mathbf{z}} - \hat{\mathbf{z}}^*\|_2 \leq \rho \|\hat{\mathbf{z}} - \hat{\mathbf{z}}^*\|_2 + \alpha L_{21} c_{\text{prox}} \|\mathbf{D} - \mathbf{D}^*\|_2 \quad (36)$$

Where $\rho \triangleq c_{\text{prox}}(1 - \alpha\mu) < 1$. Denote $q \triangleq \frac{\alpha c_{\text{prox}} L_{21}}{1 - \rho}$ which can be made to be small with proper choice of step size α .

$$\|\hat{\mathbf{z}} - \hat{\mathbf{z}}^*\|_2 \leq q \|\mathbf{D} - \mathbf{D}^*\|_2 \quad (37)$$

Hence, we get the following code error

$$\|\hat{\mathbf{z}} - \mathbf{z}^*\|_2 \leq \|\hat{\mathbf{z}} - \hat{\mathbf{z}}^*\|_2 + \|\hat{\mathbf{z}}^* - \mathbf{z}^*\|_2 \leq q \|\mathbf{D} - \mathbf{D}^*\|_2 + \hat{\delta}^* \leq \mathcal{O}(\|\mathbf{D} - \mathbf{D}^*\|_2 + \hat{\delta}^*) \quad (38)$$

■

A.4 Local backward pass proof details

In each update of the dictionary, we bound the gradient approximations as function of unrolling t (Theorem 4.7). This shows that $\mathbf{g}_t^{\text{ae-lasso}}$ converges faster than $\mathbf{g}_t^{\text{dec}}$ and $\mathbf{g}_t^{\text{ae-ls}}$, and the latter is a biased estimator of $\hat{\mathbf{g}}$. This is followed by Theorem 4.7 showing the order magnitude of the bounds is indeed tight.

Theorem 4.7 (Local convergence of gradients). *Given the convergence results from the forward pass (Theorems 4.3 and 4.5), $\exists \rho < 1, B > 0$ such that $\forall t > B$, the errors of gradients defined in Algorithm 2 w.r.t $\hat{\mathbf{g}}$ (4) satisfy*

$$\begin{aligned} \|\mathbf{g}_t^{\text{dec}} - \hat{\mathbf{g}}\|_2 &\leq \mathcal{O}(\rho^t) \\ \|\mathbf{g}_t^{\text{ae-lasso}} - \hat{\mathbf{g}}\|_2 &\leq \mathcal{O}(t\rho^{2t}) \\ \|\mathbf{g}_t^{\text{ae-ls}} - \hat{\mathbf{g}}\|_2 &\leq \mathcal{O}(t\rho^{2t} + M_J \lambda \sqrt{s}). \end{aligned} \quad (8)$$

Moreover, the order of upper bounds is tight (see Lemma A.4).

Proof. For $\mathbf{g}_t^{\text{dec}}$, with the infinite fresh samples, we have $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \nabla_2 \mathcal{L}_{\mathbf{x}^i}(\mathbf{z}_t^i, \mathbf{D}) = \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [\nabla_2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D})]$ a.s. Based on Lemma 3.3, we get

$$\begin{aligned} \|\mathbf{g}_t^{\text{dec}} - \hat{\mathbf{g}}\|_2 &= \|\mathbb{E}_{\mathbf{x} \in \mathcal{X}} [\nabla_2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D})] - \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [\nabla_2 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D})]\|_2 \\ &\leq \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [\|\nabla_2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D}) - \nabla_2 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D})\|_2] \leq \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [L_2 \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2] \leq \mathcal{O}(\rho^t). \end{aligned} \quad (39)$$

Similarly, for $\mathbf{g}_t^{\text{ae-lasso}}$ and $\mathbf{g}_t^{\text{ae-ls}}$, we replace the sample mean for gradient computations with expectation in their limit. We re-write the gradient estimation error as following

$$\begin{aligned} \mathbf{g}_t^{\text{ae-lasso}} - \hat{\mathbf{g}} &= \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [Q(\hat{\mathbf{z}}, \mathbf{J}_t)(\mathbf{z}_t - \hat{\mathbf{z}})] + \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [Q_t^{21}(\hat{\mathbf{z}})] + \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [\mathbf{J}_t Q_t^{\text{lasso-11}}(\hat{\mathbf{z}})] \\ \mathbf{g}_t^{\text{ae-ls}} - \hat{\mathbf{g}} &= \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [Q(\hat{\mathbf{z}}, \mathbf{J}_t)(\mathbf{z}_t - \hat{\mathbf{z}})] + \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [Q_t^{21}(\hat{\mathbf{z}})] + \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [\mathbf{J}_t Q_t^{\text{ls-11}}(\hat{\mathbf{z}})] \end{aligned} \quad (40)$$

where

$$\begin{aligned} Q_t^{21}(\mathbf{z}) &\triangleq \nabla_2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D}) - \nabla_2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D}) - \nabla_{21}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D})(\mathbf{z}_t - \mathbf{z}) \\ Q_t^{\text{lasso-11}}(\mathbf{z}) &\triangleq \nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D}) + \partial h(\mathbf{z}_t) - \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D})(\mathbf{z}_t - \mathbf{z}) \\ Q_t^{\text{ls-11}}(\mathbf{z}) &\triangleq \nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D}) - \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D})(\mathbf{z}_t - \mathbf{z}) \\ Q(\mathbf{z}, \mathbf{J}) &\triangleq \mathbf{J} \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D}) + \nabla_{21}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D}). \end{aligned} \quad (41)$$

We provide bounds on the above in Lemma A.5. Hence, it suffices to bound the terms on the r.h.s as follows:

$$\|\mathbf{g}_t^{\text{ae-lasso}} - \hat{\mathbf{g}}\|_2 \leq \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [L_1 \|\mathbf{J}_t - \hat{\mathbf{J}}\|_2 \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2 + (L_{21}/2) \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2^2 + M_J (L_{11}/2) \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2^2]. \quad (42)$$

Using the convergence errors from the forward pass (Theorems 4.3 and 4.5),

$$\|\mathbf{g}_t^{\text{ae-lasso}} - \hat{\mathbf{g}}\|_2 \leq L_1 \mathcal{O}(t\rho^{2t}) + (L_{21}/2 + M_J(L_{11}/2)) \mathcal{O}(\rho^{2t}) = \mathcal{O}(t\rho^{2t}). \quad (43)$$

Similarly,

$$\|\mathbf{g}_t^{\text{ae-ls}} - \hat{\mathbf{g}}\|_2 \leq \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [L_1 \|\mathbf{J}_t - \hat{\mathbf{J}}\|_2 \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2 + (L_{21}/2) \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2^2 + M_J((L_{11}/2) \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2^2 + \|\partial h(\hat{\mathbf{z}})\|_2)]. \quad (44)$$

Using the convergence errors from the forward pass (Theorems 4.3 and 4.5),

$$\|\mathbf{g}_t^{\text{ae-ls}} - \hat{\mathbf{g}}\|_2 \leq L_1 \mathcal{O}(t\rho^{2t}) + (L_{21}/2 + M_J L_{11}/2) \mathcal{O}(\rho^{2t}) + M_J \|\partial h(\hat{\mathbf{z}})\|_2 = \mathcal{O}(t\rho^{2t} + M_J \lambda \sqrt{s}). \quad (45)$$

■

Lemma A.4 (Tight local bound). *The order magnitude of the upper bounds in Theorem 4.7 is tight.*

Proof. It is sufficient to show that there exist an example such that its forward pass code and Jacobian convergences are $\mathcal{O}(\rho^t)$ and $\mathcal{O}(t\rho^t)$, respectively. The following example confirms this. Without loss of generality, let \mathbf{z}^* be 1-sparse and non-negative, $\mathbf{D} = \mathbf{D}^*$ and $\mathbf{D}_j = \mathbf{0}$ for $j \neq i$. The loss function is $\frac{1}{2} \|\mathbf{D}_i^* \mathbf{z}^{*i} - \mathbf{D}_i \mathbf{z}^i\|_2^2 + \lambda |\mathbf{z}^i|$. Given the support recovery after first iteration, the encoder forward pass implements $\mathbf{z}_{t+1}^i = \mathbf{z}_t^i - \alpha(\mathbf{D}_i^T(\mathbf{D}_i \mathbf{z}_t^i - \mathbf{D}_i^* \mathbf{z}^{*i}) + \lambda) = (1 - \alpha)\mathbf{z}_t^i + \alpha(\mathbf{z}^{*i} - \lambda)$. Hence, the forward pass convergences are

$$\begin{aligned} \mathbf{z}_t^i &= (1 - \alpha)^t \mathbf{z}_0 + \sum_{k=1}^t \alpha(1 - \alpha)^{t-k} (\mathbf{z}^{*i} - \lambda) = (1 - \alpha)^t \mathbf{z}_0 + (1 - (1 - \alpha)^t) (\mathbf{z}^{*i} - \lambda) \\ \mathbf{z}_t^i - \hat{\mathbf{z}}^i &= \rho^t (\mathbf{z}_0 - \mathbf{z}^{*i} + \lambda) = \mathcal{O}(\rho^t) \end{aligned} \quad (46)$$

and

$$\begin{aligned} \mathbf{J}_t^i &= \mathbf{J}_{t-1}^i - \alpha(\mathbf{J}_{t-1}^i + 2\mathbf{D}_i \mathbf{z}_t^i - \mathbf{D}_i^* \mathbf{z}^{*i}) = \rho \mathbf{J}_{t-1}^i + \mathcal{O}(\rho^t) + \hat{\mathbf{J}}^i \\ \mathbf{J}_t^i - \hat{\mathbf{J}}^i &= \rho^t \mathbf{J}_0^i + \sum_{k=1}^t \mathcal{O}(\rho^k) = \mathcal{O}(t\rho^t) \end{aligned} \quad (47)$$

where $\rho = 1 - \alpha$, $\hat{\mathbf{z}}^i = \mathbf{z}^{*i} - \lambda$, and $\hat{\mathbf{J}}^i = \alpha(2\mathbf{D}_i \hat{\mathbf{z}} - \mathbf{D}_i^* \mathbf{z}^{*i})$

■

Lemma A.5 (Local bounds). *From local gradient errors in Theorem 4.7, the following are satisfied*

$$\begin{aligned} \|Q_t^{21}(\hat{\mathbf{z}})\|_2 &\leq (L_{21}/2) \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2^2, & \|Q_t^{\text{lasso-11}}(\hat{\mathbf{z}})\|_2 &\leq (L_{11}/2) \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2^2 \\ \|Q(\hat{\mathbf{z}}, \mathbf{J}_t)\|_2 &\leq L_1 \|\mathbf{J}_t - \hat{\mathbf{J}}\|_2, & \|Q_t^{\text{ls-11}}(\hat{\mathbf{z}})\|_2 &\leq (L_{11}/2) \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2^2 + \|\partial h(\hat{\mathbf{z}})\|_2. \end{aligned} \quad (48)$$

Proof. For $Q_t^{21}(\hat{\mathbf{z}})$, given convexity of $\nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D})$ and its domain (Assumption 3.1) and Lemma 3.3, we achieve the quadratic upper bound. For $Q_t^{\text{lasso-11}}(\hat{\mathbf{z}})$, we add and subtract $\nabla_1 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D})$, and then use quadratic upper bound. At line four, given Lemma A.3, we use $\mathbf{0} \in \nabla_1 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D}) + \partial h(\hat{\mathbf{z}})$ and assume that \mathbf{z}_t recovers the sign entries of $\hat{\mathbf{z}}$.

$$\begin{aligned} \|Q_t^{\text{lasso-11}}\|_2 &= \|\nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D}) + \partial h(\mathbf{z}_t) - \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D})(\mathbf{z}_t - \hat{\mathbf{z}})\| \\ &= \|\nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D}) - \nabla_1 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D}) + \nabla_1 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D}) + \partial h(\mathbf{z}_t) - \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D})(\mathbf{z}_t - \hat{\mathbf{z}})\| \\ &\leq (L_{11}/2) \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2^2 + \|\partial h(\mathbf{z}_t) + \nabla_1 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D})\|_2 \\ &\leq (L_{11}/2) \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2^2 + \|\partial h(\mathbf{z}_t) - \partial h(\hat{\mathbf{z}})\|_2 \leq (L_{11}/2) \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2^2. \end{aligned} \quad (49)$$

Similarly,

$$\begin{aligned} \|Q_t^{\text{ls-11}}\|_2 &= \|\nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D}) - \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D})(\mathbf{z}_t - \hat{\mathbf{z}})\| \\ &= \|\nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D}) - \nabla_1 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D}) + \nabla_1 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D}) - \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D})(\mathbf{z}_t - \hat{\mathbf{z}})\| \\ &\leq (L_{11}/2) \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2^2 + \|\nabla_1 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D})\|_2 \leq (L_{11}/2) \|\mathbf{z}_t - \hat{\mathbf{z}}\|_2^2 + \|\partial h(\hat{\mathbf{z}})\|_2. \end{aligned} \quad (50)$$

For $Q(\hat{\mathbf{z}}, \mathbf{J}_t)$, from implicit function theorem, $Q(\hat{\mathbf{z}}, \hat{\mathbf{J}}) = 0$. Hence, we can use $\nabla_{21}^2 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D}) = -\hat{\mathbf{J}} \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D})$ in the following

$$\begin{aligned} \|Q(\hat{\mathbf{z}}, \mathbf{J}_t)\|_2 &= \|\mathbf{J}_t \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D}) + \nabla_{21}^2 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D})\|_2 = \|\mathbf{J}_t \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D}) - \hat{\mathbf{J}} \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D})\|_2 \\ &\leq \|(\mathbf{J}_t - \hat{\mathbf{J}}) \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\hat{\mathbf{z}}, \mathbf{D})\|_2 \leq L_1 \|\mathbf{J}_t - \hat{\mathbf{J}}\|_2. \end{aligned} \quad (51)$$

■

A.5 Global backward pass proof details

We re-state and proof [Theorem 4.8](#) as follows:

Theorem 4.8 (Global error of gradients). *Given the convergence results from the forward pass, ([Theorems 4.4](#) and [4.6](#)), the errors of gradients defined in [Algorithm 2](#) w.r.t global direction \mathbf{g}^* (defined in (5)) satisfy*

$$\begin{aligned} \|\mathbf{g}_{\infty}^{\text{ae-lasso}} - \mathbf{g}^*\|_2 &\leq \mathcal{O}(\|\mathbf{D} - \mathbf{D}^*\|_2^2 + \|\mathbf{D} - \mathbf{D}^*\|_2 + \|\mathbf{D} - \mathbf{D}^*\|_2 \hat{\delta}^* + \hat{\delta}^* + \hat{\delta}^{*2} + M_J \lambda \sqrt{s}) \\ \|\mathbf{g}_{\infty}^{\text{ae-ls}} - \mathbf{g}^*\|_2 &\leq \mathcal{O}(\|\mathbf{D} - \mathbf{D}^*\|_2^2 + \|\mathbf{D} - \mathbf{D}^*\|_2 + \|\mathbf{D} - \mathbf{D}^*\|_2 \hat{\delta}^* + \hat{\delta}^* + \hat{\delta}^{*2}). \end{aligned} \quad (9)$$

Proof. For $\mathbf{g}_t^{\text{dec}}$, we compute the gradient in their limit assuming infinite fresh samples $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \nabla_2 \mathcal{L}_{\mathbf{x}^i}(\mathbf{z}_t^i, \mathbf{D}) = \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [\nabla_2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D})]$ a.s.. Similar to [Theorem 4.7](#), we re-write the errors of gradients $\mathbf{g}_t^{\text{ae-lasso}}$ and $\mathbf{g}_t^{\text{ae-ls}}$ as following

$$\begin{aligned} \mathbf{g}_t^{\text{ae-lasso}} - \mathbf{g}^* &= \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [Q(\mathbf{z}^*, \mathbf{J}_t)(\mathbf{z}_t - \mathbf{z}^*)] + \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [Q_t^{21}(\mathbf{z}^*)] + \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [\mathbf{J}_t Q_t^{\text{lasso-11}}(\mathbf{z}^*)] \\ \mathbf{g}_t^{\text{ae-ls}} - \mathbf{g}^* &= \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [Q(\mathbf{z}^*, \mathbf{J}_t)(\mathbf{z}_t - \mathbf{z}^*)] + \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [Q_t^{21}(\mathbf{z}^*)] + \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [\mathbf{J}_t Q_t^{\text{ls-11}}(\mathbf{z}^*)]. \end{aligned} \quad (52)$$

where $Q_t^{21}(\mathbf{z})$, $Q_t^{\text{lasso-11}}(\mathbf{z})$, $Q_t^{\text{ls-11}}(\mathbf{z})$, and $Q(\mathbf{z}, \mathbf{J})$ are defined as in [Theorem 4.7](#). Given [Assumption 4.1](#) and [Lemma A.6](#), we find an upper bound on the *r.h.s* of the gradient errors as follows:

$$\begin{aligned} \|\mathbf{g}_t^{\text{ae-lasso}} - \mathbf{g}^*\|_2 &\leq \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [L_1 \|\mathbf{J}_t - \mathbf{J}^*\|_2 \|\mathbf{z}_t - \mathbf{z}^*\|_2 + (L_{21}/2) \|\mathbf{z}_t - \mathbf{z}^*\|_2^2] \\ &\quad + \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [M_J (L_{11}/2) \|\mathbf{z}_t - \mathbf{z}^*\|_2^2 + M_J \|\partial h(\mathbf{z}_t)\|_2 + L_D \|\mathbf{D} - \mathbf{D}^*\|_2] \\ &\leq \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [L_1 (\|\mathbf{J}_t - \hat{\mathbf{J}}\|_2 + \|\hat{\mathbf{J}} - \mathbf{J}^*\|_2) (\|\mathbf{z}_t - \hat{\mathbf{z}}\|_2 + \|\hat{\mathbf{z}} - \mathbf{z}^*\|_2)] \\ &\quad + \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [(L_{21}/2) (\|\mathbf{z}_t - \hat{\mathbf{z}}\|_2^2 + \|\hat{\mathbf{z}} - \mathbf{z}^*\|_2^2) + L_D \|\mathbf{D} - \mathbf{D}^*\|_2] \\ &\quad + \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [M_J (L_{11}/2) (\|\mathbf{z}_t - \hat{\mathbf{z}}\|_2^2 + \|\hat{\mathbf{z}} - \mathbf{z}^*\|_2^2) + M_J \|\partial h(\mathbf{z}_t)\|_2] \end{aligned} \quad (53)$$

Similarly,

$$\begin{aligned} \|\mathbf{g}_t^{\text{ae-ls}} - \mathbf{g}^*\|_2 &\leq \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [L_1 \|\mathbf{J}_t - \mathbf{J}^*\|_2 \|\mathbf{z}_t - \mathbf{z}^*\|_2 + (L_{21}/2) \|\mathbf{z}_t - \mathbf{z}^*\|_2^2] \\ &\quad + \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [M_J (L_{11}/2) \|\mathbf{z}_t - \mathbf{z}^*\|_2^2 + L_D \|\mathbf{D} - \mathbf{D}^*\|_2] \\ &\leq \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [L_1 (\|\mathbf{J}_t - \hat{\mathbf{J}}\|_2 + \|\hat{\mathbf{J}} - \mathbf{J}^*\|_2) (\|\mathbf{z}_t - \hat{\mathbf{z}}\|_2 + \|\hat{\mathbf{z}} - \mathbf{z}^*\|_2)] \\ &\quad + \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [(L_{21}/2) (\|\mathbf{z}_t - \hat{\mathbf{z}}\|_2^2 + \|\hat{\mathbf{z}} - \mathbf{z}^*\|_2^2)] \\ &\quad + \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [M_J (L_{11}/2) (\|\mathbf{z}_t - \hat{\mathbf{z}}\|_2^2 + \|\hat{\mathbf{z}} - \mathbf{z}^*\|_2^2) + L_D \|\mathbf{D} - \mathbf{D}^*\|_2]. \end{aligned} \quad (54)$$

Using the convergence errors from the forward pass ([Theorems 4.3](#) and [4.5](#)),

$$\begin{aligned} \|\mathbf{g}_t^{\text{ae-lasso}} - \mathbf{g}^*\|_2 &\leq L_1 \mathcal{O}(t \rho^{2t} + (\|\mathbf{D} - \mathbf{D}^*\|_2 + \hat{\delta}^*) t \rho^t + \rho^t (\|\mathbf{D} - \mathbf{D}^*\|_2 + \hat{\delta}^*)) \\ &\quad + L_1 \mathcal{O}(\|\mathbf{D} - \mathbf{D}^*\|_2 + \hat{\delta}^*) (\|\mathbf{D} - \mathbf{D}^*\|_2 + \hat{\delta}^*) \\ &\quad + (L_{21}/2 + M_J L_{11}/2) \mathcal{O}(\rho^t + \|\mathbf{D} - \mathbf{D}^*\|_2 + \hat{\delta}^*) + \mathcal{O}(\|\mathbf{D} - \mathbf{D}^*\|_2) + M_J \|\partial h(\mathbf{z}_t)\|_2 \end{aligned} \quad (55)$$

Hence,

$$\begin{aligned} \|\mathbf{g}_{\infty}^{\text{ae-lasso}} - \mathbf{g}^*\|_2 &\leq \mathcal{O}((\|\mathbf{D} - \mathbf{D}^*\|_2 + \hat{\delta}^*) (\|\mathbf{D} - \mathbf{D}^*\|_2 + \hat{\delta}^* + 1) + M_J \lambda \sqrt{s}) \\ &= \mathcal{O}(\|\mathbf{D} - \mathbf{D}^*\|_2^2 + \|\mathbf{D} - \mathbf{D}^*\|_2 + \|\mathbf{D} - \mathbf{D}^*\|_2 \hat{\delta}^* + \hat{\delta}^* + \hat{\delta}^{*2} + M_J \lambda \sqrt{s}) \end{aligned} \quad (56)$$

Similarly,

$$\|g_{\infty}^{\text{ae-ls}} - \mathbf{g}^*\|_2 \leq \mathcal{O}(\|\mathbf{D} - \mathbf{D}^*\|_2^2 + \|\mathbf{D} - \mathbf{D}^*\|_2 + \|\mathbf{D} - \mathbf{D}^*\|_2 \hat{\delta}^* + \hat{\delta}^* + \hat{\delta}^{*2}) \quad (57)$$

Lemma A.6 (Global bounds). *From global gradient errors in Theorem 4.8, the following are satisfied*

$$\begin{aligned} \|Q_t^{21}(\mathbf{z}^*)\|_2 &\leq (L_{21}/2)\|\mathbf{z}_t - \mathbf{z}^*\|_2^2 \\ \|Q_t^{\text{lasso-11}}(\mathbf{z}^*)\|_2 &\leq (L_{11}/2)\|\mathbf{z}_t - \mathbf{z}^*\|_2^2 + L_D\|\mathbf{D} - \mathbf{D}^*\|_2 + \|\partial h(\mathbf{z}_t)\|_2 \\ \|Q_t^{\text{ls-11}}(\mathbf{z}^*)\|_2 &\leq (L_{11}/2)\|\mathbf{z}_t - \mathbf{z}^*\|_2^2 + L_D\|\mathbf{D} - \mathbf{D}^*\|_2 \\ \|Q(\mathbf{z}^*, \mathbf{J}_t)\|_2 &\leq L_1\|\mathbf{J}_t - \mathbf{J}^*\|_2. \end{aligned} \quad (58)$$

Proof. For $Q_t^{21}(\mathbf{z}^*)$, we achieve the quadratic bound using convexity of $\nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D})$ and its domain (Assumption 3.1) and Lemma 3.3. For $Q_t^{\text{lasso-11}}(\mathbf{z}^*)$, we add and subtract $\nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D})$, and use quadratic upper bound similar to Lemma A.5. At line four, we use $\mathbf{0} \in \nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D}^*)$ (Lemma A.3) and assume that \mathbf{z}_t recovers the sign entries of \mathbf{z}^* (see Theorem 4.1 and Theorem 4.2).

$$\begin{aligned} \|Q_t^{\text{lasso-11}}(\mathbf{z}^*)\|_2 &= \|\nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D}) + \partial h(\mathbf{z}_t) - \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D})(\mathbf{z}_t - \mathbf{z}^*)\|_2 \\ &= \|\nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D}) - \nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D}) + \nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D}) + \partial h(\mathbf{z}_t) - \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D})(\mathbf{z}_t - \mathbf{z}^*)\|_2 \\ &\leq (L_{11}/2)\|\mathbf{z}_t - \mathbf{z}^*\|_2^2 + \|\partial h(\mathbf{z}_t) + \nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D})\|_2 \\ &\leq (L_{11}/2)\|\mathbf{z}_t - \mathbf{z}^*\|_2^2 + \|\partial h(\mathbf{z}_t) + \nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D}) - \nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D}^*)\|_2 \\ &\leq (L_{11}/2)\|\mathbf{z}_t - \mathbf{z}^*\|_2^2 + L_D\|\mathbf{D} - \mathbf{D}^*\|_2 + \|\partial h(\mathbf{z}_t)\|_2. \end{aligned} \quad (59)$$

Similarly,

$$\begin{aligned} \|Q_t^{\text{ls-11}}(\mathbf{z}^*)\|_2 &= \|\nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D}) - \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D})(\mathbf{z}_t - \mathbf{z}^*)\|_2 \\ &= \|\nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D}) - \nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D}) + \nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D}) - \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D})(\mathbf{z}_t - \mathbf{z}^*)\|_2 \\ &\leq (L_{11}/2)\|\mathbf{z}_t - \mathbf{z}^*\|_2^2 + \|\nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D})\|_2 \leq (L_{11}/2)\|\mathbf{z}_t - \mathbf{z}^*\|_2^2 + \|\nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D}) - \nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D}^*)\|_2 \\ &\leq (L_{11}/2)\|\mathbf{z}_t - \mathbf{z}^*\|_2^2 + L_D\|\mathbf{D} - \mathbf{D}^*\|_2. \end{aligned} \quad (60)$$

For $Q(\mathbf{z}^*, \mathbf{J}_t)$, from implicit function theorem, $Q(\mathbf{z}^*, \mathbf{J}^*) = 0$. Hence, we can use $\nabla_{21}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D}) = -\mathbf{J}^* \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D})$ in the following

$$\begin{aligned} \|Q(\mathbf{z}^*, \mathbf{J}_t)\|_2 &= \|\mathbf{J}_t \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D}) + \nabla_{21}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D})\|_2 = \|\mathbf{J}_t \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D}) - \mathbf{J}^* \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D})\|_2 \\ &\leq \|(\mathbf{J}_t - \mathbf{J}^*) \nabla_{11}^2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D})\|_2 \leq L_1\|\mathbf{J}_t - \mathbf{J}^*\|_2. \end{aligned} \quad (61)$$

A.6 Interpretability

Theorem 5.1 (Interpretable unrolled network). *Consider the dictionary learning optimization of the form $\min_{\mathbf{Z}, \mathbf{D}} \frac{1}{2}\|\mathbf{X} - \mathbf{D}\mathbf{Z}\|_2^2 + \lambda\|\mathbf{Z}\|_1 + \omega/2\|\mathbf{D}\|_F^2$, where $\mathbf{X} = [\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^n] \in \mathbb{R}^{m \times n}$ and $\mathbf{Z} = [\mathbf{z}^1, \mathbf{z}^2, \dots, \mathbf{z}^n] \in \mathbb{R}^{p \times n}$. Let $\tilde{\mathbf{Z}}$ be the given converged sparse codes, then stationary points of the problem w.r.t the network weights (dictionary) follows $\tilde{\mathbf{D}} = \mathbf{X}\mathbf{G}^{-1}\tilde{\mathbf{Z}}^T$, where we denote $\mathbf{G} := (\tilde{\mathbf{Z}}^T \tilde{\mathbf{Z}} + \omega \mathbf{I})$.*

Proof. For all stationary points, the objective gradient is $\mathbf{0}$ with respect to the dictionary, i.e.,

$$\mathbf{0} = (\mathbf{X} - \tilde{\mathbf{D}}\tilde{\mathbf{Z}})\tilde{\mathbf{Z}}^T + \omega\tilde{\mathbf{D}} \quad (62)$$

where $\tilde{\mathbf{D}}$ is the learned dictionary at convergence. Re-arranging the terms, we get

$$\tilde{\mathbf{D}} = \mathbf{X}\tilde{\mathbf{Z}}^T(\tilde{\mathbf{Z}}\tilde{\mathbf{Z}}^T + \omega\mathbf{I})^{-1} \quad (63)$$

Using the relation $\mathbf{A}^T(\mathbf{A}\mathbf{A}^T + \omega\mathbf{I})^{-1} = (\mathbf{A}^T\mathbf{A} + \omega\mathbf{I})^{-1}\mathbf{A}^T$, we can re-write the solution as

$$\tilde{\mathbf{D}} = \mathbf{X}\mathbf{G}^{-1}\tilde{\mathbf{Z}}^T \quad (64)$$

where we denote $\mathbf{G} := (\tilde{\mathbf{Z}}^T \tilde{\mathbf{Z}} + \omega \mathbf{I})$. ■

B Appendix - future works and limitations

Beyond dictionary learning Our results are founded on three main properties: Lipschitz differentiability of the loss relating to the parameter of interest, proximal gradient descent, and strong convexity in finite-iteration. The findings can be applied to other min-min optimization problems, e.g., ridge regression and logistic regression, following such properties. For example, our analysis generalizes to the unrolled network in (Tolooshams et al., 2020) for learning dictionaries using data from the natural exponential family. In this case, the least-squares loss is replaced with negative log-likelihood, and the dictionary models the expectation of the data.

Limitations Finite-iteration support selection (Proposition 4.1) (Hale et al., 2007) and strong convexity may seem stringent going beyond dictionary learning. Ablin et al. (2020) discuss generalization of local gradient convergence by relaxing strong convexity to the p -Łojasiewicz property (Ablin et al., 2020; Attouch & Bolte, 2009). We considered the noiseless setting and conjecture that the relative comparison of the gradients in the presence of noise still holds, where the upper bounds will involve an additional noise term. We focused on infinite sample convergence to highlight the relative differences between the gradients. We leave for future work the derivation of finite-sample upper bounds, a step similar to (Chatterji & Bartlett, 2017; Arora et al., 2015).

C Appendix - details of experiments

PUDLE is developed using PyTorch (Paszke et al., 2017) on Python. We used one GeForce GTX 1080 Ti GPU.

C.1 Numerical experiments for theories

Dataset We generated $n = 10,000$ samples following the model (2). We sampled $\mathbf{D}^* \in \mathbb{R}^{50 \times 100}$ from zero-mean Gaussian distribution, and then normalized the columns. The codes are 5-sparse with their support uniformly chosen at random and their amplitudes are sampled from Uniform(1, 2).

Training We let $T = 200$, $\lambda = 0.2$, and $\alpha = 0.2$. The dictionary is initialized to $\mathbf{D} = \mathbf{D}^* + \tau_B \mathbf{B}$ with $\mathbf{B} \sim \mathcal{N}(\mathbf{0}, \frac{1}{m} \mathbf{I})$. For Figures 2, 3 and 4a, we set $\tau_B \approx 0.55/\log m$. We picked this closeness to highlight the gradient directions in a closer neighbourhood of the dictionary. For Figures 4b and 4c, we chose much larger noise level, $\tau_B \approx 2.8/\log m$. The network is trained for 600 epochs with full-batch gradient descent using Adam optimizer (Kingma & Ba, 2014) with learning rate of 10^{-3} and $\epsilon = 10^{-8}$. The learned dictionary is evaluated based on the error $\|\mathbf{D} - \mathbf{D}^*\|_2 / \|\mathbf{D}^*\|_2$. The results and conclusion were consistent across various realizations of the dataset and across various optimizers. Hence, in the main paper, the figures visualize results of one realization.

C.2 Dictionary learning

Dataset We generated $n = 50,000$ samples following (2). We let $m = 1000$ and $p = 1500$, and sample \mathbf{D}^* from zero-mean Gaussian distribution, and then normalized the columns. The sparse codes \mathbf{z}^i are 10, 20, 40-sparse, where their supports are chosen uniformly at random and their amplitudes are sampled from Uniform(1, 2).

Training The dictionary is initialized to $\mathbf{D} = \mathbf{D}^* + \tau_B \mathbf{B}$ with $\mathbf{B} \sim \mathcal{N}(\mathbf{0}, \frac{1}{m} \mathbf{I})$ where $\tau_B \approx 1/\log m$. We let $\lambda = 0.2$, and $\alpha = 0.2$, and $T = 100$. The network is trained for 1,000 iterative updates with batch-size of 50 using Adam (Kingma & Ba, 2014) with learning rate of 10^{-3} and $\epsilon = 10^{-3}$. For decay method, ν is decreased in value by 0.005 every 100 update iterations. Each filter is normalized after every update. The learned dictionary is evaluated based on the relative error $\|\mathbf{D} - \mathbf{D}^*\|_2 / \|\mathbf{D}^*\|_2$.

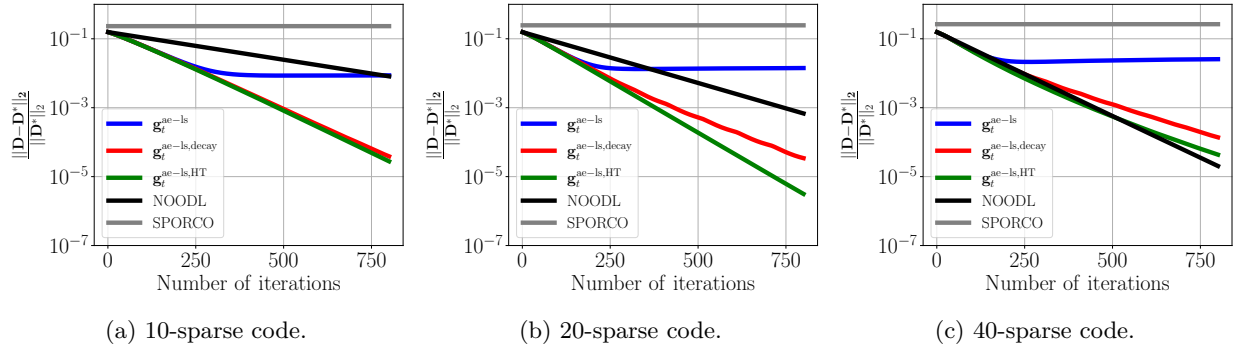


Figure 10: Dictionary learning convergence using $\mathbf{g}_t^{\text{ae-ls}}$ compared to NOODL (Rambhatla et al., 2018) and SPORCO (Wohlberg, 2017).

C.3 Image denoising

Training We trained PUDLE where the dictionary is convolutional with 64 filters of size 9×9 and strides of 4. The encoder unrolls for $T = 15$, and the step size is set to $\alpha = 0.1$. Unlike the theoretical analysis where full-batch gradient descent is studied, the network is trained stochastically with Adam optimizer (Kingma & Ba, 2014) with a learning rate of 10^{-4} and $\epsilon = 10^{-3}$ for 250 epochs. At every training iteration, a random 129×129 patch is cropped and a zero-mean Gaussian noise with a standard deviation of 25 is added. We report results in terms of the peak signal-to-noise ratio (PSNR). The standard deviation of the test PSNR across multiple noise realizations was lower than 0.02 dB for all the methods. Hence, we only reported the mean PSNR of the test set.

C.4 Interpretable sparse coding and dictionary learning

We focused on digits of $\{0, 1, 2, 3, 4\}$ of MNIST. We set $T = 15$, $\lambda = 0.7$, and $\alpha = 1$. The dictionary dimensions are $m = 784$ and $p = 500$. We trained the network for 200 epochs using Adam optimizer with a learning rate of 10^{-4} and batch size of 32. For construction of \mathbf{G} , ω is set to 0.001. For Figure 7, we computed the image contributions using 6,000 randomly chosen training images. The Gram matrix used in Figure 8, is constructed by 6,000 training examples, and the reconstruction is from the 200 most contributed training images.