

Interpretable Unrolled Dictionary Learning Networks

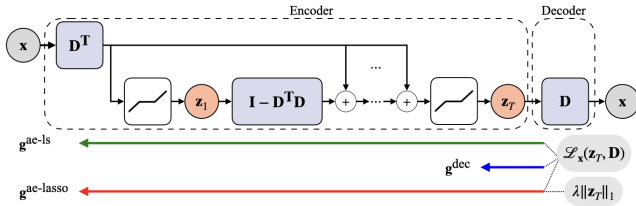
Bahareh Tolooshams and Demba Ba

School of Engineering and Applied Sciences, Harvard University, Cambridge, MA

I. INTRODUCTION

The dictionary learning problem, representing data $\mathbf{x} \in \mathbb{R}^m$ as a combination of a few atoms from a dictionary $\mathbf{D} \in \mathbb{R}^{m \times p}$, has long stood as a popular method for learning representations in statistics and signal processing [1, 2, 3, 4]. The most popular dictionary learning algorithm alternates between sparse coding and dictionary update steps. Sparse coding has been utilized to construct neural architectures through recurrent sparsifying encoders [5], initiating a growing literature on constructing interpretable unrolled networks [6, 7]. We offer the theoretical analysis of unrolled sparse coding. We address the following challenge; the vanilla unrolled sparse coding computes a biased code estimate; this results in a biased estimate of the backward gradient. We reduce this bias and demonstrate unrolled interpretability.

Given \mathbf{x} and \mathbf{D} , the problem of recovering the sparse coefficients $\mathbf{z} \in \mathbb{R}^p$ is referred to as sparse coding, and can be solved through the lasso [8] $\ell_{\mathbf{x}}(\mathbf{D}) := \min_{\mathbf{z} \in \mathbb{R}^p} \mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D}) + h(\mathbf{z})$ where $\mathcal{L}_{\mathbf{x}}(\mathbf{z}, \mathbf{D}) = \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|_2^2$, and $h(\mathbf{z}) = \lambda \|\mathbf{z}\|_1$. The problem aims to recover a dictionary \mathbf{D}^* that generates the data, i.e., $\mathbf{x} = \mathbf{D}^* \mathbf{z}^*$ where \mathbf{z}^* is sparse. Prior to unrolled networks, gradient-based dictionary learning relied on analytic gradients. With unrolled networks, backpropagation gained attention for parameter estimation [9, 10, 11]¹.



1: Unrolled dictionary learning.

Unrolled dictionary learning is constructed as following: sparse coding is converted into an encoder by unfolding T iterations of ISTA ($\mathbf{z}_{t+1} = \Phi(\mathbf{z}_t, \mathbf{D}) = \mathcal{P}_{\alpha\lambda}(\mathbf{z}_t - \alpha \nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_t, \mathbf{D}))$) with $\mathcal{P}_b(v) \triangleq \text{sign}(v) \max(|v| - b, 0)$, [12, 13]. The decoder is $\hat{\mathbf{x}} = \mathbf{D}\mathbf{z}_T$. We recover \mathbf{D}^* by backpropagated gradient (i.e., $\mathbf{D}^{(l+1)} = \mathbf{D}^{(l)} - \eta \mathbf{g}_T^{(l)}$) (See Fig. 1). Backpropagation through the decoder results in the analytic gradient $\mathbf{g}_T^{\text{dec}} \triangleq \nabla_2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_T, \mathbf{D})$. The gradients $\mathbf{g}_T^{\text{ae-lasso}} \triangleq \nabla_2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_T, \mathbf{D}) + \frac{\partial \mathbf{z}_T}{\partial \mathbf{D}} (\nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_T, \mathbf{D}) + \partial h(\mathbf{z}_T))$ and $\mathbf{g}_T^{\text{ae-ls}} \triangleq \nabla_2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_T, \mathbf{D}) + \frac{\partial \mathbf{z}_T}{\partial \mathbf{D}} \nabla_1 \mathcal{L}_{\mathbf{x}}(\mathbf{z}_T, \mathbf{D})$ are computed by backpropagation through the autoencoder using the lasso and least-squares objectives. We show how using $\mathbf{g}_T^{\text{ae-ls}}$ is a better gradient estimator to recover \mathbf{D}^* . The desired direction is $\mathbf{g}^* \triangleq \mathbb{E}_{\mathbf{x} \in \mathcal{X}} [\nabla_2 \mathcal{L}_{\mathbf{x}}(\mathbf{z}^*, \mathbf{D})]$.

II. MAIN RESULTS

Assumptions: The code \mathbf{z}^* is at most s -sparse with the support $S^* = \text{supp}(\mathbf{z}^*)$. Given the support, $\mathbf{z}_{S^*}^*$ is i.i.d, $\mathbb{E}[\mathbf{z}_{(j)}^* | j \in S^*] = 0$ and $\mathbb{E}[\mathbf{z}_{(S)}^* \mathbf{z}_{(S^*)}^T | S^*] = \mathbf{I}$. \mathbf{D}^* is μ -incoherent with $\mu = \mathcal{O}(\log(m))$. $\|\mathbf{D}_j^*\|_2 = 1$ and $\|\mathbf{D}^*\|_2 = \mathcal{O}(\sqrt{p/m})$, and $p = \mathcal{O}(m)$. $\forall j \|\mathbf{D}_j^{(0)} - \mathbf{D}_j^*\|_2 \leq \delta$ and $\|\mathbf{D}^{(0)} - \mathbf{D}^*\|_2 \leq 2\|\mathbf{D}^*\|_2$. $\mathbf{D}^{(l)}$ is μ_l -incoherent and $\|\mathbf{D}_j^{(l)} - \mathbf{D}_j^*\|_2 \leq \delta_l$ with $\delta_l = \mathcal{O}^*(1/\log p)$.

Results: Thm. II.1 provides an upper bound (as a function of dictionary error, amount of unrolling, and sparse regularizer) on the

error between the true code \mathbf{z}^* and the sparse latent code \mathbf{z}_t . The λ term in the upper bound shows that the code error when we strictly perform ℓ_1 -norm based sparse coding does not go to zero.

Theorem II.1. *If $s = \mathcal{O}^*(\sqrt{m}/\mu \log m)$, and the regularizer and step size are $\lambda_t^{(l)} = \lambda = \frac{\mu_t}{\sqrt{m}} \|\mathbf{z}^* - \mathbf{z}_0\|_1 + \alpha_\gamma = \Omega(\frac{s \log m}{\sqrt{m}})$ and $\alpha^{(l)} \leq 1 - \frac{2\lambda_t^{(l)} - (1 - \frac{\delta_t^2}{2})C_{\min}}{\lambda_{t-1}^{(l)}}$, then with high probability, $|\mathbf{z}_{t,(j)}^{(l)} - \mathbf{z}_{(j)}^*| \leq \mathcal{O}(\sqrt{s\|\mathbf{D}_j^{(l)} - \mathbf{D}_j^*\|_2 + e_{t,j}^{(l)} + \lambda)}$ where $e_{t,j}^{(l)} \rightarrow 0$ as $t \rightarrow \infty$.*

Given the forward pass, Thm. II.2 shows that training unrolled dictionary learning network using gradient descent with $\mathbf{g}_T^{\text{dec}}$ results in a contractive mapping, hence recovery of the dictionary up to a \mathbf{D}^* neighborhood mainly characterized by the regularization parameter λ .

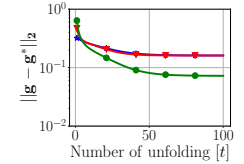
Theorem II.2. *If $s = \mathcal{O}(\sqrt{m})$, $\eta = \mathcal{O}(\frac{p}{s(1-\delta_t^2/2)})$, and the regularizer λ and α are set according to above, then with high probability $\|\mathbf{D}_j^{(l+1)} - \mathbf{D}_j^*\|_2^2 \leq (1 - \psi)\|\mathbf{D}_j^{(l)} - \mathbf{D}_j^*\|_2^2 + \epsilon_\lambda^{(l)}$ where $\epsilon_\lambda^{(l)} := \eta \frac{2p}{s(1 - (\mathbf{D}_j^{(l)} - \mathbf{D}_j^*))} \lambda^2$.*

To reduce this bias in the dictionary update, we propose to use backpropagation using the reconstruction loss $\mathbf{g}_T^{\text{ae-ls}}$ instead of the analytic gradient $\mathbf{g}_T^{\text{dec}}$. Theorem II.3 compares the gradients for appropriately large T ; it shows that $\mathbf{g}_T^{\text{ae-ls}}$ is a better estimator of the desired direction \mathbf{g}^* .

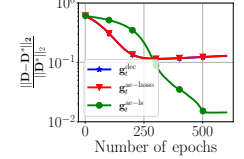
Theorem II.3. *$\mathbf{g}_T^{\text{ae-lasso}}$ is equivalent to $\mathbf{g}_T^{\text{dec}}$ as $T \rightarrow \infty$ (Fig. 2a). $\|\mathbf{g}_T^{\text{ae-lasso}} - \mathbf{g}^*\|_2 \leq \mathcal{O}(\|\mathbf{D} - \mathbf{D}^*\|_2 + \delta^* + C_{\text{lasso}})$ and $\|\mathbf{g}_T^{\text{ae-ls}} - \mathbf{g}^*\|_2 \leq \mathcal{O}(\|\mathbf{D} - \mathbf{D}^*\|_2 + \delta^*)$, where δ^* is proportional to the biased code estimate, and $C_{\text{lasso}} := \mathcal{O}(\lambda\sqrt{s})$. Hence, $\mathbf{g}_T^{\text{ae-ls}}$ is a better estimator of \mathbf{g}^* (Fig. 2a), and \mathbf{D}^* neighbourhood at which $\mathbf{g}_T^{\text{ae-ls}}$ is guaranteed to converge to is smaller than of the $\mathbf{g}_T^{\text{ae-lasso}}$ and $\mathbf{g}_T^{\text{dec}}$ (Fig. 2b).*

Interpretability: We build a mathematical relation between the network weights, training data, and test reconstruction. Thm. II.4 characterizes stationary points of the trained network and proves that the dictionary interpolates the training data, i.e., $\tilde{\mathbf{D}}_j = \mathbf{X}(\mathbf{G}^{-1}\mathbf{w}_j) = \sum_{k=1}^n (\mathbf{G}^{-1}\mathbf{w}_j)_k \mathbf{x}^k$ (Fig. 3a, green for high and red for low contribution). We write the reconstruction of a new example \mathbf{x}^j as a linear combination of all the training examples, i.e., $\hat{\mathbf{x}}^j = \tilde{\mathbf{D}}\hat{\mathbf{z}}^j = \sum_{k=1}^n \beta_k^j \mathbf{x}^k$ where $\beta_k^j = \sum_{a=1}^n \mathbf{G}_{ka}^{-1} \langle \tilde{\mathbf{z}}^a, \hat{\mathbf{z}}^j \rangle$ (Fig. 3b, images with high contribution (green) are similar to the input image, and those with low (red) are different).

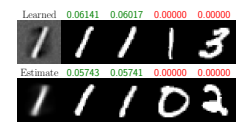
Theorem II.4. *Consider $\min_{\mathbf{Z}, \mathbf{D}} \frac{1}{2} \|\mathbf{X} - \mathbf{D}\mathbf{Z}\|_F^2 + \lambda \|\mathbf{Z}\|_1 + \omega/2 \|\mathbf{D}\|_F^2$, where $\mathbf{X} = [\mathbf{x}^1, \dots, \mathbf{x}^n] \in \mathbb{R}^{m \times n}$ and $\mathbf{Z} = [\mathbf{z}^1, \dots, \mathbf{z}^n] \in \mathbb{R}^{p \times n}$. Let $\tilde{\mathbf{Z}}$ be the given converged codes, then network stationary points follows $\tilde{\mathbf{D}} = \mathbf{X}\mathbf{G}^{-1}\tilde{\mathbf{Z}}^T$, where we denote $\mathbf{G} := (\tilde{\mathbf{Z}}^T\tilde{\mathbf{Z}} + \omega\mathbf{I})$.*



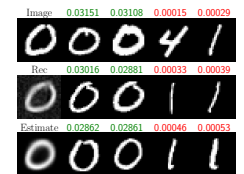
2a: Gradient.



2b: Dictionary.



3a: Contribution to $\tilde{\mathbf{D}}$.



3b: Contribution for $\hat{\mathbf{x}}^j$.

¹This work is presented as a talk at the Conference on the Mathematical Theory of Deep Neural Networks, 2022.

REFERENCES

- [1] T. Hastie, R. Tibshirani, and M. Wainwright, *Statistical learning with sparsity: the lasso and generalizations*. CRC press, 2015.
- [2] M. Elad, *Sparse and redundant representations: from theory to applications in signal and image processing*. Springer Science & Business Media, 2010.
- [3] R. Jenatton, J. Mairal, G. Obozinski, and F. Bach, “Proximal methods for hierarchical sparse coding,” *The Journal of Machine Learning Research*, vol. 12, pp. 2297–2334, 2011.
- [4] J. Mairal, J. Ponce, G. Sapiro, A. Zisserman, and F. Bach, “Supervised dictionary learning,” in *Proceedings of Advances in Neural Information Processing Systems*, vol. 21, 2009, pp. 1–8.
- [5] K. Gregor and Y. LeCun, “Learning fast approximations of sparse coding,” in *Proceedings of international conference on international conference on machine learning*, 2010, pp. 399–406.
- [6] J. R. Hershey, J. L. Roux, and F. Weninger, “Deep unfolding: Model-based inspiration of novel deep architectures,” *arXiv:1409.2574*, pp. 1–27, 2014.
- [7] V. Monga, Y. Li, and Y. C. Eldar, “Algorithm unrolling: Interpretable, efficient deep learning for signal and image processing,” *arXiv:1912.10557*, pp. 1–27, 2019.
- [8] R. Tibshirani, “Regression shrinkage and selection via the lasso,” *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.
- [9] B. Tolooshams, S. Dey, and D. Ba, “Deep residual autoencoders for expectation maximization-inspired dictionary learning,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 6, pp. 2415–2429, 2021.
- [10] B. Tolooshams, A. Song, S. Temereanca, and D. Ba, “Convolutional dictionary learning based auto-encoders for natural exponential-family distributions,” in *Proceedings of the 37th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, H. D. III and A. Singh, Eds., vol. 119. PMLR, 13–18 Jul 2020, pp. 9493–9503.
- [11] B. Malézieux, T. Moreau, and M. Kowalski, “Understanding approximate and unrolled dictionary learning for pattern recovery,” in *Proceedings of International Conference on Learning Representations*, 2022.
- [12] 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*, vol. 57, no. 11, pp. 1413–1457, 2004.
- [13] T. Blumensath and M. E. Davies, “Iterative thresholding for sparse approximations,” *Journal of Fourier analysis and Applications*, vol. 14, no. 5-6, pp. 629–654, 2008.