
Sparse L^1 -Autoencoders for Scientific Data Compression

Matthias Chung*

Department of Mathematics
Emory University
Atlanta, GA 30322, USA
matthias.chung@emory.edu

Rick Archibald

Computer Science and Mathematics Division
Oak Ridge National Laboratory
Oak Ridge, TN 37830, USA
archibaldrk@ornl.gov

Paul Atzberger

Department of Mathematics
University of California Santa Barbara
Santa Barbara, CA 93106, USA
atzberg@gmail.com

Jack Michael Solomon

Department of Mathematics
Emory University
Atlanta, GA 30322
jack.michael.solomon@emory.edu

Abstract

Scientific datasets present unique challenges for machine learning-driven compression methods, including more stringent requirements on accuracy and mitigation of potential invalidating artifacts. Drawing on results from compressed sensing and rate-distortion theory, we introduce effective data compression methods by developing autoencoders using high dimensional latent spaces that are L^1 -regularized to obtain sparse low dimensional representations. We show how these information-rich latent spaces can be used to mitigate blurring and other artifacts to obtain highly effective data compression methods for scientific data. We demonstrate our methods for short angle scattering (SAS) datasets showing they can achieve compression ratios around two orders of magnitude and in some cases better. Our compression methods show promise for use in addressing current bottlenecks in transmission, storage, and analysis in high-performance distributed computing environments. This is central to processing the large volume of SAS data being generated at shared experimental facilities around the world to support scientific investigations. Our approaches provide general ways for obtaining specialized compression methods for targeted scientific datasets.

1 Introduction & Background

Autoencoders are one of the most prominent and highly successful types of neural networks [3, 18, 24], they are attractive by their conceptual simplicity (learning parameterized encoder e and decoder d such that $x \approx d(e(x))$), their strong connection to well-established mathematical concepts [41], and versatile when utilized as generative models [23]. Autoencoders have wide applicability in unsupervised learning environments ranging from denoising [51], anomaly detection [45], image and audio compression [26, 48], to generic recommender systems [55].

Most commonly, with their inherent structure of low dimensional latent spaces, autoencoders are a natural and data-driven machine learning technique for model reduction and data compression [10, 31]. Autoencoders' nonlinear characteristics make them a valuable complement to established techniques

*

like principal component analysis and singular value decomposition ([52]), Fourier analysis ([28]), reduced order models ([54]), and dictionary learning approaches ([47]).

Mathematically, an autoencoder is a nonlinear parameterized mapping $a : \mathcal{X} \rightarrow \mathcal{X}$ with a functional composition $a = d \circ e$, into *encoder* $e : \mathcal{X} \times \Theta_e \rightarrow \mathcal{Z}$ and *decoder* $d : \mathcal{Z} \times \Theta_d \rightarrow \mathcal{X}$ with trainable parameters $\theta_e \in \Theta_e \subset \mathbb{R}^{n_e}$ and $\theta_d \in \Theta_d \subset \mathbb{R}^{n_d}$. The feature space $\mathcal{Z} \subset \mathbb{R}^\ell$ is referred to as *latent space* while $\mathcal{X} \subset \mathbb{R}^n$ is the *data space*.

While terminology varies, autoencoders are classified based on their model configuration. Standard autoencoders with small dimensional latent spaces $\ell < n$ are referred to as *undercomplete* and are widely utilized, while autoencoders with large dimensional latent spaces $\ell > n$ are referred to as *overcomplete* and are less common. A *shallow* autoencoder is characterized by having only one, while a *deep* autoencoder has more than one hidden layer. Autoencoders typically maintain network symmetry, i.e., the structure of e and d are mirrored such that the decoder transformations mimic the encoder in reverse order. Its general structure makes autoencoders flexible, for instance, giving up the mapping back into the input space has led to encoder-decoder networks which are widely used as likelihood-free surrogate models for physical forward propagation [56] and even inverse modeling [1].

Despite its versatility, various drawbacks associated with autoencoders have been identified. For instance, the undercomplete “hour-glass” autoencoder shape compresses input signals $x \in \mathcal{X}$ into the low-dimensional latent space $z = e(x; \theta_e) \in \mathcal{Z}$ and may lead to corrupted reconstructions, e.g., blurring artifacts in the reconstructed images [36]. This drawback has been highlighted in the generative process of autoencoders, i.e., variational autoencoders. Autoencoders with a large number of trainable network parameters, such as overcomplete autoencoder, tend to overfit, resulting in “identity mappings” countering one of the main purposes of autoencoders: removing unwanted artifacts from the input x [22]. A further drawback includes incorporating scientific and physical features into the network remains a major hurdle, where some initial research is making strides towards this goal [2, 9, 29].

To mitigate the challenges outlined above, we break with one common assumption of autoencoders, that is, we consider utilizing an overcomplete autoencoder framework with sparsity promoting mappings of the latent variable, illustrated in Figure 1. Overcomplete autoencoders are prone to severely overfit without taking mitigating measures. Hence, imposing sparsity onto the latent space variable is a regularizing measure and various strategies have been proposed. Despite being tremendously successful, we recognize that overcomplete autoencoders with sparsity-promoting features in the latent space are severely underutilized.

Sparse autoencoders have first been introduced in the 2010s with pioneering work from various research groups including [21, 33, 34, 39]. Here, various strategies have been entertained to promote sparsity of the latent variable. For instance, by selecting a fixed amount of k nonzero elements of the latent variable z with maximal reconstruction features [34]. Another approach limits the number of active latent components by utilizing a binary Bernoulli random variable model realized through a Kullback-Leibner divergence penalty [39]. A third approach, which we will follow here, is to use the compressed sensing framework via L^1 regularization [21]. Recent years have brought advances, various extensions of sparse autoencoders have been developed, and scientific applications considered [4, 22, 30, 32, 35, 40, 46], however, despite its successes sparse autoencoder have yet to find its way into mainstream applications. We like to point out that the term sparse autoencoders may not refer to sparsity induced onto the the latent variable, but sparsity imposed on the network parameters θ , e.g., see [32, 46].

To the best of our knowledge, a core application on the compressive feature of sparse autoencoder has not yet been fully addressed. Hence our proposed method utilizes sparse latent space signals to efficiently store input signals. Furthermore, compared to prior work, we consider not only the sparsity of the latent variable z through L^1 regularization but promote sparsity on a signal $f(z)$. The functional f provides the possibility to promote structure within the latent variable z , e.g., through a total variation type (generalized lasso).

Sparse autoencoder shares limitations of the general class of autoencoder, that is, the ability to effectively generalize to novel data instances, particularly when the training dataset does not accurately reflect the characteristics of the testing dataset.

Our work is structured as follows, in Section 2 we introduce our proposed sparsity promoting autoencoder for data compression tasks, discuss its numerical applications in Section 3, and provide concluding remarks and discuss future work in Section 4

2 Sparse Autoencoders for Scientific Data Compression

Scientific datasets present distinct challenges for machine learning-driven compression methods given how they are used applications. Scientific investigations often require more stringent requirements on individual sample reconstruction accuracy and mitigation of artifacts such as blurring. To help address these challenges, we use large dimensional information-rich latent spaces that are reduced by using sparsity regularizations to obtain representations amenable to further compression. Historically, embedding sparse signals into large dimensional vector spaces has had a major impact on signal processing starting in the 1990s with the *compressed sensing framework* [5, 6, 15, 49]. We further develop this strategy for autoencoders using latent space dimensions larger than the feature space of the data, e.g., $\ell > m$. Without imposing additional restrictions, learning the encoder $e(\cdot; \theta_e)$ and decoder $d(\cdot; \theta_d)$ would be an ill-posed problem without advantages for later compression. Instead, we leverage the large dimensional spaces to allow our encoders to further process information and obtain well-posed compressed signals for the data $x \in \mathcal{X}$ in \mathcal{Z} . We enforce sparsity on features of the latent vector z . Let θ_e and θ_d be the trainable network parameter of the encoder and decoder, respectively, ideally, we may formulate the network training as

$$\min_{(\theta_e, \theta_d) \in \Theta_e \times \Theta_d} \mathbb{E} \|f[e(x; \theta_e)]\|_0 \quad \text{subject to} \quad \|d(e(x; \theta_e); \theta_d) - x\|_2 \leq \delta, \quad (1)$$

where \mathbb{E} denotes the expectation over the data x , $\|\cdot\|_2$ the L^2 -norm, $\|e(x; \theta_e)\|_0$ is defined as the cardinality of nonzero elements in $e(x; \theta_e)$, and $\delta > 0$ represents a desired reconstruction quality. We further let $f: \mathcal{Z} \rightarrow \mathcal{F}$ be a predefined operator we refer to as the *sparse structure selector*.

Solving (1) is NP-hard and efficient approximation approaches need to be utilized, [17]. Under the restricted isometry properties [6], we may reformulate using an L^1 convex relaxation of (1) which leads to the generalized lasso problem

$$\min_{\theta_e, \theta_d \in \Theta} \mathbb{E} \|d(e(x; \theta_e); \theta_d) - x\|_2^2 + \lambda \|f[e(x; \theta_e)]\|_1, \quad (2)$$

for suitable sparsity enforcing $\lambda > 0$, see [5, 11] for details. This can be viewed as a rate-distortion objective, where $\lambda \|\cdot\|_1$ serves as a measure of the compression rate and $\|\cdot\|_2$ for the reconstruction distortion [13].

What are the benefits of introducing an operator f and not directly inducing sparsity on the latent variable z with a standard L^1 regularization on the latent variable, e.g., $\|e(x; \theta_e)\|_1$? With sparsity enforced only on each component of z individually, the latent variable carries only minimal interpretability or structure [35]. Now, a mapping f may enforce structure to the latent space variable z . The function f may remedy this disadvantage and add additional geometric interpretability. For simplicity, we use a common total variation regularization approach $f[\cdot] = \nabla(\cdot)$, where ∇ is the gradient operation. In practice, for finite dimensional spaces \mathcal{Z} , we approximate this operation by the finite difference operator $(f[z])_i = (z_{i+1} - z_i)/h$, e.g., with $h = 1$. This is also used to help in clustering of information in the latent space.

Only because our autoencoder has a large dimensional latent space with a large amount of network parameter our approach does not just encode and decode the each training image individually within

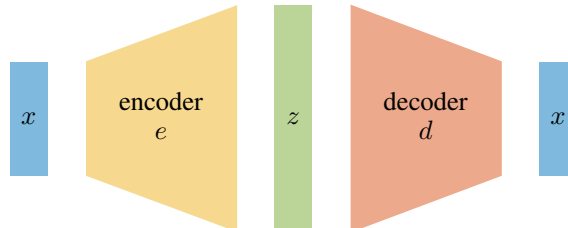


Figure 1: Overcomplete autoencoder architecture, where the latent space dimension ℓ is bigger than the input dimension n . Sparsity on the latent variable is imposed via L^1 -type regularization.

the network. But the numerical results in 3 clearly show that our approach generalizes well to testing data.

Further, note that the regularization parameter λ balances potential over- and underfitting. Small λ values may generate autoencoder with identity mappings for training data but may not generate any sparsity within the sparse structure selector. On the other hand large λ may produce significant sparsity while missing to reconstruct the input signal x . Cross validation techniques are readily available for calibrating the hyperparameter λ but is subject to further investigations.

Given a representative set of (unsupervised) training samples $\{x_j\}_{j=1}^m$ we minimize the empirical generalized lasso

$$\min_{\theta_e, \theta_d \in \Theta} \frac{1}{m} \sum_{j=1}^m \|d(e(x_j; \theta_e); \theta_d) - x_j\|_2^2 + \lambda \|f(e(x_j; \theta_e))\|_1, \quad (3)$$

Our developed methods leverage results in compressed sensing showing promise for having a significant impact on scientific compression techniques. Under mild assumptions, compressed sensing has shown high compression rates, far below theoretical Nyquist rates [7, 8]. Benefiting from its advantages in a trainable deep neural networks provides significant compression rates while maintaining high accuracy of the signal itself Section 3. Theoretically our methods have also connections to dictionary learning frameworks. Sparse dictionary learning provides good reconstruction of a sparse selection of dictionary atoms [16, 25, 38, 50]. Here, since linear, zeros in its dictionary representation does not carry any information. However, utilizing sparsifying autoencoders allows for nonlinear transformations and therefore enriches information carried by the latent variable. Consequently, even “zero” elements in the latent variable z carry information of the underlying signal. Furthermore, the generic design of the neural network architecture provides an additional level of flexibility toward efficient encoding and decoding of the underlying signals x . For (b), under mild assumptions, compressed sensing has shown high compression rates, far below theoretical Nyquist rates [7, 8]. Benefiting from its advantages in a trainable deep neural networks may provide significant compression rates while maintaining high accuracy of the signal itself.

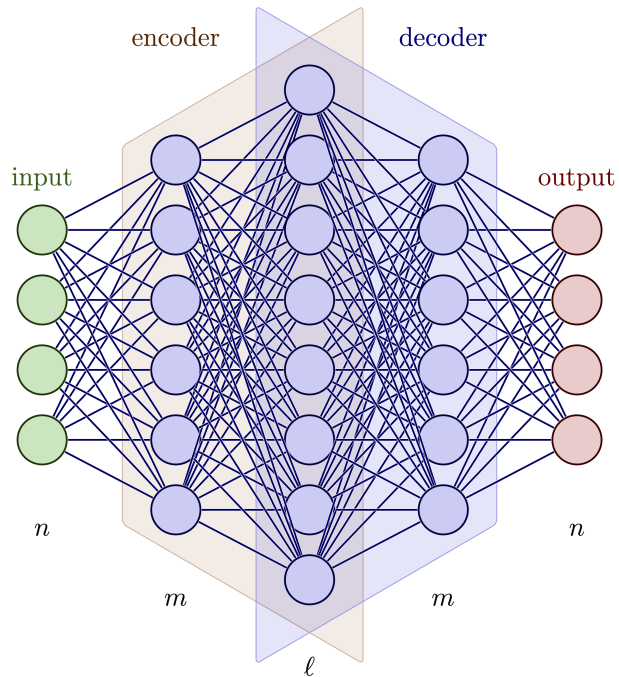


Figure 2: Autoencoder architecture for numerical examples. Network is a fully connected five layer symmetric neural network with hidden layer size (m, ℓ, m) , ReLU activation function between each layer, and input/output size n .

We also develop methods to further compress our representations z by combining our approaches with lossy quantization and lossless entropy encoding [13, 19]. We represent the information in z as a list of the m indices of nonzero entries i_1, \dots, i_m and the weights at these locations $w_k = z_{i_k}$. We represent this by storing the distances between successive indices $\delta_k = i_{k+1} - i_k$ along with a termination symbol ι to obtain the sequence $\delta_1, \dots, \delta_{m-1}, \iota$. We expect in practice for most datasets that the probability distribution $\rho(\delta_i)$ over the differences δ_i will tend to skew toward the smaller values, such as having $\delta_i < \ell/2$ for most entries. By modeling this distribution we can obtain gains in the compression using an entropy encoder. To obtain a lossless entropy encoding for our model probability distribution $\rho(\delta)$, we develop a lossless arithmetic coding method \mathcal{A} to obtain $c = \mathcal{A}(\delta; \rho(\cdot))$, [27, 44, 53]. To compress the weights $\{w_k\}$ we use lossy quantization of

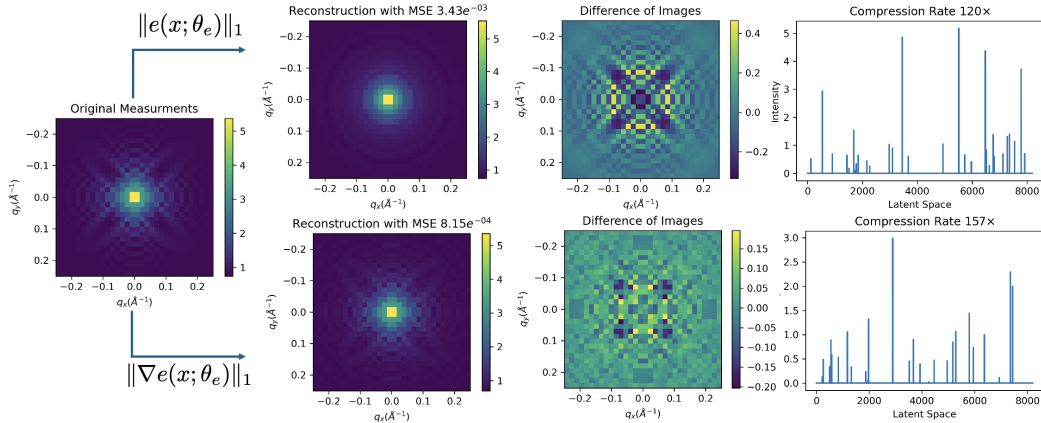


Figure 3: We show a representative testing input x_j for the SAS application in the first column, while its reconstructions using sparse autoencoder networks $(1, 2, 10, z)$ on the top and $(1, 2, 10, \nabla z)$ on the bottom are presented in the second column. The reconstruction errors $\|d(e(x_j; \theta_e); \theta_d) - x\|_2$ are 3.43×10^{-3} and 8.15×10^{-4} , respectively. We show the latent space variable $z_j = e(x_j; \theta_e)$ in the third column. The latent variable z_j each contains 34 and 26 non-zero elements. With an original image size of 64×64 the resulting compression ratio $\|x_j\|_0$ to $\|e(x_j; \theta_e)\|_0$ and $\|x_j\|_0$ to $\|\nabla e(x_j; \theta_e)\|_0$ are $120 : 1$ and $157 : 1$.

the latent weights $\tilde{w} = Q(w)$, such as using 16-bit floating-points, [19, 37, 42]. This provides for z the compressed representation (c, \tilde{w}) . These methods provide further ways to compress the data in addition to the sparsity.

3 Numerical Investigations

We illustrate the significant advantages our novel approach carries on simulated small angle scattering (SAS) data, a technique that is ubiquitous across the world’s X-ray light and neutron facilities. We utilize the tool SASView [43], a community-based tool used at the experimental facilities to analyze and simulate SAS experiments. For all SAS experiments simulated, we set the number of sensors to be uniformly spaced with $n = 64 \times 64$. Measurement sensors for SAS experiments are some form of charge-coupled device (CCD), so uniform spacing is common. All networks in this section use the same autoencoder architecture depicted in Figure 2 and this architecture follows that depicted in Figure 1. For a given input of size n and loss defined in Equation (3), we adopt the notation $\left(\frac{m}{n}, \frac{\ell}{n}, \lambda, f(z)\right)$ to uniquely define all networks used in this paper. We further report that all networks were trained using 1,000 epochs with a batch size of 512 on Oak Ridge National Laboratory’s (ORNL) Compute and Data Environment for Science (CADES) cluster [12]. In the initial phases of our investigation, we used convolution-type architectures which generally use convolution operators at the beginning and end of the network. We observed that the filtering aspect of the convolutional neural network dominated and hindered all information from reaching the latent space layer that connects the encoder and decoder.

We demonstrate for realistic configurations of SASView [20], we are able to highly compress all simulations from this package. We begin by randomly generating 50,000 images using the aforementioned sensor configuration of $n = 64 \times 64$, which is characteristic of the range of SAS experimental data collected at scattering facilities. The representative result of this investigation is presented in Figure 3, which demonstrates high compression rates with high accuracy for the networks $(1, 2, 10, z)$ and $(1, 2, 10, \nabla z)$. The tested data averaged a relative reconstruction error of $7.75 \cdot 10^{-2}\%$, and an averaged compression rate of $525 \times$ with minimum rate $205 \times$ and maximum rate $1365 \times$. For comparison, we tested a similar fully connected encoder-decoder, with the same number of parameters, but had the typical hourglass framework without the sparsity promoting L^1 norm in (3). With the same training and testing procedures, we were only able to obtain an average

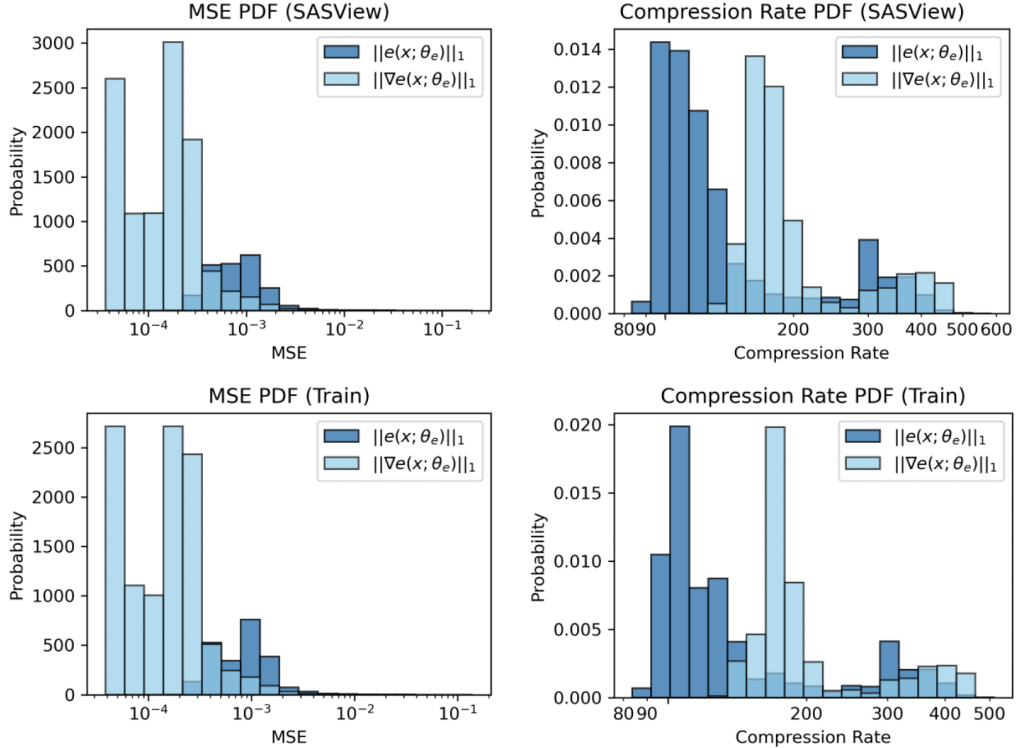


Figure 4: We show the error and compression rate partial distribution functions (PDFs) for the sparse autoencoder networks $(1, 2, 10, z)$ on the top and $(1, 2, 10, \nabla z)$ on the bottom. The top left plot displays the distributions of the entire training dataset, which consists of 50,000 random SAS images generated by SASView. The bottom left plot displays the distributions in the prediction of 150,000 independently random SAS images generated by SASView post-training. The mean training errors of both approaches are 2.48×10^{-3} and 8.00×10^{-4} , respectively. Correspondingly, the mean training compression rates are $153\times$ and $216\times$. Note that these values do only alter insignificantly for the testing set.

of $8\times$ compression rates with comparable accuracy. Using hyperparameter tuning, we found with the standard contracting encoder-decoder architecture the best size for the latent space was $\ell = 650$. This brings to light a significant difference with our sparse encoder-decoder, whereas the standard encoder-decoder requires hyperparameter tuning of the network architecture that results in training many different networks for set targets of accuracy and sparsity. Our sparse encoder-decoder is able to dynamically control the balance between sparsity and accuracy during training by adaptively changing the sparsity enforcing parameter λ in (3) on a single network architecture.

We demonstrate that we have captured all realistic configurations of SASView [20], using our train networks $(1, 2, 10, z)$ and $(1, 2, 10, \nabla z)$ in Figure 4. Here we sample again a much denser set of measurements using 150,000 images for testing. The takeaway from this analysis is that we can maintain the same level of compression and accuracy for both testing and training data. Additionally, it is demonstrated for the same number of network parameters, the sparsity promoting function $f(z) = \nabla z$ significantly improved compression rates and accuracy. Again this is visually represented in Figure 3 where this increased accuracy is able to maintain a more complex scattering pattern that can occur in SAS experiments.

Our methods can also be combined with further lossy and lossless methods to obtain further compression. As discussed above, we represent the information in z as a list of differences in the m indices of non-zero entries to obtain the sequence $\delta_1, \dots, \delta_{m-1}, \ell$ and the weights w_k at these indices to obtain $z \rightarrow (\delta, w)$. Here, we use arithmetic coding \mathcal{A} to develop for δ lossless compression methods

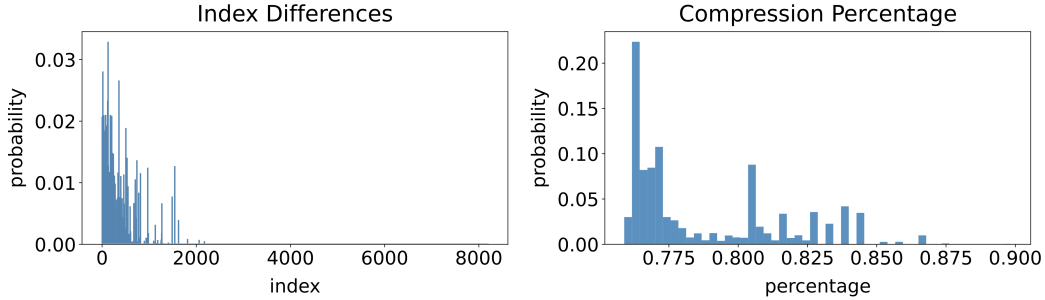


Figure 5: For further compression of z with our arithmetic entropy encoding, we show the distribution of index differences $\rho(\delta_k)$ for our representation $z \rightarrow (\delta, w)$ (left). For the SAS scattering data, we show the further compression reductions in percentage obtained for the index differences δ (right).

$c = \mathcal{A}(\delta; \rho(\cdot))$ [27, 44, 53]. We also can quantize \mathcal{Q} for lossy compression of w as $\tilde{w} = \mathcal{Q}(w)$, such as using lower-precision floating-points [19, 37]. For δ , we leverage that the probability distribution $\rho(\delta)$ will tend to skew to the left, for the SAS data see Figure 5. As an initial model for this distribution, we use a Gaussian-like form $\rho(\delta) = q(\delta)/Z$, where $Z = \sum_{\delta} q(\delta)$ where $q(\delta)$ is normally distributed with density $P(\delta; 0, \sigma^2)$ where $P(x; \mu, \sigma^2) = (2\pi\sigma^2)^{-1/2} \exp(-(x - \mu)^2/2\sigma^2)$. To help ensure efficient encoding on future samples we used conservative parameters $\sigma^2 = 10^3$ and $c_0 = 10^{-3}$. We found the lossless compression methods provide on average a compressed representation 79% of the uncompressed δ . Combining this with lossy quantization of the weights from 64-bit floating-points to 16-bit floating-points [19, 37, 42], yields an overall compressed representation (c, \tilde{w}) that is 52% of the uncompressed case. These methods provide an additional factor of around $2\times$ to the already favorable compression ratios achieved by the sparsity.

With respect to large scatter experiments across the world, there is starting to be a paradigm shift in how analysis of these measurements is being performed. Traditionally, analysis has been performed one experiment at a time on local clusters. However, as the datasets grow at these institutes, more complex analysis is needed at high-performance computing facilities that are not geographically collocated with experimental facilities. This presents a new challenge for scientific data reduction that requires guarantees of the reconstruction accuracy. Further, there is a goal of reducing analysis time, by doing this analysis in the compressed latent space. For this reason, we further investigate the properties of our sparse autoencoder using the well-studied MNIST (Modified National Institute of Standards and Technology database) database [14].

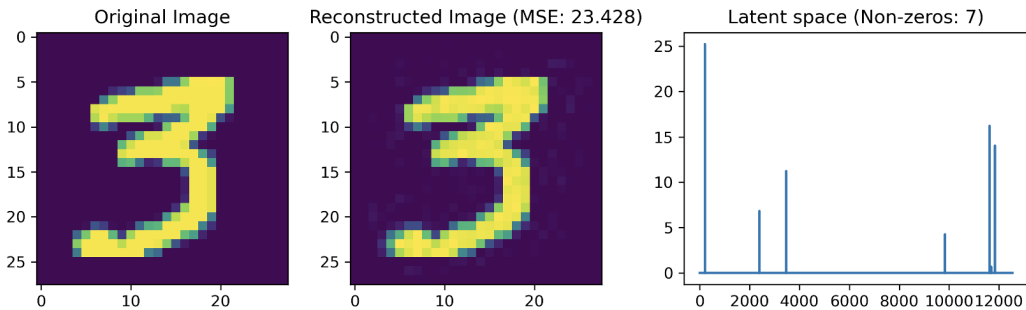


Figure 6: We show a representative testing input x_j for the MNIST dataset in the first column at a compression with an MSE of 23.428 and a compression ratio of 120 : 1.

We first demonstrate the high level of compression possible using a sparse network $(8, 16, 5000, z)$ using the given 60,000 testing and 10,000 handwritten images of size $n = 28 \times 28$. The purpose of this shift of dataset attention is to determine if the latent space can be used to maintain classification

accuracy and how latent space size affects this classification problem. Here, we are not concerned with the accuracy of any classification procedure per se, but rather the comparison of classification rates as a function of compression. For this reason, we consider the classification procedure of k-nearest neighbors (KNN) using the Cosine similarity measure. Table 1 gives a synopsis of the results.

Architecture $(\frac{m}{n}, \frac{\ell}{n}, \lambda, z)$	L^2		L^1		L^0		KNN
	Test	Train	Test	Train	Test	Train	
$(1/2, 1/4, 0, z)$	$1.9(2) \times 10^2$	$1.8(2) \times 10^2$	$1.8(1) \times 10^2$	$1.8(1) \times 10^2$	$1.57(1) \times 10^2$	$1.57(1) \times 10^2$	$9.54(2) \times 10^{-1}$
$(3/4, 1/2, 50, z)$	$1.6(1) \times 10^2$	$1.3(1) \times 10^2$	1.2037(6)	1.196(5)	$1.59(5) \times 10^2$	$1.57(5) \times 10^2$	$9.42(1) \times 10^{-1}$
$(11/2, 2, 200, z)$	$2.2(2) \times 10^2$	$1.6(2) \times 10^2$	$1.34(2) \times 10^{-1}$	$1.34(2) \times 10^{-1}$	$1.07(5) \times 10^2$	$1.06(5) \times 10^2$	$9.44(1) \times 10^{-1}$
$(2, 4, 800, z)$	$4(1) \times 10^2$	$2.6(9) \times 10^2$	$8.8(3) \times 10^{-2}$	$8.8(2) \times 10^{-2}$	$7(4) \times 10^1$	$7(4) \times 10^1$	$9.51(6) \times 10^{-1}$
$(4, 8, 3200, z)$	$3(1) \times 10^2$	$2.5(8) \times 10^2$	$1.9(1) \times 10^{-2}$	$1.914(1) \times 10^{-2}$	$5.12(9) \times 10^1$	$5(1) \times 10^1$	$9.57(5) \times 10^{-1}$

Table 1: This table displays the average L^2 norm for the difference of the input image and output image, $\|d(e(x_j; \theta_e); \theta_d) - x_j\|_2$, with the L^1 and L^0 norm for the latent space, $\|f(e(x_j; \theta_e))\|_1$ and $\|f(e(x_j; \theta_e))\|_0$, for different autoencoder architectures. The statistical notation, for example $1.57(1) \times 10^2$, implies 157 ± 1 . The uncertainty was determined by training five autoencoders on a random sample of 20% on all of the MNIST dataset (combining the 60,000 training and 10,000 testing into 70,000 samples that were randomly split into 14,000 training and 56,000 testing). To measure the informational content of the latent space we report the KNN accuracy by training on the 14,000 set latent space representation and testing on the derived 56,000 latent space. For reference the KNN accuracy in the image space is $9.60(9) \times 10^{-1}$.

We train multiple different sparse autoencoders on a random sample of 20% on all of the MNIST dataset. We know for KNN on the full MNIST training set that KNN accuracy in the image domain is 97%. Multiple training and test of the KNN accuracy in the image space with a 20% - 80% split on MNIST dataset produces a benchmark distribution of $9.60(9) \times 10^{-1}$. Table 1 shows the reconstruction error and latent space compression rate for latent spaces that range from $\frac{\ell}{n} \in [\frac{1}{4} - 8]$. For reference the network $(1/2, 1/4, 0, z)$ has no sparsity enforcement in the latent space, but rather sparsity is enforced through the hourglass architecture. The sparsity parameter λ for every other row was set to approximate the reconstruction accuracy achieved by $(1/2, 1/4, 0, z)$. The KNN result was trained on by the latent space each sparse autoencoder trained on and tested on the derived testing latent space produced by the same sparse autoencoder. The takeaway from this analysis is that the compression rate is increasing as a function of latent space size. Accuracy of reconstruction is maintained for each sparse autoencoder on data that the network never trained upon, and finally KNN prediction accuracy in the latent space increases with latent space size. The implications of this to analysis on scientific data reduction are positive in the sense that given a small random sample of a scientific dataset, sparse autoencoders with large latent spaces can maximize compression rates and this compressed representation provides equivalent information for analysis in the latent space in comparison to analysis in the data space.

4 Conclusion and Future Work

Our introduced sparse autoencoder methods provide natural extensions of compressed sensing approaches for use in the lossy compression of scientific data. Our introduced sparsity-promoting regularizations on the mappings of the latent variable were demonstrated to provide significant benefits for encoding scientific short-angle scattering data. Our numerical investigations indicate that the use of information-rich large dimensional latent spaces provides significant advantages in preserving features of signals during compression. Our methods provide ways for obtaining sparse representations by separating the structure of the encoded signals from the latent variable representations. Our methods introduce robust learning strategies enabling significant compression ratios allowing for the accurate and efficient storage, transmission, and analysis of scientific datasets. Our introduced methods also can be combined with other autoencoder strategies and lossy/lossless compression methods for handling diverse types of data in scientific applications.

Acknowledgement: This work was partially supported by the National Science Foundation (NSF) under grant DMS-2152661 (M. Chung). Author P.J.A. would like to acknowledge support from NSF Grant DMS-2306101. Author R.A. would like to acknowledge support the Office of Advanced Scientific Computing Research and performed at the Oak Ridge National Laboratory, which is managed by UT-Battelle, LLC for the US Department of Energy under Contract No. DE-AC05-00OR22725. Research used resources of the Oak Ridge Leadership Computing Facility at Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725

Bibliography

- [1] B. M. Afkham, J. Chung, and M. Chung. “Goal-oriented Uncertainty Quantification for Inverse Problems via Variational Encoder-Decoder Networks”. In: *arXiv preprint arXiv:2304.08324* (2023) (cited on page 2).
- [2] C. Bonneville et al. “A Comprehensive Review of Latent Space Dynamics Identification Algorithms for Intrusive and Non-Intrusive Reduced-Order-Modeling”. In: *arXiv preprint arXiv:2403.10748* (2024) (cited on page 2).
- [3] H. Bourlard and Y. Kamp. “Auto-association by multilayer perceptrons and singular value decomposition”. In: *Biological cybernetics* 59.4 (1988), pages 291–294 (cited on page 1).
- [4] T. Bricken et al. “Towards monosemanticity: Decomposing language models with dictionary learning”. In: *Transformer Circuits Thread* (2023), page 2 (cited on page 2).
- [5] E. J. Candes, J. K. Romberg, and T. Tao. “Stable signal recovery from incomplete and inaccurate measurements”. In: *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences* 59.8 (2006), pages 1207–1223 (cited on page 3).
- [6] E. J. Candes and T. Tao. “Decoding by linear programming”. In: *IEEE transactions on information theory* 51.12 (2005), pages 4203–4215 (cited on page 3).
- [7] E. J. Candès, J. Romberg, and T. Tao. “Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information”. In: *IEEE Transactions on information theory* 52.2 (2006), pages 489–509 (cited on page 4).
- [8] X. Chen et al. “A sub-Nyquist rate sampling receiver exploiting compressive sensing”. In: *IEEE Transactions on Circuits and Systems I: Regular Papers* 58.3 (2010), pages 507–520 (cited on page 4).
- [9] N. Cheng et al. “Bi-fidelity variational auto-encoder for uncertainty quantification”. In: *Computer Methods in Applied Mechanics and Engineering* 421 (2024), page 116793 (cited on page 2).
- [10] Z. Cheng et al. “Deep convolutional autoencoder-based lossy image compression”. In: *2018 Picture Coding Symposium (PCS)*. IEEE, 2018, pages 253–257 (cited on page 1).
- [11] M. Chung and R. A. Renaut. “A variable projection method for large-scale inverse problems with ℓ^1 regularization”. In: *Applied Numerical Mathematics* 192 (2023), pages 297–318 (cited on page 3).
- [12] T. Compute and D. E. for Science (CADES). *www.cades.ornl.gov*. 2024 (cited on page 5).
- [13] T. M. Cover and J. A. Thomas. *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)*. USA: Wiley-Interscience, 2006. ISBN: 0471241954. URL: <http://dx.doi.org/10.1002/047174882X> (cited on pages 3, 4).
- [14] L. Deng. “The mnist database of handwritten digit images for machine learning research”. In: *IEEE Signal Processing Magazine* 29.6 (2012), pages 141–142 (cited on page 7).
- [15] D. L. Donoho. “For most large underdetermined systems of linear equations the minimal ℓ^1 -norm solution is also the sparsest solution”. In: *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences* 59.6 (2006), pages 797–829 (cited on page 3).
- [16] B. Dumitrescu and P. Irofti. *Dictionary learning algorithms and applications*. Springer, 2018 (cited on page 4).
- [17] M. Elad. *Sparse and redundant representations: from theory to applications in signal and image processing*. Volume 2. Springer, 2010 (cited on page 3).
- [18] I. Goodfellow, Y. Bengio, and A. Courville. *Deep learning*. MIT press, 2016 (cited on page 1).

- [19] R. M. Gray and D. L. Neuhoff. “Quantization”. In: *IEEE transactions on information theory* 44.6 (1998), pages 2325–2383 (cited on pages 4, 5, 7).
- [20] W. T. Heller, M. Doucet, and R. K. Archibald. “Sas-temper: Software for fitting small-angle scattering data that provides automated reproducibility characterization”. In: *SoftwareX* 16 (2021), page 100849. ISSN: 2352-7110. DOI: <https://doi.org/10.1016/j.softx.2021.100849>. URL: <https://www.sciencedirect.com/science/article/pii/S235271102100128X> (cited on pages 5, 6).
- [21] X. Jiang et al. “A novel sparse auto-encoder for deep unsupervised learning”. In: *2013 Sixth international conference on advanced computational intelligence (ICACI)*. IEEE, 2013, pages 256–261 (cited on page 2).
- [22] S. H. Kabil and H. Bourlard. “From Undercomplete to Sparse Overcomplete Autoencoders to Improve LF-MMI Speech Recognition”. In: *Interspeech 2022* (2022), pages 1061–1065 (cited on page 2).
- [23] D. P. Kingma and M. Welling. “Auto-encoding variational Bayes”. In: *arXiv preprint arXiv:1312.6114* (2013) (cited on page 1).
- [24] M. A. Kramer. “Autoassociative neural networks”. In: *Computers & chemical engineering* 16.4 (1992), pages 313–328 (cited on page 1).
- [25] K. Kreutz-Delgado et al. “Dictionary learning algorithms for sparse representation”. In: *Neural computation* 15.2 (2003), pages 349–396 (cited on page 4).
- [26] R. Kumar et al. “High-fidelity audio compression with improved RVQGAN”. In: *Advances in Neural Information Processing Systems* 36 (2024) (cited on page 1).
- [27] G. G. Langdon. “An introduction to arithmetic coding”. In: *IBM Journal of Research and Development* 28.2 (1984), pages 135–149 (cited on pages 4, 6).
- [28] D. Lappas, V. Argyriou, and D. Makris. “Fourier transformation autoencoders for anomaly detection”. In: *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pages 1475–1479 (cited on page 2).
- [29] J. Lee, A. Rangarajan, and S. Ranka. “Nonlinear-by-Linear: Guaranteeing Error Bounds in Compressive Autoencoders”. In: *Proceedings of the 2023 Fifteenth International Conference on Contemporary Computing*. 2023, pages 552–561 (cited on page 2).
- [30] Q. Li et al. “Deep sparse autoencoder and recursive neural network for EEG emotion recognition”. In: *Entropy* 24.9 (2022), page 1187 (cited on page 2).
- [31] J. Liu et al. “Exploring autoencoder-based error-bounded compression for scientific data”. In: *2021 IEEE International Conference on Cluster Computing (CLUSTER)*. IEEE, 2021, pages 294–306 (cited on page 1).
- [32] C. Louizos, M. Welling, and D. P. Kingma. “Learning sparse neural networks through L_0 regularization”. In: *arXiv preprint arXiv:1712.01312* (2017) (cited on page 2).
- [33] A. Majumdar. “An autoencoder based formulation for compressed sensing reconstruction”. In: *Magnetic resonance imaging* 52 (2018), pages 62–68 (cited on page 2).
- [34] A. Makhzani and B. Frey. “K-sparse autoencoders”. In: *arXiv preprint arXiv:1312.5663* (2013) (cited on page 2).
- [35] G. Martino, D. Moroni, and M. Martinelli. “Are We Using Autoencoders in a Wrong Way?” In: *arXiv preprint arXiv:2309.01532* (2023) (cited on pages 2, 3).
- [36] Q. Meng et al. “Relational autoencoder for feature extraction”. In: *2017 International joint conference on neural networks (IJCNN)*. IEEE, 2017, pages 364–371 (cited on page 2).
- [37] J.-M. Muller et al. *Handbook of floating-point arithmetic*. Springer, 2018 (cited on pages 5, 7).
- [38] E. Newman, J. M. Solomon, and M. Chung. “Image reconstructions using sparse dictionary representations and implicit, non-negative mappings”. In: *arXiv preprint arXiv:2312.03180* (2023) (cited on page 4).
- [39] A. Ng et al. “Sparse autoencoder”. In: *CS294A Lecture notes* 72.2011 (2011), pages 1–19 (cited on page 2).
- [40] H.-A. T. Nguyen, T. H. Le, and T. D. Bui. “A deep wavelet sparse autoencoder method for online and automatic electrooculographical artifact removal”. In: *Neural Computing and Applications* 32.24 (2020), pages 18255–18270 (cited on page 2).
- [41] E. Plaut. “From principal subspaces to principal components with linear autoencoders”. In: *arXiv preprint arXiv:1804.10253* (2018) (cited on page 1).

- [42] A. Polino, R. Pascanu, and D. Alistarh. “Model compression via distillation and quantization”. In: *arXiv preprint arXiv:1802.05668* (2018) (cited on pages 5, 7).
- [43] S. project. *sasview.org*. 2024 (cited on page 5).
- [44] J. Rissanen and G. G. Langdon. “Arithmetic coding”. In: *IBM Journal of research and development* 23.2 (1979), pages 149–162 (cited on pages 4, 6).
- [45] M. Sakurada and T. Yairi. “Anomaly detection using autoencoders with nonlinear dimensionality reduction”. In: *Proceedings of the MLSDA 2014 2nd workshop on machine learning for sensory data analysis*. 2014, pages 4–11 (cited on page 1).
- [46] S. Scardapane et al. “Group sparse regularization for deep neural networks”. In: *Neurocomputing* 241 (2017), pages 81–89 (cited on page 2).
- [47] S. Tariyal et al. “Deep dictionary learning”. In: *IEEE Access* 4 (2016), pages 10096–10109 (cited on page 2).
- [48] L. Theis et al. “Lossy image compression with compressive autoencoders”. In: *International conference on learning representations*. 2022 (cited on page 1).
- [49] R. Tibshirani. “Regression shrinkage and selection via the lasso”. In: *Journal of the Royal Statistical Society Series B: Statistical Methodology* 58.1 (1996), pages 267–288 (cited on page 3).
- [50] I. Tošić and P. Frossard. “Dictionary learning”. In: *IEEE Signal Processing Magazine* 28.2 (2011), pages 27–38 (cited on page 4).
- [51] P. Vincent et al. “Extracting and composing robust features with denoising autoencoders”. In: *Proceedings of the 25th international conference on Machine learning*. 2008, pages 1096–1103 (cited on page 1).
- [52] Y. Wang, H. Yao, and S. Zhao. “Auto-encoder based dimensionality reduction”. In: *Neurocomputing* 184 (2016), pages 232–242 (cited on page 2).
- [53] I. H. Witten, R. M. Neal, and J. G. Cleary. “Arithmetic coding for data compression”. In: *Communications of the ACM* 30.6 (1987), pages 520–540 (cited on pages 4, 6).
- [54] P. Wu et al. “Reduced order model using convolutional auto-encoder with self-attention”. In: *Physics of Fluids* 33.7 (2021) (cited on page 2).
- [55] G. Zhang, Y. Liu, and X. Jin. “A survey of autoencoder-based recommender systems”. In: *Frontiers of Computer Science* 14 (2020), pages 430–450 (cited on page 1).
- [56] Y. Zhu et al. “Physics-constrained deep learning for high-dimensional surrogate modeling and uncertainty quantification without labeled data”. In: *Journal of Computational Physics* 394 (2019), pages 56–81 (cited on page 2).