

How Learning by Reconstruction Produces Uninformative Features for Perception

Randall Balestriero¹ Yann LeCun²

Abstract

Input space reconstruction is an attractive representation learning paradigm. Despite interpretability benefit of reconstruction and generation, we identify a misalignment between learning to reconstruct, and learning for perception. We show that the former allocates a model’s capacity towards a subspace of the data explaining the observed variance—a subspace with uninformative features for the latter. For example, the supervised TinyImagenet task with images projected onto the top subspace explaining 90% of the pixel variance can be solved with 45% test accuracy. Using the bottom subspace instead, accounting for only 20% of the pixel variance, reaches 55% test accuracy. Learning by reconstruction is also wasteful as the features for perception are learned last, pushing the need for long training schedules. We finally prove that learning by denoising can alleviate that misalignment for some noise strategies, e.g., masking. While tuning the noise strategy without knowledge of the perception task seems challenging, we provide a solution to detect if a noise strategy is never beneficial regardless of the perception task, e.g., additive Gaussian noise.

1. Introduction

One of the far reaching mandate of deep learning is to provide a self-contained methodology to learn intelligible and universal representations of data (LeCun et al., 2015). That is, to learn a nonlinear transformation of the data producing a parsimonious representation able to solve numerous downstream tasks. Significant progress has been made through the lens of supervised learning, i.e., by learning a representation that maps the observed data to provided labels of an

¹Brown University ²NYU. Correspondence to: Randall Balestriero <rbalestr@brown.edu>.

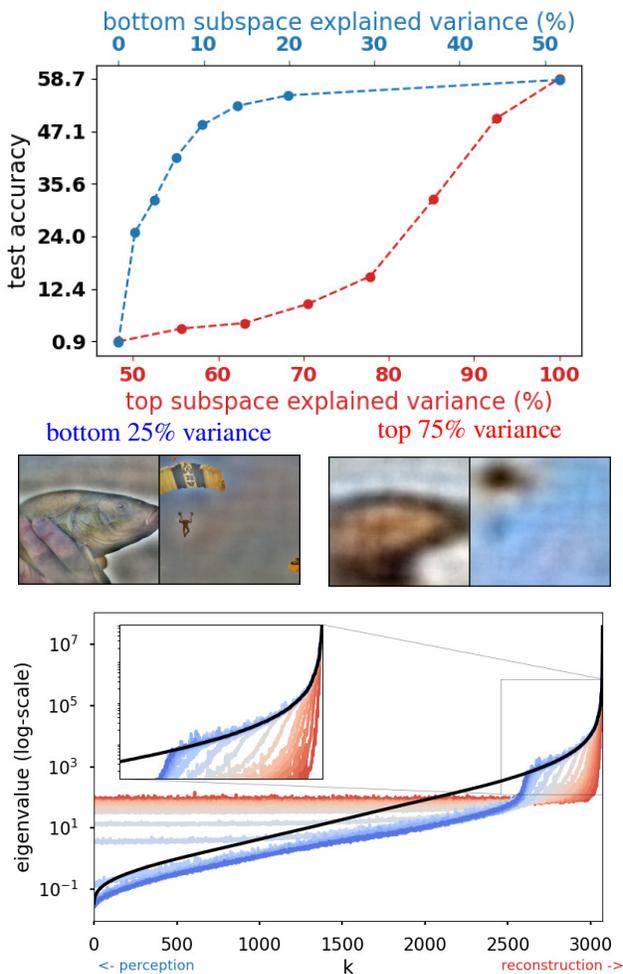


Figure 1. Features for reconstruction are uninformative for perception (top): TinyImagenet ResNet9 top-1 accuracy when trained and validated on images projected on the top-subspace (red) or bottom subspace (blue) of explained variance, corresponding images displayed in the **middle** and in Fig. 10. **Perception features are learned last (bottom):** training loss evolution (red to blue) of reconstructed training images from a deep Autoencoder projected onto the eigenspace of the original data (black). The top eigenspace (right) is learned first, and then, if training lasts long enough, the features most useful for perception (left) are finally learned. This explains why learning by performances on perception task keep increasing long after reconstructed samples look appealing.

priori known downstream task (Krizhevsky et al., 2012). Labels being costly and over-specialized, much progress was also made through the lens of unsupervised learning (Barlow, 1989; Ghahramani, 2003) roughly falling into three camps. First, reconstruction-based methods that produce (compressed) latent representations of the data nonetheless sufficient to recover most of the original data, e.g., Denoising/Variational/Masked Auto-Encoders (Vincent et al., 2010; Kingma & Welling, 2013; He et al., 2022) and Autoregressive models (Van den Oord et al., 2016; Chen et al., 2018). Second, score matching which is often solved by setting up a surrogate supervised task of classifying observed samples from noise (Hyvärinen & Dayan, 2005). Third, Self-Supervised Learning (SSL) (Chen et al., 2020; Zbontar et al., 2021; Bardes et al., 2021; Balestrierio et al., 2023)—whose contrastive methods can be thought of as generalized versions of score matching and Noise Contrastive Estimation (Gutmann & Hyvärinen, 2010)—combining an invariance term bringing together representations of data known to be semantically similar, or generated to be so, and an anti-collapse term making sure that not all representations become similar.

In recent years, SSL emerged as the preferred solution—regularly reaching new state-of-the-arts through careful experimental design (Garrido et al., 2023). Yet, reconstruction-based methods maintain a large presence due to their ability to provide reconstructed samples that are human interpretable, enabling informed quality assessment of a model (Selvaraju et al., 2016; Bordes et al., 2021; Brynjolfsson et al., 2023; Baidoo-Anu & Ansah, 2023). Despite that benefit, reconstruction-based learning falls behind SSL as it requires fine-tuning to reach the state-of-the-art. One of the most popular reconstruction-based learning strategy that emerged as the solution of choice in recent years is the Masked-Autoencoder (He et al., 2022).

We ask the following question: *Why reconstruction based learning easily produces compelling reconstructed samples but fails to provide competitive latent representations for perception?* We can pinpoint that observation to at least three reasons.

R1: Misaligned. The features with most reconstructive power are the least informative for perceptual tasks as depicted in Fig. 1 (top). Instead, images projected onto the subspace least useful for reconstruction still contain sufficient information for perception.

R2: Ill-conditioned. From the misalignment, the features useful for perception are learned last for a reconstruction task as depicted in Fig. 1 (bottom), meaning that learning by reconstruction is wasteful as it requires long training schedules to capture informative features for perception.

R3: Ill-posed. Due to the “irrelevance” of perception fea-

tures to affect reconstruction error, learning by reconstruction can produce drastically different representation quality between training settings. For example, we can find a same model with two different set of parameters that produce the same train and test reconstruction error while exhibiting significant performance gaps for perceptual tasks as depicted in Fig. 7, where for a given reconstruction error the top-1 accuracy on Imagenet-10 can vary from 50% to almost 90%.

These findings provide first clues as to why learning by reconstruction requires long training time and fine-tuning—observations that are common knowledge within the representation learning community but still lack theoretical justification. Those findings alone do not answer the following question: *Why Masked Autoencoders were able to provide a significant improvement in the quality of the learned representation for solving perception tasks?* We will prove that some denoising tasks can push back some of the limitations of learning by reconstructions. In particular, we will demonstrate that masking is provably beneficial while other noise distributions, e.g., additive Gaussian noise, are not. We summarize our technical contributions below

- We provide a novel measure of alignment between reconstruction and classification task from first principles (Section 3.1, Eq. (7), and Fig. 2)
- We identify the root cause for reconstruction and perceptual classification tasks to be misaligned in term of the spectral property of the data covariance matrix (Section 3.2 and Fig. 4). The misalignment increases with presence of background and increased number of objects in the images
- We prove that denoising tasks introduce a mechanism through the noise distribution to improve alignment between reconstruction and perception: masking provides great alignment benefits while additive Gaussian noise is provably ineffective(Fig. 6 and Corollary 3.1)

We hope that our findings will help skew further research in learning by reconstruction to explore alternative noise distributions as they are the main driver of learning useful representations for perception. Codebase provided at github.com/RandallBalestrierio/LearningByReconstruction.

2. Background and Notations

Notations: We denote by $x_n \in \mathbb{R}^D, n = 1, \dots, N$ the n^{th} input sample, e.g., a (H, W, C) image flattened to a $D = H \times W \times C$ dimensional vector. The entire training set is collected into the matrix $\mathbf{X} \triangleq [x_1, \dots, x_N] \in \mathcal{M}_{D,N}(\mathbb{R})$ where $\mathcal{M}_{n,m}(\mathbb{R})$ is the vector space of real $n \times m$ matrices. Throughout our study, vectors will always

be column-vectors, and matrices are built by horizontally stacking column-vectors, i.e., they are column major. Unless specified otherwise, we assume that the input matrix \mathbf{X} is full-rank. In practice, if this is not the case, one can easily disregard the subspace associated with 0 singular values and apply our analysis on that filtered matrix instead.

Learning by reconstruction. Learning representations by fitting a model’s parameters θ to produce a reconstruction of presented inputs as in

$$\min_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{\mathbf{x}}} [d(g_{\theta}(f_{\theta}(\mathbf{x})), \mathbf{x})], \quad (1)$$

is common (Bottou, 2012; Kingma & Ba, 2014; LeCun et al., 2015), given some distance measure d . The reconstruction provides a qualitative signal enabling one to easily assess the quality of the model, and even interpret trained classification models with g_{θ} being trained and f_{θ} the (frozen) pretrained classifier (Zeiler & Fergus, 2014; Mahendran & Vedaldi, 2015; Olah et al., 2017; Shen et al., 2020). In its simplest form, the encoder $f_{\theta} : \mathbb{R}^D \mapsto \mathbb{R}^K$ and the decoder $g_{\theta} : \mathbb{R}^K \mapsto \mathbb{R}^D$ are linear, possibly with shared parameters. In such settings, the optimal parameters are obtained from Principal Component Analysis (Wold et al., 1987). Many variants of Eq. (1) have emerged, such as denoising and masked autoencoders (MAEs)(Vincent et al., 2010; He et al., 2022). The objective remains similar: learn a low-dimensional latent embedding of the data that is able to reconstruct the original samples while being robust to some noise perturbation added onto the samples as in

$$\min_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{\mathbf{x}}} \left[\mathbb{E}_{\mathbf{x}' \sim p_{\mathbf{x}'|\mathbf{x}}} [d(g_{\theta}(f_{\theta}(\mathbf{x}')), \mathbf{x})] \right], \quad (2)$$

with $p_{\mathbf{x}'|\mathbf{x}}$ applying some (conditional) noise transformation to the original input, e.g., $\mathbf{x}' \sim \mathcal{N}(\mathbf{x}, \epsilon \mathbf{I})$, $\epsilon > 0$.

Known limitations. Learning by reconstruction is widely popular and thus heavily studied. Major axes of research evolve around (i) deriving novel loss functions for specific datasets that better align with semantic distance (Wang et al., 2004; Kulis et al., 2013; Ballé et al., 2016; Bao et al., 2017), (ii) explaining the learned embedding dimensions (Tran et al., 2017; Esmaeili et al., 2019; Mathieu et al., 2019), and (iii) imposing structure in the embedding space such as clustered embeddings (Jiang et al., 2016; Dilokthanakul et al., 2016; Lim et al., 2020; Seydoux et al., 2020). Despite the rich literature, the current solutions of choice to learn representation in computer vision still rely on the Mean Squared Error loss in pixel space with the possible application of structured noise ($p_{\mathbf{x}'|\mathbf{x}}$ in Eq. (2)), e.g., as employed by the current state-of-the-art solution of masked-autoencoders (MAEs) (He et al., 2022). Yet, even that solution learns a representation that needs to be fine-tuned to compete with state-of-the-art, e.g., MAE’s performances are drastically determined by two parameters (i) the training time, i.e.,

evaluation performance of the learned representation does not plateau even after 1600 epochs on Imagenet, and (ii) the need to fine tune, i.e., evaluation performance with and without finetuning have a significant gap going from 70% to 84% on top-1 Imagenet classification task.

Closer to our study, (Locatello et al., 2019) showed a first possible theoretical limitation of learning disentangled representations through reconstruction. That limitation is also reinforced by the large empirical study of (Zamir et al., 2018) concluding that learning by reconstruction and learning by denoising produce similar representations that are misaligned from supervised learning.

Our study will propose some first hints as to why even current state-of-the-art reconstruction-based solutions are poised with slow training, and the need to finetune (Section 3). We will conclude by proving that MAEs’s masking strategy partially alleviates those limitations (Section 4), showing that the most rewarding findings for learning by reconstruction may emerge from novel denoising strategies, hence bringing some nuance into the alignment results from (Zamir et al., 2018).

3. Rich Features for Reconstruction are Poor Features For Perception

This section provides the theoretical ground and empirical validation of **R1** and **R2** from Section 1, namely, that learning by reconstruction learns features that are misaligned with common perception tasks. We start by deriving a closed form alignment measure between those two tasks in Section 3.1 and conclude by empirically measuring that mismatch in Section 3.2.

3.1. How To Measure The Alignment Between Reconstruction and Supervised Tasks

As a starting point to our study, we will build intuition and obtain theoretical results in the linear regime. As we will see at the end of this Section 3.1, this seemingly simplified setting turns out to be informative of practical cases.

Let’s consider an encoder mapping $\mathbf{V} \in \mathcal{M}_{D,K}(\mathbb{R})$, a decoder mapping $\mathbf{Z} \in \mathcal{M}_{K,D}(\mathbb{R})$, and a predictor head $\mathbf{W} \in \mathcal{M}_{K,C}(\mathbb{R})$, where C is the number of target dimensions, or classes. The targets for $\mathbf{X} \in \mathcal{M}_{D,N}(\mathbb{R})$ are given by $\mathbf{Y} \in \mathcal{M}_{C,N}(\mathbb{R})$. The combination of the supervised and reconstruction losses is given by

$$\mathcal{L}(\mathbf{V}, \mathbf{W}, \mathbf{Z}) = \|\mathbf{W}^{\top} \mathbf{V}^{\top} \mathbf{X} - \mathbf{Y}\|_F^2 + \lambda \|\mathbf{Z}^{\top} \mathbf{V}^{\top} \mathbf{X} - \mathbf{X}\|_F^2, \quad (3)$$

where the latent representation, $\mathbf{V}^{\top} \mathbf{X}$, is shared between the two losses, $\lambda \geq 0$ controls the trade-off between the two terms. Quantifying how the optimal parameters of Eq. (3)

vary with λ will be key to assess how the two losses are aligned. As a starting point, let's formalize below the optimal parameters of this loss function using the notations $M = P_M D_M P_M^\top$ for the eigendecomposition of a symmetric semi-definite positive matrix M . We will also denote, to lighten notations, $A \triangleq X (Y^\top Y + \lambda X^\top X) X^\top$.

Theorem 1. *The loss function from Eq. (3) is minimized for*

$$V^* \text{ spans } P_{XX^\top} D_{XX^\top}^{-\frac{1}{2}} (P_H)_{:,1:K}, \quad (4)$$

$$W^* = (V^{*\top} X X^\top V^*)^{-1} V^{*\top} X Y^\top, \quad (5)$$

$$Z^* = (V^{*\top} X X^\top V^*)^{-1} V^{*\top} X X^\top, \quad (6)$$

where $H \triangleq D_{XX^\top}^{-\frac{1}{2}} P_{XX^\top}^\top A P_{XX^\top} D_{XX^\top}^{-\frac{1}{2}}$. (Proof in Section 5.1, empirical validation in Fig. 9.)

We observe that the optimal solutions from Theorem 1 continuously interpolates between the standard least square (OLS) problem ($\lambda = 0$) and the unsupervised linear auto-encoder or Principal Component Analysis (PCA) setting ($\lambda \rightarrow \infty$). We formalize below that we recover the optimal solutions for each of those extreme cases from Theorem 1.

Corollary 1.1. *The solution from Theorem 2 recovers the OLS solution for $W^{*\top} V^{*\top}$ as $\lambda \rightarrow 0$, and the PCA solution for $Z^{*\top} V^{*\top}$ as $\lambda \rightarrow \infty$. (Proof in Section 5.2.)*

That observation should comfort the reader that Eq. (3) accurately conveys both ends of the spectrum from supervised learning to reconstruction based learning, while continuously interpolating in-between.

Condition for perfect alignment. The result from Theorem 1 enables us to formalize the condition for perfect alignment between the two tasks, i.e., under which condition the solution V^* is not impacted by λ .

Proposition 1. *The supervised and reconstruction tasks are aligned (the optimal solutions do not depend on λ) iff the intersection of the top- K eigenspaces of $X^\top X$ and $Y^\top Y$ is of dimension K .*

In other words, whenever the condition in Proposition 1 holds, the matrix P_H (recall Theorem 1) will include the same eigenvectors (up to rotation) for any λ , making the optimal parameters (Eqs. (4) to (6)) independent of λ . In practice, we will see that Proposition 1's alignment condition is never fulfilled, pushing the need to define a more precise measure of alignment between the two tasks.

Continuous measure of alignment. As we aim to measure the tasks alignment more precisely that in a yes/no setting (Proposition 1), we propose the following continuous measure

$$\text{alignment}(k) \triangleq \frac{\|Y^\top Y (P_{XX^\top})_{1:k}\|_F^2}{\|Y^\top Y P_{XX^\top}\|_F^2}, \quad (7)$$

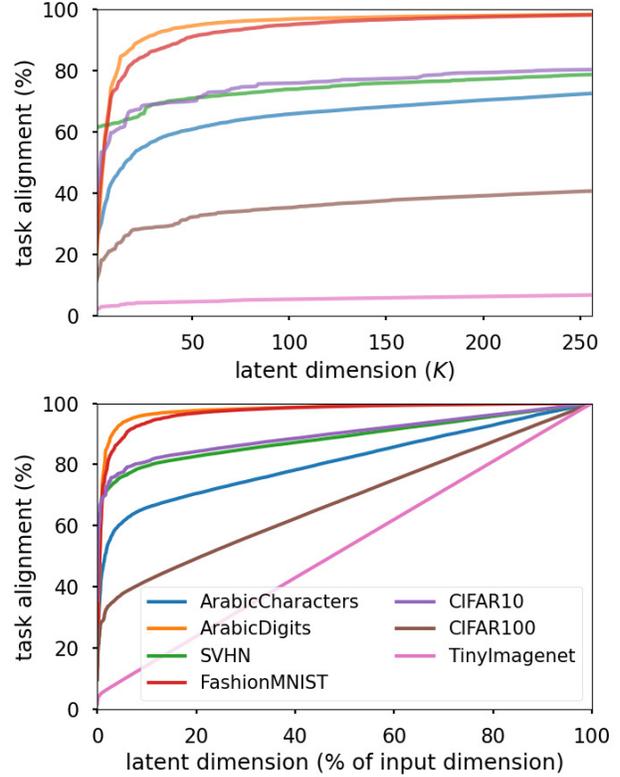


Figure 2. Depiction of the closed form alignment measure from Eq. (7) measuring the minimum supervised training error achievable given the optimal reconstruction parameters, as per Theorem 1 and Corollary 1.2. **Top:** depiction in terms of the latent dimension K . **Bottom:** depiction in terms of the ratio of the latent dimension K to the input dimension D . We clearly observe that as the dataset becomes more realistic (going from background-free images to CIFAR and TinyImagenet), as the alignment between the reconstruction and supervised task lessens. In particular, when going to TinyImagenet, we observe that the alignment only increases linearly with respect to the latent space dimension.

where $\|Y^\top Y P_{XX^\top}\|_F^2$ simplifies as $\|Y^\top Y\|_F^2$ when $D = N$. We assume that the supervised task can be at least partially solved from X , ensuring $\|Y^\top Y P_{XX^\top}\|_F^2 > \epsilon$. In words, Eq. (7) is the (scaled and negated) supervised training error achieved given the representation $(V^\top X)$ minimizing the reconstruction loss, which is measured by how much of the matrix $Y^\top Y$ can be reconstructed from the top- k subspace of $X^\top X$, as formalized below.

Corollary 1.2. *alignment(k) from Eq. (7) increases with k , has value 0 iff the two losses are misaligned, and has value 1 iff the two losses are aligned. (Proof in Section 5.3.)*

Although the measure from Eq. (7) is motivated from the linear setting of Theorem 1, we will demonstrate at the end of this Section 3.1 that it actually aligns well with practical

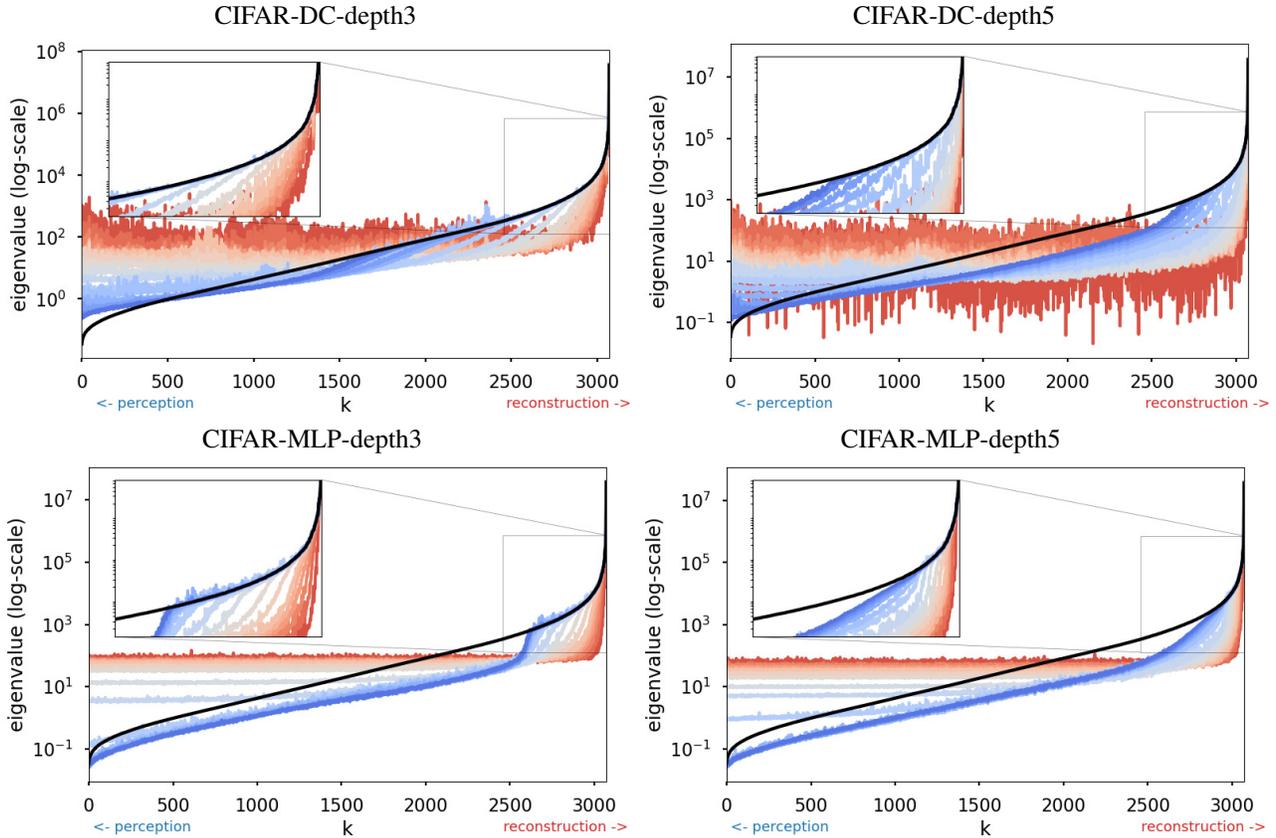


Figure 3. Reprise of Fig. 1 for additional autoencoder architectures: convolutional encoder and deconvolutional decoder (**top**) and MLP encoder and decoder (**bottom**). We clearly observe that the top subspace is learned first during training, which is the one that best minimize the reconstruction loss but that contains the least informative features for perception, as per Fig. 4.

settings.

Findings. We now propose to evaluate the closed form alignment metric (Eq. (7)) on a few datasets. Note that it can be implemented efficiently as detailed in Section 5.8. In Fig. 2, we measure the metric of Eq. (7) for a sweep of latent dimension K over 7 different datasets. We observe three striking regimes. First, for images without background, reconstruction and classification tasks are very much aligned even for small latent dimension, as low as 20% of the input dimension. Second, when comparing datasets with same image distributions but different number of classes (CIFAR10 to CIFAR100) the misalignment increases between the two task, especially for small embedding dimension. This follows our intuition that additional budget must be devoted to separate more classes. And as the subspace of the data used for reconstruction and classification do not align (recall Proposition 1), a greater misalignment is measured. And third, when looking at more realistic images (higher resolution and more diverse) such as tiny-imagenet, we observe that the alignment only increases linearly with the latent space dimension K , requiring $K = D$ in that case to ensure alignment. We thus conclude that *the presence*

of background, finer classification tasks, and higher resolution images, are all factors that drastically decrease the alignment between learning features for perception tasks and learning features that reconstruct those images.

Linear regime results are informative. We have focused on the linear encoder, decoder, and classification head of Eq. (3). Albeit insightful, one might wonder how much of those insights transfer to the more realistic setting of employing nonlinear mappings. We note that it remains common to keep the classification head linear therefore leading to the following generalization of Eq. (3) as

$$\mathcal{L}(\mathbf{W}, \theta, \gamma) = \|\mathbf{W}^\top f_\theta(\mathbf{X}) - \mathbf{Y}\|_F^2 + \lambda \|g_\gamma(f_\theta(\mathbf{X})) - \mathbf{X}\|_F^2. \quad (8)$$

The encoder is now the nonlinear mapping $f_\theta : \mathbb{R}^D \mapsto \mathbb{R}^K$ and the decoder is the nonlinear mapping $g_\gamma : \mathbb{R}^K \mapsto \mathbb{R}^D$. We formalize below a result that will reinforce the legitimacy of our linear regime analysis (Theorem 1, Proposition 1, and Corollary 1.2) by showing that it is (i) a correct model during the early phase of training, and even (ii) a correct model throughout training when the decoder being employed is under-parametrized.

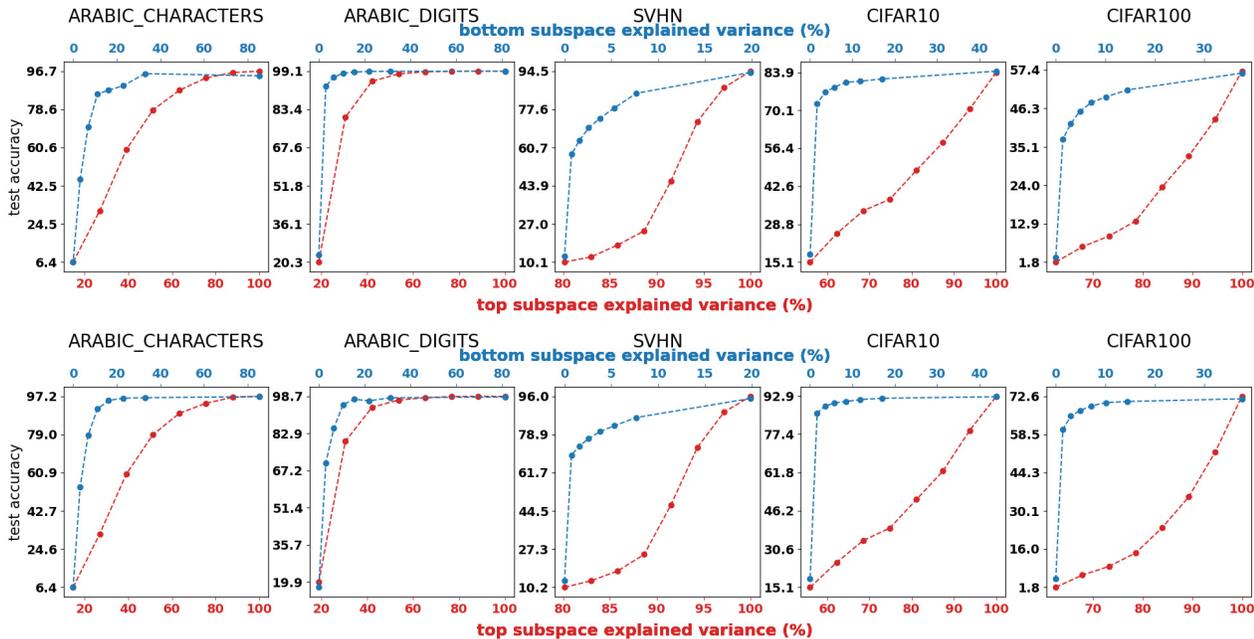


Figure 4. We depict the classification accuracy of a ResNet9 DNN when trained and tested on images that have been projected onto the top (red) and bottom (blue) subspace as ordered per the eigenvalues of the data covariance matrix, without data-augmentation (top) and with data-augmentation (bottom). We clearly observe that except for datasets without background and for which reconstruction and classification are better aligned (recall Fig. 2), the final performance is greater when employing the subspace of the data that explains the least the pixel variation, i.e., the bottom subspace. Training was done with cross-entropy loss, Fig. 8 provides the MSE case with one-hot labels as target showcasing the same trends.

Theorem 2. For any high-capacity encoder f_θ , studying Eq. (3) and Eq. (8) is equivalent at initialization for any decoder, and is always equivalent when the decoder is linear. (Proof in Section 5.6.)

Combining Theorems 1 and 2, we obtain that even with DNs, during the early stages of learning, the encoder-decoder mapping focuses on the principal subspace of the data, i.e., the space that explains most of the reconstruction error in the linear regime. As our study strongly hinges on that claim, we propose to empirically validate it in the following Section 3.2.

3.2. Reconstruction and Perception Features Live In Different Subspaces of the Data

We characterized in the previous Section 3.1 how classification and reconstruction tasks fail to align when it comes to learning common features. In particular, Section 3.1 and Fig. 2 validated how training focuses first on the top subspace of the data. We now reinforce our claim by showing that supervised tasks can not be solved when restricting the images on the subspace that is learned first by reconstruction.

Perception can not be solved from the principal subspace of the data. We first propose a controlled experiment where we artificially remove some of the original data subspace. In particular, we consider two settings. First, we gradually remove the subspace associated to the top eigenvectors of the data covariance matrix effectively removing what is most useful for reconstruction but also what we claim to be least useful for perception. Second, we gradually remove the subspace associated to the bottom eigenvectors (the one least useful for reconstruction but that we claim to be most useful for perception). This procedure is applied to the entire dataset (train and test images) before any DN training occurs. Hence the DN is only presented with the filtered images. We report the top-1 accuracy over numerous datasets in Fig. 1 (top) and in Figs. 4 and 8 for cross-entropy and MSE loss respectively. Recall that it has been previously observed that cross-entropy or MSE with one-hot labels produce similar representations (Hui & Belkin, 2020). We obtain a few key observations. First, for any % of filtering, keeping the bottom subspace of the data produces higher test accuracy than when keeping the top subspace. That is, the subspace that is most useful for reconstruction (top) is least useful for perception. Second, the accuracy gap is impacted by the presence of background, finer-grained classes, and higher resolution images. This further validates our theoretical

observations from Fig. 2 and the result from Theorem 2. We will now focus on validating the second part of our claim that the subspace used for perception (bottom) is learned last, and slowly.

Useful features for perception are learned last. The above results demonstrate that the top subspace of the data—explaining most of the pixel variance—is not aligned with the perception tasks. Yet, perfect reconstruction implies capturing both the top and bottom subspace. Albeit correct, we demonstrate in Fig. 1 (bottom) and in Fig. 3 that the rate at which the bottom subspace is learned is exponentially slower than the rate at which the top subspace is learned. This empirically validates Theorem 2. For the reader familiar with optimization (Benzi, 2002), or power iteration methods, this observation is akin to their converge rate which is a function of the eigengap (Booth, 2006; Xu et al., 2018), i.e., the difference between λ_i and λ_{i+1} , where λ are the sorted eigenvalues. Because natural images have an exponential decay of their eigenvalues (Van der Schaaf & van Hateren, 1996; Ruderman, 1997), the rate at which the top subspace is approximated is exponentially faster than the bottom one, therefore making the learning of useful features for perception occur only late during training. This finding also resonates with previous studies on the spectral bias of DNs in classification and generative settings (Chakrabarty & Maji, 2019; Rahaman et al., 2019; Schwarz et al., 2021).

Combining the observations from this section supports **R1** and **R2** from Section 1. It remains to study **R3** which states that since features for perception lie within a negligible subspace (as measured by the reconstruction loss), one can find two separate models that equally solve the reconstruction task (same train and test loss values) but provide drastically different perception task performances.

4. Learning By Reconstruction Needs Guidance

We now turn to **R3** stating that features for perceptions lie within a space with negligible impact on the reconstruction error—therefore motivating the need to add additional guidance to the training, e.g., through denoising tasks. To do so, we will show that it is possible to obtain two DNs with same reconstruction error but one with perception capabilities far greater than the other (Section 4.1). Lastly we will prove that some guidance can be provided to the learned representation to reduce that gap and focus towards more useful features through careful design of the denoising task (Section 4.2).

4.1. Learning By Reconstruction Can Produce Optimal Representations

One interesting benefit emerging from the observations made for **R1** and **R2** is that guiding a DN to focus on the subspace containing informative features for perception has minimal impact on the reconstruction loss—as they focus on different subspaces. Therefore, we now propose a simple experiment to demonstrate the above argument. We take a resnet34 autoencoder train it with the usual reconstruction loss (MSE) on Imagenette. This gives us a model that (as per R1 and R2) fails to properly focus on discriminative feature. To obtain the second DN with improved classification performance, we simply add a classification head on top of the embedding of the encoder. That is, the same embedding that is fed to the decoder for reconstruction, is also fed to the linear classifier with supervised training loss (recall Eq. (8)). We obtain the key insight of **R3** which is that one can produce two DNs with same training loss (reconstruction) and validation loss (reconstruction), but with significantly different classification performance, as reported in Figs. 5 and 7.

To further understand that observation, we can recall the results from Theorem 1. We demonstrated that the encoder (V) having K dimensions at its disposal, is optimal when selecting the top singular vectors of the data matrix X . If K is large enough that it encompasses both the top subspace of the data (which is learned first and has greatest impact on the reconstruction loss) and the bottom subspace of the data (which is useful for perception as per Figs. 1 and 4), then both objective can coexist (recall Proposition 1) as long as enough capacity is given to the encoder. Therefore, we obtain the following key insight. *Whenever the capacity of the autoencoder is large, the encoder embedding can (and will at the end of training) include features useful for perception all while being able to reconstruct its inputs.* Again, for this to happen one requires the capacity of the encoder to grow with the image resolution (as more and more dimensions will be taken up by the top subspace), and with the complexity of the image background (again taking more dimensions in the top subspace) (recall Section 3.2).

The above observation demonstrates that learning to reconstruct needs an additional training signal to focus towards discriminative features. As we will prove below, learning by denoising offers such a solution.

4.2. Provable Benefits of Learning by Denoising

Recalling the Denoising Autoencoder setting from Eq. (2), we aim to obtain a closed form solution of the linear loss Eq. (3) in order to find some hints as to why masking and additive Gaussian noise produce representations of different quality for perception.

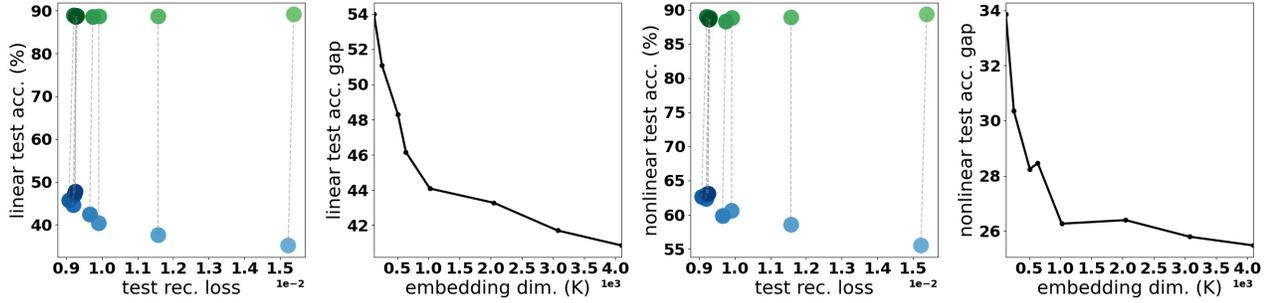


Figure 5. Depiction of multiple resnet34 autoencoders with varying embedding dimensions (**light to dark**) some trained only to reconstruct the input samples with data-augmentations (**blue**) and others with an additional supervised loss signal (as per Eq. (8)) (**green**). We report the test set accuracy and the relative difference (**y-axis**) for each of the “paired” models, i.e., the ones with every training setting identical except for the use of the supervised signal, as a function of the train and test rec loss. We clearly observe that for any embedding dimension and reconstruction loss, one can find two set of parameters with drastically different ability to solve perception tasks. Reconstructed samples and training curves are provided in Fig. 7.

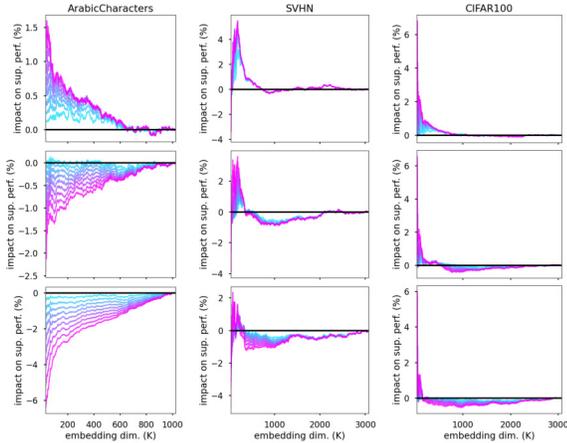


Figure 6. Depiction of the relative alignment difference when employing denoising tasks (recall Eq. (9)) with masking noise, with probability of dropping ranging from 0% to 99% (**cyan to pink**) for patch size of (1, 1) recovering multiplicative dropout (**top**), (2, 2) (**middle**), and (4, 4) (**bottom**) on various datasets. A positive number indicates a beneficial impact of using the denoising loss on the supervised performance of the learned representation. We observe that for datasets such as ArabicDigits that already have a strong alignment between the two tasks (recall Fig. 2), the use of any form of masking is detrimental except with shape (1, 1). However for datasets such as CIFAR100 (**right column**) with originally poor alignment, masking is beneficial and increases the alignment between the two tasks. As the original alignment increases with K , as the benefit of masking reduces.

Our goal is therefore to study the misalignment metric (Eq. (7)) under the denoising setting which, as per Corollary 1.2, is given by

$$\text{alignment}(k) \triangleq \min_{\mathbf{W}} \|\mathbf{W}^\top \mathbf{V}^{*\top} \mathbf{X} - \mathbf{Y}\|_F^2, \\ \mathbf{V}^* = \arg \min_{\mathbf{V}} \min_{\mathbf{Z}} \mathbb{E}_{\mathbf{X}'|\mathbf{X}} [\|\mathbf{Z}^\top \mathbf{V}^\top \mathbf{X}' - \mathbf{X}\|_F^2], \quad (9)$$

which is the minimum supervised loss that can be attained using the representation from \mathbf{V}^* that minimizes the denoising loss. We ought to highlight that we can obtain a closed form solution under the expectation over the noise distribution $(\mathbf{X}'|\mathbf{X})$ as formalized below, where we denote $\mathbf{G} \triangleq \mathbb{E}_{\mathbf{X}'|\mathbf{X}} [\mathbf{X}' \mathbf{X}'^\top]$ and $\mathbf{S} \triangleq \mathbb{E}_{\mathbf{X}'|\mathbf{X}} [\mathbf{X}']$.

Theorem 3. *The closed form solution for \mathbf{V}^* from Eq. (9) is given by \mathbf{V}^* spans $\mathbf{P}_G \mathbf{D}_G^{-\frac{1}{2}} (\mathbf{P}_H)_{:,1:K}$, where $\mathbf{H} \triangleq \mathbf{D}_G^{-\frac{1}{2}} \mathbf{P}_G^\top \mathbf{S} \mathbf{X}^\top \mathbf{X} \mathbf{S}^\top \mathbf{P}_G \mathbf{D}_G^{-\frac{1}{2}}$. (Proof in Section 5.4.)*

The above result demonstrates that even when employing denoising autoencoders with additive Gaussian noise or masking, we can obtain a closed form solution for \mathbf{V} , and from that obtain all the alignment metrics studied so far. In particular, Section 5.4 also provides the closed form solutions for \mathbf{G} and \mathbf{S} . We illustrate the alignment between reconstruction and perception tasks in the denoising autoencoder setting (Eq. (9)) in Fig. 6 for the case of random masking, as per the MAE setting. We clearly observe that the denoising task has the ability to increase the alignment between the two tasks—especially for small embedding dimension (K). We however observe that the size of the masked patches that provides the best gains vary with the dataset, hinging at another challenge of denoising autoencoders: the cross-validation of the denoising task. Another formal result we propose below will reinforce that point.

Denoting by $\mathbf{V}^*(\sigma)$ the optimal denoising autoencoder parameters when employing additive isotropic Gaussian noise with standard deviation σ , we obtain the following statement showcasing that this type of denoising task does not help supervised tasks.

Corollary 3.1. *Under the settings of Theorem 3, additive Gaussian noise has no impact in the supervised task performance as $\mathbf{W}^{*\top} \mathbf{V}^*(\sigma)^\top = \mathbf{W}^{*\top} \mathbf{V}^*(0)^\top, \forall \sigma \geq 0$, regardless of the supervised task. (Proof in Section 5.5.)*

We therefore obtain the following insights. *Denoising tasks offer a powerful guidance to skew learned representations to better align with perception tasks (Fig. 6) but some noise distributions such as additive Gaussian noise are provably unable to help.* A challenge that naturally emerges is in selecting the adequate denoising tasks, e.g., to avoid Corollary 3.1 in a setting where labels are not available and the supervised tasks to be tackled may not be known a priori. An interesting portal that we obtained (Corollary 3.1) in our study is that it is possible to assess if a denoising task has any impact on the perception task without yet knowing what is that supervised task. That alone could help in at least focusing on denoising tasks that do have an impact, albeit it will remain unknown if that impact will be beneficial or not.

5. Conclusion

We proposed to study the transferability of representations learned by reconstruction towards perception tasks. In particular, we obtain that the two objectives are fundamentally misaligned, with a degree of misalignment that grows with the presence of complicated background, with greater number of classes for the perception task, and with higher image resolutions. While our study focused on bringing those limitations forward from a theoretical and empirical angle, we also opened new avenues to reduce those limitations in the future. For example, we obtained a closed form solution to measure the impact of noise distributions to better align the learned representation to the downstream perception task. This novel methodology opens the door to a priori selecting noise distribution candidates. Even when the downstream task is unknown, we found that some noise distributions, such as additive Gaussian noise, are effectively unable to provide any benefit for better aligning reconstruction and perception tasks. On the opposite, we validated that masking is a valid strategy, albeit requiring some per-dataset tuning. That finding is in line with MAE’s performances going from about 50% to 74% on Imagenet top-1 accuracy when masking is employed. We hope that our study will also open new avenues to study reconstruction methods for other modalities such as time-series and NLP.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

Baidoo-Anu, D. and Ansah, L. O. Education in the era of generative artificial intelligence (ai): Understanding the potential benefits of chatgpt in promoting teaching and learning. *Journal of AI*, 7(1):52–62, 2023.

- Balestriero, R., Misra, I., and LeCun, Y. A data-augmentation is worth a thousand samples: Analytical moments and sampling-free training. *Advances in Neural Information Processing Systems*, 35:19631–19644, 2022.
- Balestriero, R., Ibrahim, M., Sobal, V., Morcos, A., Shekhar, S., Goldstein, T., Bordes, F., Bardes, A., Mialon, G., Tian, Y., et al. A cookbook of self-supervised learning. *arXiv preprint arXiv:2304.12210*, 2023.
- Ballé, J., Laparra, V., and Simoncelli, E. P. End-to-end optimization of nonlinear transform codes for perceptual quality. In *2016 Picture Coding Symposium (PCS)*, pp. 1–5. IEEE, 2016.
- Bao, J., Chen, D., Wen, F., Li, H., and Hua, G. Cvae-gan: fine-grained image generation through asymmetric training. In *Proceedings of the IEEE international conference on computer vision*, pp. 2745–2754, 2017.
- Bardes, A., Ponce, J., and LeCun, Y. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021.
- Barlow, H. B. Unsupervised learning. *Neural computation*, 1(3): 295–311, 1989.
- Benzi, M. Preconditioning techniques for large linear systems: a survey. *Journal of computational Physics*, 182(2):418–477, 2002.
- Booth, T. E. Power iteration method for the several largest eigenvalues and eigenfunctions. *Nuclear science and engineering*, 154(1):48–62, 2006.
- Bordes, F., Balestriero, R., and Vincent, P. High fidelity visualization of what your self-supervised representation knows about. *arXiv preprint arXiv:2112.09164*, 2021.
- Bottou, L. Stochastic gradient descent tricks. In *Neural Networks: Tricks of the Trade: Second Edition*, pp. 421–436. Springer, 2012.
- Brynjolfsson, E., Li, D., and Raymond, L. R. Generative ai at work. Technical report, National Bureau of Economic Research, 2023.
- Chakrabarty, P. and Maji, S. The spectral bias of the deep image prior. *arXiv preprint arXiv:1912.08905*, 2019.
- Chen, T., Kornblith, S., Swersky, K., Norouzi, M., and Hinton, G. E. Big self-supervised models are strong semi-supervised learners. *Advances in neural information processing systems*, 33:22243–22255, 2020.
- Chen, X., Mishra, N., Rohaninejad, M., and Abbeel, P. Pixelsnail: An improved autoregressive generative model. In *International Conference on Machine Learning*, pp. 864–872. PMLR, 2018.
- Dilokthanakul, N., Mediano, P. A., Garnelo, M., Lee, M. C., Salimbeni, H., Arulkumaran, K., and Shanahan, M. Deep unsupervised clustering with gaussian mixture variational autoencoders. *arXiv preprint arXiv:1611.02648*, 2016.
- Esmaili, B., Wu, H., Jain, S., Bozkurt, A., Siddharth, N., Paige, B., Brooks, D. H., Dy, J., and Meent, J.-W. Structured disentangled representations. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 2525–2534. PMLR, 2019.

- Garrido, Q., Balestriero, R., Najman, L., and Lecun, Y. Rankme: Assessing the downstream performance of pretrained self-supervised representations by their rank. In *International Conference on Machine Learning*, pp. 10929–10974. PMLR, 2023.
- Ghahramani, Z. Unsupervised learning. In *Summer school on machine learning*, pp. 72–112. Springer, 2003.
- Gutmann, M. and Hyvärinen, A. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pp. 297–304. JMLR Workshop and Conference Proceedings, 2010.
- He, K., Chen, X., Xie, S., Li, Y., Dollár, P., and Girshick, R. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16000–16009, 2022.
- Hui, L. and Belkin, M. Evaluation of neural architectures trained with square loss vs cross-entropy in classification tasks. *arXiv preprint arXiv:2006.07322*, 2020.
- Hyvärinen, A. and Dayan, P. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4), 2005.
- Jiang, Z., Zheng, Y., Tan, H., Tang, B., and Zhou, H. Variational deep embedding: An unsupervised and generative approach to clustering. *arXiv preprint arXiv:1611.05148*, 2016.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- Kulis, B. et al. Metric learning: A survey. *Foundations and Trends® in Machine Learning*, 5(4):287–364, 2013.
- LeCun, Y., Bengio, Y., and Hinton, G. Deep learning. *nature*, 521(7553):436–444, 2015.
- Lim, K.-L., Jiang, X., and Yi, C. Deep clustering with variational autoencoder. *IEEE Signal Processing Letters*, 27:231–235, 2020.
- Locatello, F., Bauer, S., Lucic, M., Raetsch, G., Gelly, S., Schölkopf, B., and Bachem, O. Challenging common assumptions in the unsupervised learning of disentangled representations. In *international conference on machine learning*, pp. 4114–4124. PMLR, 2019.
- Mahendran, A. and Vedaldi, A. Understanding deep image representations by inverting them. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5188–5196, 2015.
- Mathieu, E., Rainforth, T., Siddharth, N., and Teh, Y. W. Disentangling disentanglement in variational autoencoders. In *International conference on machine learning*, pp. 4402–4412. PMLR, 2019.
- Olah, C., Mordvintsev, A., and Schubert, L. Feature visualization. *Distill*, 2(11):e7, 2017.
- Rahaman, N., Baratin, A., Arpit, D., Draxler, F., Lin, M., Hamprecht, F., Bengio, Y., and Courville, A. On the spectral bias of neural networks. In *International Conference on Machine Learning*, pp. 5301–5310. PMLR, 2019.
- Ruderman, D. L. Origins of scaling in natural images. *Vision research*, 37(23):3385–3398, 1997.
- Schwarz, K., Liao, Y., and Geiger, A. On the frequency bias of generative models. *Advances in Neural Information Processing Systems*, 34:18126–18136, 2021.
- Selvaraju, R. R., Das, A., Vedantam, R., Cogswell, M., Parikh, D., and Batra, D. Grad-cam: Why did you say that? *arXiv preprint arXiv:1611.07450*, 2016.
- Seydoux, L., Balestriero, R., Poli, P., Hoop, M. d., Campillo, M., and Baraniuk, R. Clustering earthquake signals and background noises in continuous seismic data with unsupervised deep learning. *Nature communications*, 11(1):3972, 2020.
- Shen, Y., Gu, J., Tang, X., and Zhou, B. Interpreting the latent space of gans for semantic face editing. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9243–9252, 2020.
- Tran, L., Yin, X., and Liu, X. Disentangled representation learning gan for pose-invariant face recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1415–1424, 2017.
- Van den Oord, A., Kalchbrenner, N., Espeholt, L., Vinyals, O., Graves, A., et al. Conditional image generation with pixelcnn decoders. *Advances in neural information processing systems*, 29, 2016.
- Van der Schaaf, v. A. and van Hateren, J. v. Modelling the power spectra of natural images: statistics and information. *Vision research*, 36(17):2759–2770, 1996.
- Vincent, P., Larochelle, H., Lajoie, I., Bengio, Y., Manzagol, P.-A., and Bottou, L. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of machine learning research*, 11(12), 2010.
- Wang, Z., Bovik, A. C., Sheikh, H. R., and Simoncelli, E. P. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- Wold, S., Esbensen, K., and Geladi, P. Principal component analysis. *Chemometrics and intelligent laboratory systems*, 2(1-3): 37–52, 1987.
- Xu, P., He, B., De Sa, C., Mitliagkas, I., and Re, C. Accelerated stochastic power iteration. In *International Conference on Artificial Intelligence and Statistics*, pp. 58–67. PMLR, 2018.
- Zamir, A. R., Sax, A., Shen, W., Guibas, L. J., Malik, J., and Savarese, S. Taskonomy: Disentangling task transfer learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3712–3722, 2018.
- Zbontar, J., Jing, L., Misra, I., LeCun, Y., and Deny, S. Barlow twins: Self-supervised learning via redundancy reduction. In *International Conference on Machine Learning*, pp. 12310–12320. PMLR, 2021.
- Zeiler, M. D. and Fergus, R. Visualizing and understanding convolutional networks. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part I 13*, pp. 818–833. Springer, 2014.

SUPPLEMENTARY MATERIALS

5.1. Proof of Theorem 1

Proof. The first part of the proof finds the optimum \mathbf{W}^* and \mathbf{Z}^* as a function of \mathbf{V} which is direct since we are in a least-square style setting for each of them. The second part will consist in showing that the optimal \mathbf{V} can be found as the solution of a generalized eigenvalue problem. The third and final step will be to express the solution for \mathbf{V} in close-form that is also friendly for computations.

Step 1. Recall that our loss function is given by

$$\mathcal{L} = \|\mathbf{W}^\top \mathbf{V}^\top \mathbf{X} - \mathbf{Y}\|_F^2 + \lambda \|\mathbf{Z}^\top \mathbf{V}^\top \mathbf{X} - \mathbf{X}\|_F^2,$$

recalling that $\|\mathbf{M}\|_F^2 = \text{Tr}(\mathbf{M}^\top \mathbf{M})$, the above simplifies to

$$\begin{aligned} \mathcal{L} = & \|\mathbf{Y}\|_F^2 - 2 \text{Tr}(\mathbf{X}^\top \mathbf{V} \mathbf{W} \mathbf{Y}) + \text{Tr}(\mathbf{X}^\top \mathbf{V} \mathbf{W} \mathbf{W}^\top \mathbf{V}^\top \mathbf{X}) \\ & + \lambda \|\mathbf{X}\|_F^2 - 2\lambda \text{Tr}(\mathbf{X}^\top \mathbf{V} \mathbf{Z} \mathbf{X}) + \lambda \text{Tr}(\mathbf{X}^\top \mathbf{V} \mathbf{Z} \mathbf{Z}^\top \mathbf{V}^\top \mathbf{X}), \end{aligned}$$

we are now going to find the optimal \mathbf{W} and \mathbf{Z} which are unique by convexity of the loss and of their domain. Recall that we assume \mathbf{Y} and \mathbf{X} to be full-rank (therefore also making \mathbf{V} full-rank). Recalling the derive of traces, we obtain

$$\begin{aligned} \nabla_{\mathbf{W}} \mathcal{L} &= -2\mathbf{V}^\top \mathbf{X} \mathbf{Y}^\top + 2\mathbf{V}^\top \mathbf{X} \mathbf{X}^\top \mathbf{V} \mathbf{W}, \\ \nabla_{\mathbf{Z}} \mathcal{L} &= -2\lambda \mathbf{V}^\top \mathbf{X} \mathbf{X}^\top + 2\lambda \mathbf{V}^\top \mathbf{X} \mathbf{X}^\top \mathbf{V} \mathbf{Z}, \end{aligned}$$

setting it to zero (we assume here $\lambda > 0$ otherwise we can not solve for \mathbf{Z} since its value does not impact the loss) and solving leads

$$\begin{aligned} \mathbf{W}^* &= (\mathbf{V}^\top \mathbf{X} \mathbf{X}^\top \mathbf{V})^{-1} \mathbf{V}^\top \mathbf{X} \mathbf{Y}^\top, \\ \mathbf{Z}^* &= (\mathbf{V}^\top \mathbf{X} \mathbf{X}^\top \mathbf{V})^{-1} \mathbf{V}^\top \mathbf{X} \mathbf{X}^\top. \end{aligned}$$

We now have solved for \mathbf{W} , \mathbf{Z} as a function of \mathbf{V} , i.e., the loss is now only a function of \mathbf{V} , which we are going to solve for now.

Step 2. We will first proceed by plugging the values for \mathbf{W}^* , \mathbf{Z}^* back into the loss, which will now be only a function of \mathbf{V} . Let's first simplify our derivations by noticing that

$$\begin{aligned} \text{Tr}(\mathbf{X}^\top \mathbf{V} \mathbf{W}^* \mathbf{W}^{*\top} \mathbf{V}^\top \mathbf{X}) &= \text{Tr}(\mathbf{X}^\top \mathbf{V} (\mathbf{V}^\top \mathbf{X} \mathbf{X}^\top \mathbf{V})^{-1} \mathbf{V}^\top \mathbf{X} \mathbf{Y}^\top \mathbf{Y} \mathbf{X}^\top \mathbf{V} (\mathbf{V}^\top \mathbf{X} \mathbf{X}^\top \mathbf{V})^{-1} \mathbf{V}^\top \mathbf{X}) \\ &= \text{Tr}(\mathbf{X}^\top \mathbf{V} (\mathbf{V}^\top \mathbf{X} \mathbf{X}^\top \mathbf{V})^{-1} \mathbf{V}^\top \mathbf{X} \mathbf{Y}^\top \mathbf{Y}) \\ &= \text{Tr}(\mathbf{X}^\top \mathbf{V} \mathbf{W}^* \mathbf{Y}), \end{aligned}$$

and similarly

$$\text{Tr}(\mathbf{X}^\top \mathbf{V} \mathbf{Z}^* \mathbf{Z}^{*\top} \mathbf{V}^\top \mathbf{X}) = \text{Tr}(\mathbf{X}^\top \mathbf{V} \mathbf{Z}^* \mathbf{X}),$$

finally making the entire loss simplify as follows

$$\begin{aligned} \mathcal{L} &= \|\mathbf{Y}\|_F^2 - 2 \text{Tr}(\mathbf{X}^\top \mathbf{V} \mathbf{W}^* \mathbf{Y}) + \text{Tr}(\mathbf{X}^\top \mathbf{V} \mathbf{W}^* \mathbf{W}^{*\top} \mathbf{V}^\top \mathbf{X}) \\ &\quad + \lambda \|\mathbf{X}\|_F^2 - 2\lambda \text{Tr}(\mathbf{X}^\top \mathbf{V} \mathbf{Z}^* \mathbf{X}) + \lambda \text{Tr}(\mathbf{X}^\top \mathbf{V} \mathbf{Z}^* \mathbf{Z}^{*\top} \mathbf{V}^\top \mathbf{X}) \\ &= \|\mathbf{Y}\|_F^2 - \text{Tr}(\mathbf{X}^\top \mathbf{V} \mathbf{W}^* \mathbf{Y}) + \lambda \|\mathbf{X}\|_F^2 - \lambda \text{Tr}(\mathbf{X}^\top \mathbf{V} \mathbf{Z}^* \mathbf{X}) \\ &= \|\mathbf{Y}\|_F^2 - \text{Tr}(\mathbf{X}^\top \mathbf{V} (\mathbf{V}^\top \mathbf{X} \mathbf{X}^\top \mathbf{V})^{-1} \mathbf{V}^\top \mathbf{X} \mathbf{Y}^\top \mathbf{Y}) \\ &\quad + \lambda \|\mathbf{X}\|_F^2 - \lambda \text{Tr}(\mathbf{X}^\top \mathbf{V} (\mathbf{V}^\top \mathbf{X} \mathbf{X}^\top \mathbf{V})^{-1} \mathbf{V}^\top \mathbf{X} \mathbf{X}^\top \mathbf{X}) \\ &= \|\mathbf{Y}\|_F^2 + \lambda \|\mathbf{X}\|_F^2 - \text{Tr}(\mathbf{X}^\top \mathbf{V} (\mathbf{V}^\top \mathbf{X} \mathbf{X}^\top \mathbf{V})^{-1} \mathbf{V}^\top \mathbf{X} (\mathbf{Y}^\top \mathbf{Y} + \lambda \mathbf{X}^\top \mathbf{X})). \end{aligned}$$

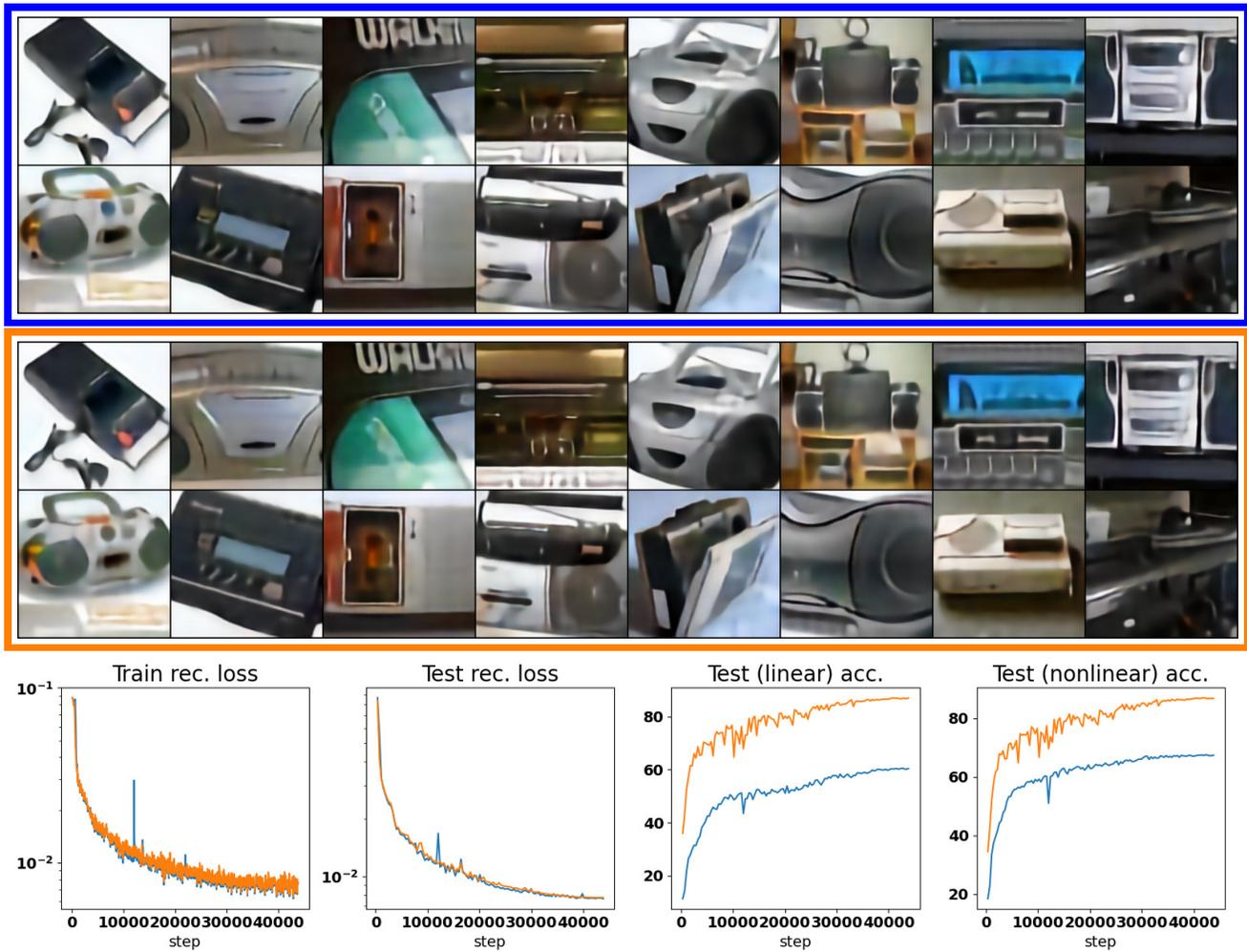


Figure 7. Depiction of two resnet34 autoencoders trained on Imagenette (Imagenet-10) images, one (orange) with an additional training signal that favors latent representations suited for classification, and the other (blue) that is only the reconstruction loss. As per **R1** and **R2** the latter naturally focuses on suboptimal features as showcases in the test accuracy, both when using a linear or a nonlinear probe. Crucially, the autoencoder with the additional signal produces representations with much greater discriminative power in both the linear and nonlinear setting. Yet, and despite popular belief, doing so has no impact on the reconstruction losses on the train or test set, and thus no impact on the quality of the reconstruction presented at the **top**. Therefore validating **R3**.

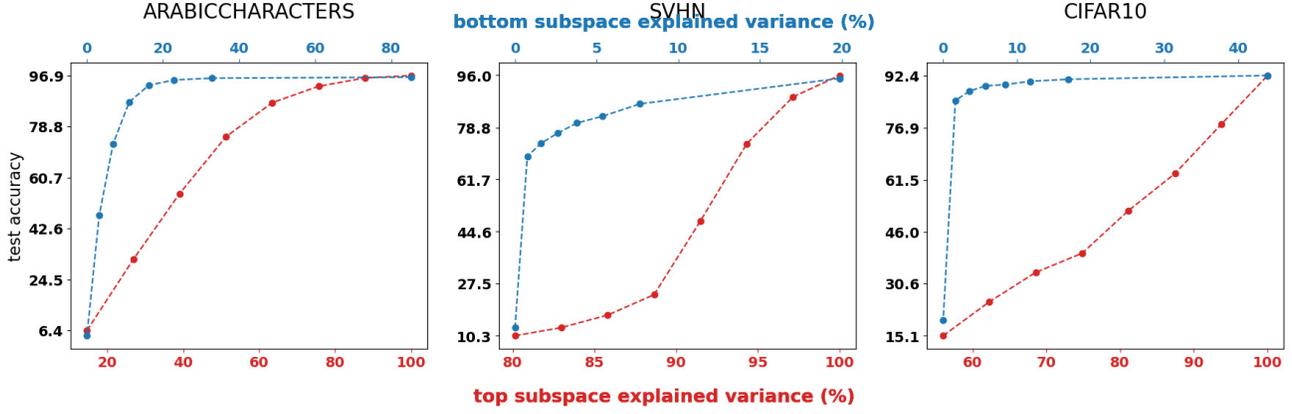


Figure 8. Reprise of Fig. 4 but now with MSE loss, the same observations can be made.

First, notice that both $\mathbf{X}(\mathbf{Y}^\top \mathbf{Y} + \lambda \mathbf{X}^\top \mathbf{X})\mathbf{X}^\top$ and $\mathbf{X}\mathbf{X}^\top$ are symmetric. Therefore, we can minimize the loss by solving the following generalized eigenvalue problem:

$$\text{find } \mathbf{V} \in \mathcal{M}_{D,D}(\mathbb{R}) \text{ so that } \mathbf{X}(\mathbf{Y}^\top \mathbf{Y} + \lambda \mathbf{X}^\top \mathbf{X})\mathbf{X}^\top \mathbf{V} = \mathbf{v}^\top \mathbf{X}\mathbf{X}^\top \mathbf{V} \Lambda,$$

where \mathbf{V} are the eigenvectors of the generalized eigenvalue problem, and Λ the eigenvalues, then the solution to our problem will be any rotation of any K eigenvectors, but the minimum will be achieved for the top- K ones.

Step 3. We will first demonstrate the general solution for the generalized eigenvalue problem. Given that solution, it will be easy to take the top- K eigenvectors that solve the considered problem. Denoting $\mathbf{A} \triangleq \mathbf{X}(\mathbf{Y}^\top \mathbf{Y} + \lambda \mathbf{X}^\top \mathbf{X})\mathbf{X}^\top$ and $\mathbf{B} \triangleq \mathbf{X}\mathbf{X}^\top$, and $\mathbf{H} \triangleq \mathbf{D}_B^{-\frac{1}{2}} \mathbf{P}_B^\top \mathbf{A} \mathbf{P}_B \mathbf{D}_B^{-\frac{1}{2}}$, we have

$$\begin{aligned} \mathbf{A} \mathbf{P}_B \mathbf{D}_B^{-\frac{1}{2}} \mathbf{P}_H &= \mathbf{B} \mathbf{P}_B \mathbf{D}_B^{-\frac{1}{2}} \mathbf{P}_H \Lambda \\ \Leftrightarrow \mathbf{P}_H^\top \mathbf{D}_B^{-\frac{1}{2}} \mathbf{P}_B^\top \mathbf{A} \mathbf{P}_B \mathbf{D}_B^{-\frac{1}{2}} \mathbf{P}_H &= \mathbf{P}_H^\top \mathbf{D}_B^{-\frac{1}{2}} \mathbf{P}_B^\top \mathbf{B} \mathbf{P}_B \mathbf{D}_B^{-\frac{1}{2}} \mathbf{P}_H \Lambda && (\mathbf{P}_H^\top \mathbf{D}_B^{-\frac{1}{2}} \mathbf{P}_B^\top \text{ bijective}) \\ \Leftrightarrow \mathbf{P}_H^\top \mathbf{D}_B^{-\frac{1}{2}} \mathbf{P}_B^\top \mathbf{A} \mathbf{P}_B \mathbf{D}_B^{-\frac{1}{2}} \mathbf{P}_H &= \Lambda \\ \Leftrightarrow \mathbf{D}_H &= \Lambda, \end{aligned}$$

therefore the eigenvalues are given by \mathbf{D}_H and the eigenvectors are given by $\mathbf{P}_B \mathbf{D}_B^{-\frac{1}{2}} \mathbf{P}_H$ or equivalently $(\mathbf{P}_{\mathbf{X}\mathbf{X}^\top} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{-\frac{1}{2}} \mathbf{P}_H)_{:,1:K}$. So the optimal \mathbf{V} is any rotation of the top- K eigenvectors. The above is simple to use as-is whenever $N > D$, if not, then we can obtain a solution without having to compute any $D \times D$ matrix, thus making the process more efficient. To that end, we can obtain

$$\begin{aligned} \mathbf{M} &\triangleq \mathbf{Y}^\top \mathbf{Y} + \lambda \mathbf{X}^\top \mathbf{X} \\ \mathbf{S} &\triangleq \mathbf{D}_M^{\frac{1}{2}} \mathbf{P}_M^\top \mathbf{X}^\top \\ \mathbf{P}_A &= \mathbf{V}_S, \mathbf{D}_A = \Sigma_S^2, \end{aligned}$$

that only involves $D \times \min(D, N)$ matrices instead of $D \times D$. □

5.2. Proof of Corollary 1.1

Proof. We will start with the fully supervised (least-square) proof obtained when $\lambda = 0$. Also notice that in any case, we have that

$$\mathbf{X} = \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{\frac{1}{2}} \mathbf{V}_X^\top.$$

Ordinary Least Square recovery. since $\lambda = 0$ we also have

$$A = X (Y^\top Y + \lambda X^\top X) X^\top = XY^\top YX^\top$$

when then lead to

$$\begin{aligned} H &= D_{XX^\top}^{-\frac{1}{2}} P_{XX^\top}^\top A P_{XX^\top} D_{XX^\top}^{-\frac{1}{2}} \\ &= D_{XX^\top}^{-\frac{1}{2}} P_{XX^\top}^\top XY^\top YX^\top P_{XX^\top} D_{XX^\top}^{-\frac{1}{2}} \\ &= D_{XX^\top}^{-\frac{1}{2}} P_{XX^\top}^\top P_{XX^\top} D_{XX^\top}^{\frac{1}{2}} V_X^\top Y^\top Y V_X D_{XX^\top}^{\frac{1}{2}} P_{XX^\top}^\top P_{XX^\top} D_{XX^\top}^{-\frac{1}{2}} \\ &= V_X^\top Y^\top Y V_X, \end{aligned}$$

from the above, we can simply plug those values in the analytical form for V^* from Eq. (4) to obtain

$$V = P_{XX^\top} D_{XX^\top}^{-\frac{1}{2}} P_H = P_{XX^\top} D_{XX^\top}^{-\frac{1}{2}} V_X^\top V_Y = X^\dagger(V_Y)_{:,1:K}$$

since we easily see that $P_H = V_X^\top V_Y$. We also have the optimum for W from Eq. (5) to be

$$\begin{aligned} W &= (V^{*\top} X X^\top V^*)^{-1} V^{*\top} X Y^\top \\ &= (V_Y^\top V_Y)^{-1} V_Y^\top V_X D_{XX^\top}^{-\frac{1}{2}} P_{XX^\top}^\top X Y^\top \\ &= V_Y^\top V_X V_X^\top Y^\top \\ &= \Sigma_Y U_Y^\top \end{aligned}$$

and finally the product of both matrices (which produce the supervised linear model) is obtained as

$$W^\top V^\top = U_Y \Sigma_Y (V_Y)_{:,1:K}^\top (X^\dagger)^\top = Y X^\top (X X^\top)^{-1}, K \geq C,$$

therefore recovering the OLS optimal solution. Note that if $K < C$ then we have a bottleneck and we therefore obtain an interesting alternative solution that looks at the top subspace of Y (this is however never the case in OLS settings).

Principal Component Analysis recovery. We now consider the case where we only employ the unsupervised loss (akin to $\lambda \rightarrow \infty$). In this setting we get

$$A = X X^\top X X^\top,$$

and we directly obtain

$$\begin{aligned} H &= D_{XX^\top}^{-\frac{1}{2}} P_{XX^\top}^\top A P_{XX^\top} D_{XX^\top}^{-\frac{1}{2}} \\ &= D_{XX^\top}^{-\frac{1}{2}} P_{XX^\top}^\top X X^\top X X^\top P_{XX^\top} D_{XX^\top}^{-\frac{1}{2}} \\ &= D_{XX^\top}, \end{aligned}$$

therefore the optimal form for V will be

$$V = P_{XX^\top} D_{XX^\top}^{-\frac{1}{2}} (P_H)_{:,1:K} = P_{XX^\top} (D_{XX^\top}^{-\frac{1}{2}})_{:,1:K},$$

which will select the top- K subspace of X (recall that the eigenvalues of H are D_{XX^\top} and therefore its top- K eigenvectors are selected the top- K dimension of the subspace. Then the solution for Z from Eq. (6) gives

$$\begin{aligned} Z &= (V^{*\top} X X^\top V^*)^{-1} V^{*\top} X X^\top \\ &= V^{*\top} X X^\top \\ &= ((D_{XX^\top}^{\frac{1}{2}})_{:,1:K})^\top P_{XX^\top}^\top, \end{aligned}$$

and lastly the product of Z and V (which produce the final linear transformation processing X takes the form

$$Z^\top V^\top = P_{XX^\top} (D_{XX^\top}^{\frac{1}{2}})_{:,1:K} ((D_{XX^\top}^{-\frac{1}{2}})_{:,1:K})^\top P_{XX^\top}^\top = (P_{XX^\top})_{:,1:K} ((P_{XX^\top})_{:,1:K})^\top,$$

which is the projection matrix onto the top- K subspace of the data X i.e. recovering the optimal solution of Principal Component Analysis. \square

5.3. Proof of Corollary 1.2

Proof. Recall that in the $\lambda \rightarrow \infty$ regime, we have that $\mathbf{V}^* = \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} (\mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{-\frac{1}{2}})_{.,1:K}$ and $\mathbf{W}^* = \mathbf{P}_{\mathbf{X}\mathbf{X}^\top}^\top \mathbf{Y}^\top$. We thus develop

$$\begin{aligned}
 \|\mathbf{W}^{*\top} \mathbf{V}^{*\top} \mathbf{X} - \mathbf{Y}\|_F^2 &= \|\mathbf{Y} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} ((\mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{-\frac{1}{2}})_{.,1:K})^\top \mathbf{P}_{\mathbf{X}\mathbf{X}^\top}^\top \mathbf{X} - \mathbf{Y}\|_F^2 \\
 &= \|\mathbf{Y} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} ((\mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{-\frac{1}{2}})_{.,1:K})^\top \mathbf{P}_{\mathbf{X}\mathbf{X}^\top}^\top \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{\frac{1}{2}} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top}^\top - \mathbf{Y}\|_F^2 \\
 &= \|\mathbf{Y} (\mathbf{P}_{\mathbf{X}\mathbf{X}^\top})_{.,1:K} ((\mathbf{P}_{\mathbf{X}\mathbf{X}^\top})_{.,1:K})^\top - \mathbf{Y}\|_F^2 \\
 &= \|\mathbf{Y}\|_F^2 - 2 \operatorname{Tr}(\mathbf{Y} (\mathbf{P}_{\mathbf{X}\mathbf{X}^\top})_{.,1:K} ((\mathbf{P}_{\mathbf{X}\mathbf{X}^\top})_{.,1:K})^\top \mathbf{Y}^\top) + \|\mathbf{Y} (\mathbf{P}_{\mathbf{X}\mathbf{X}^\top})_{.,1:K} ((\mathbf{P}_{\mathbf{X}\mathbf{X}^\top})_{.,1:K})^\top\|_F^2 \\
 &= \|\mathbf{Y}\|_F^2 - \operatorname{Tr}(\mathbf{Y} (\mathbf{P}_{\mathbf{X}\mathbf{X}^\top})_{.,1:K} ((\mathbf{P}_{\mathbf{X}\mathbf{X}^\top})_{.,1:K})^\top \mathbf{Y}^\top) \\
 &= \|\mathbf{Y}\|_F^2 - \|\mathbf{Y} (\mathbf{P}_{\mathbf{X}\mathbf{X}^\top})_{.,1:K}\|_F^2,
 \end{aligned}$$

as $\|\mathbf{Y}\|_F^2$ is a constant with respect to the parameters, we consider $\|\mathbf{Y} (\mathbf{P}_{\mathbf{X}\mathbf{X}^\top})_{.,1:K}\|_F^2$ as our alignment measure (the greater, the better the supervised loss can be minimized from the parameters). Since this quantity lives in the range $[0, \|\mathbf{Y} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top}\|_F^2]$, we see that by using the reparametrization from Eq. (7) we obtain the proposed measure of alignment rescaled to $[0, 1]$. \square

5.4. Proof of Theorem 3

We need to find the optimal reconstruction solution $(\mathbf{V}^*, \mathbf{Z}^*)$ first (corresponding to the case of $\lambda \rightarrow \infty$, and then plug it into the supervised loss with optimal \mathbf{W} .

$$\mathbb{E}_{\mathbf{X}'_n \sim p_{\mathbf{X}'|\mathbf{X}}} \|\mathbf{W}^\top \mathbf{V}^\top \mathbf{X}' - \mathbf{X}\|_2^2,$$

from which we obtain

$$\min_{\mathbf{W}, \mathbf{V}} \sum_{n=1}^N \mathbb{E}_{\mathbf{x}'_n \sim p_{\mathbf{x}'|\mathbf{x}_n}} \operatorname{Tr}(\mathbf{W}^\top \mathbf{V}^\top \mathbf{X}' \mathbf{X}'^\top \mathbf{V} \mathbf{W}) - 2 \operatorname{Tr}(\mathbf{W}^\top \mathbf{V}^\top \mathbf{X}' \mathbf{X}^\top) + \text{cst},$$

leading to $\mathbf{W}^* = (\mathbf{V}^\top \mathbb{E}[\mathbf{X}' \mathbf{X}'^\top] \mathbf{V})^{-1} \mathbf{V}^\top \mathbb{E}[\mathbf{X}'] \mathbf{X}^\top$ which we can plug back into the loss to obtain

$$\begin{aligned}
 &\operatorname{Tr}(\mathbf{X} \mathbb{E}[\mathbf{X}']^\top \mathbf{V} (\mathbf{V}^\top \mathbb{E}[\mathbf{X}' \mathbf{X}'^\top] \mathbf{V})^{-1} \mathbf{V}^\top \mathbb{E}[\mathbf{X}' \mathbf{X}'^\top] \mathbf{V} (\mathbf{V}^\top \mathbb{E}[\mathbf{X}' \mathbf{X}'^\top] \mathbf{V})^{-1} \mathbf{V}^\top \mathbb{E}[\mathbf{X}'] \mathbf{X}^\top) \\
 &\quad - 2 \operatorname{Tr}(\mathbf{X} \mathbb{E}[\mathbf{X}']^\top \mathbf{V} (\mathbf{V}^\top \mathbb{E}[\mathbf{X}' \mathbf{X}'^\top] \mathbf{V})^{-1} \mathbf{V}^\top \mathbb{E}[\mathbf{X}'] \mathbf{X}^\top) + \text{cst}, \\
 &= - \operatorname{Tr}(\mathbf{X} \mathbb{E}[\mathbf{X}']^\top \mathbf{V} (\mathbf{V}^\top \mathbb{E}[\mathbf{X}' \mathbf{X}'^\top] \mathbf{V})^{-1} \mathbf{V}^\top \mathbb{E}[\mathbf{X}'] \mathbf{X}^\top) + \text{cst}, \\
 &= - \operatorname{Tr}(\mathbf{V}^\top \mathbb{E}[\mathbf{X}'] \mathbf{X}^\top \mathbf{X} \mathbb{E}[\mathbf{X}']^\top \mathbf{V} (\mathbf{V}^\top \mathbb{E}[\mathbf{X}' \mathbf{X}'^\top] \mathbf{V})^{-1}) + \text{cst},
 \end{aligned}$$

whose solution is given (assuming $\mathbb{E}[\mathbf{X}' \mathbf{X}'^\top]$ is full-rank) by the solution of the generalized eigenvalue problem

$$\arg \max_{\mathbf{v}} \frac{\mathbf{v}^\top \mathbb{E}[\mathbf{X}'] \mathbf{X}^\top \mathbf{X} \mathbb{E}[\mathbf{X}']^\top \mathbf{v}}{\mathbf{v}^\top \mathbb{E}[\mathbf{X}' \mathbf{X}'^\top] \mathbf{v}}.$$

Given those optimum we can thus obtain the alignment measure same as before as the supervised loss obtain from \mathbf{V}^* from the unsupervised loss and \mathbf{Z}^* from the supervised one:

$$\mathbf{Z}^* = \arg \min_{\mathbf{Z}} \mathbb{E}_{\mathbf{X}'_n \sim p_{\mathbf{X}'|\mathbf{X}}} \|\mathbf{Z}^\top \mathbf{V}^{*\top} \mathbf{X}' - \mathbf{Y}\|_2^2,$$

whose optimum is therefore given by $\mathbf{Z}^* = (\mathbf{V}^{*\top} \mathbb{E}[\mathbf{X}' \mathbf{X}'^\top] \mathbf{V}^*)^{-1} \mathbf{V}^{*\top} \mathbb{E}[\mathbf{X}'] \mathbf{Y}^\top$ which can then be plugged back to produce the (unnormalized) measure. The analytical form of can be obtained (following the derivations of (Balestriero et al., 2022)) as follows for example for the case of additive Gaussian noise which is trivial and commonly done before e.g. showing the link between ridge regression and additive dropout $\mathbb{E}[\mathbf{X}'] = \mathbf{X}$ and $\mathbb{E}[\mathbf{X}' \mathbf{X}'^\top] = \mathbf{X} \mathbf{X}^\top + \sigma \mathbf{I}$. The perhaps more interesting derivations concern the masking as employed by MAE.

5.5. Proof of Corollary 3.1

In the case of additive, centered Gaussian noise, we have $\mathbb{E}[\mathbf{X}'] = \mathbf{X}$ and $\mathbb{E}[\mathbf{X}'\mathbf{X}'^\top] = \mathbf{X}\mathbf{X}^\top + \sigma\mathbf{I}$. Therefore the optimal value for \mathbf{V} , is given by solving the generalized eigenvalue problem $(\mathbf{X}\mathbf{X}^\top\mathbf{X}\mathbf{X}^\top, \mathbf{X}\mathbf{X}^\top + \sigma\mathbf{I})$. Recalling the derivations of the optimal solution for such problem from Section 5.1, we have that

$$\mathbf{V} = \mathbf{P}_{\mathbf{X}\mathbf{X}^\top + \sigma\mathbf{I}} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top + \sigma\mathbf{I}}^{-\frac{1}{2}} (\mathbf{P}_H)_{:,1:K} = \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} (\mathbf{D}_{\mathbf{X}\mathbf{X}^\top} + \sigma\mathbf{I})^{-\frac{1}{2}} (\mathbf{P}_H)_{:,1:K},$$

with

$$\begin{aligned} \mathbf{H} &= \mathbf{D}_{\mathbf{X}\mathbf{X}^\top + \sigma\mathbf{I}}^{-\frac{1}{2}} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top + \sigma\mathbf{I}}^\top (\mathbf{X}\mathbf{X}^\top\mathbf{X}\mathbf{X}^\top) \mathbf{P}_{\mathbf{X}\mathbf{X}^\top + \sigma\mathbf{I}} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top + \sigma\mathbf{I}}^{-\frac{1}{2}} \\ &= \mathbf{D}_{\mathbf{X}\mathbf{X}^\top + \sigma\mathbf{I}}^{-\frac{1}{2}} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top}^\top (\mathbf{X}\mathbf{X}^\top\mathbf{X}\mathbf{X}^\top) \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top + \sigma\mathbf{I}}^{-\frac{1}{2}} \\ &= (\mathbf{D}_{\mathbf{X}\mathbf{X}^\top})^2 (\mathbf{D}_{\mathbf{X}\mathbf{X}^\top} + \sigma\mathbf{I})^{-1}, \end{aligned}$$

and the important property to notice is that the ordering of the eigenvalues of \mathbf{H} which are given by $(\mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^\top)^2 (\mathbf{D}_{\mathbf{X}\mathbf{X}^\top} + \sigma\mathbf{I})^{-1}$ are the same as the ordering of $\mathbf{X}\mathbf{X}^\top$ which are given by $\mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^\top$. That is, the top- K subspace that will be picked up by \mathbf{P}_H are the same for any noise standard deviation σ . Now given the close form for \mathbf{V} we can obtain the close form of the classifier weights \mathbf{W} from Eq. (5) to be

$$\begin{aligned} \mathbf{W} &= (\mathbf{V}^{*\top} \mathbf{X}\mathbf{X}^\top \mathbf{V}^*)^{-1} \mathbf{V}^{*\top} \mathbf{X}\mathbf{Y}^\top \\ \mathbf{W} &= \left(((\mathbf{P}_H)_{:,1:K})^\top (\mathbf{D}_{\mathbf{X}\mathbf{X}^\top} + \sigma\mathbf{I})^{-\frac{1}{2}} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top}^\top \mathbf{X}\mathbf{X}^\top \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} (\mathbf{D}_{\mathbf{X}\mathbf{X}^\top} + \sigma\mathbf{I})^{-\frac{1}{2}} (\mathbf{P}_H)_{:,1:K} \right)^{-1} \\ &\quad \times ((\mathbf{P}_H)_{:,1:K})^\top (\mathbf{D}_{\mathbf{X}\mathbf{X}^\top} + \sigma\mathbf{I})^{-\frac{1}{2}} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top}^\top \mathbf{X}\mathbf{Y}^\top \\ \mathbf{W} &= \left(((\mathbf{P}_H)_{:,1:K})^\top (\mathbf{D}_{\mathbf{X}\mathbf{X}^\top} + \sigma\mathbf{I})^{-1} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top} (\mathbf{P}_H)_{:,1:K} \right)^{-1} ((\mathbf{P}_H)_{:,1:K})^\top (\mathbf{D}_{\mathbf{X}\mathbf{X}^\top} + \sigma\mathbf{I})^{-\frac{1}{2}} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{\frac{1}{2}} \mathbf{V}_X^\top \mathbf{Y}^\top \\ \mathbf{W} &= ((\mathbf{P}_H)_{:,1:K})^\top (\mathbf{D}_{\mathbf{X}\mathbf{X}^\top} + \sigma\mathbf{I}) \mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{-1} (\mathbf{P}_H)_{:,1:K} ((\mathbf{P}_H)_{:,1:K})^\top (\mathbf{D}_{\mathbf{X}\mathbf{X}^\top} + \sigma\mathbf{I})^{-\frac{1}{2}} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{\frac{1}{2}} \mathbf{V}_X^\top \mathbf{Y}^\top \\ \mathbf{W} &= ((\mathbf{P}_H)_{:,1:K})^\top \left((\mathbf{D}_{\mathbf{X}\mathbf{X}^\top} + \sigma\mathbf{I})^{\frac{1}{2}} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{-\frac{1}{2}} \right)_s \mathbf{V}_X^\top \mathbf{Y}^\top, \end{aligned}$$

where the s subscript indicates that only the top $K \times K$ part of the diagonal matrix is nonzero. Lastly, the product of both matrices (producing the supervised linear model) is obtained as

$$\begin{aligned} \mathbf{W}^\top \mathbf{V}^\top &= \mathbf{Y} \mathbf{V}_X \left((\mathbf{D}_{\mathbf{X}\mathbf{X}^\top} + \sigma\mathbf{I})^{\frac{1}{2}} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{-\frac{1}{2}} \right)_s (\mathbf{P}_H)_{:,1:K} ((\mathbf{P}_H)_{:,1:K})^\top (\mathbf{D}_{\mathbf{X}\mathbf{X}^\top} + \sigma\mathbf{I})^{-\frac{1}{2}} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top}^\top \\ &= \mathbf{Y} \mathbf{V}_X (\mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{-\frac{1}{2}})_s \mathbf{P}_{\mathbf{X}\mathbf{X}^\top}^\top \end{aligned}$$

therefore recovering the OLS optimal solution $\mathbf{Y}\mathbf{X}^\top(\mathbf{X}\mathbf{X}^\top)^{-1}$ whenever $K \geq D$, and otherwise recovers the projection onto the top subspace of \mathbf{X} —in any case the final parameters are invariant to the choice of the standard deviation of the additive Gaussian noise (σ) during the denoising autoencoder pre-training phase.

5.6. Proof of Theorem 2

Proof. We prove separately the two cases of linear and nonlinear decoder. The linear decoder setting relies on less assumptions about the data distributions but is also the less general statement. The nonlinear decoder setting is only valid near initialization as it assumes that the Jacobian matrix of the decoder is at random initialization, i.e., not yet aligned with the data manifold.

Linear decoder setting. The first part of the proof is to rewrite the joint objective classification and reconstruction objective with an arbitrary encoder network f_θ

$$\min_{\theta \in \mathbb{R}^P, \mathbf{W} \in \mathcal{M}_{C,K}(\mathbb{R}), \mathbf{V} \in \mathcal{M}_{D,K}(\mathbb{R})} \|\mathbf{W}f_\theta(\mathbf{X}) - \mathbf{Y}\|_F^2 + \lambda \|\mathbf{V}f_\theta(\mathbf{X}) - \mathbf{X}\|_F^2,$$

as the nonparametric version

$$\min_{\mathbf{Z} \in \mathcal{M}_{K,N}(\mathbb{R}), \mathbf{W} \in \mathcal{M}_{C,K}(\mathbb{R}), \mathbf{V} \in \mathcal{M}_{D,K}(\mathbb{R})} \|\mathbf{W}\mathbf{Z} - \mathbf{Y}\|_F^2 + \lambda \|\mathbf{V}\mathbf{Z} - \mathbf{X}\|_F^2,$$

both being identical if we assume that the encoder is powerful enough to reach any representation, which is a realistic assumption given current architectures. Given that nonparametric objective, we can now solve for both the optimal decoder weight \mathbf{V} and the optimal representation \mathbf{Z} as follows

$$\begin{aligned}
 & \min_{\mathbf{Z} \in \mathcal{M}_{K,N}(\mathbb{R}), \mathbf{W} \in \mathcal{M}_{C,K}(\mathbb{R}), \mathbf{V} \in \mathcal{M}_{D,K}(\mathbb{R})} \|\mathbf{W}\mathbf{Z} - \mathbf{Y}\|_F^2 + \lambda \|\mathbf{V}\mathbf{Z} - \mathbf{X}\|_F^2 \\
 = & \min_{\mathbf{Z} \in \mathcal{M}_{K,N}(\mathbb{R})} \|\mathbf{Y}\mathbf{Z}^\top (\mathbf{Z}\mathbf{Z}^\top)^{-1} \mathbf{Z} - \mathbf{Y}\|_F^2 + \lambda \|\mathbf{X}\mathbf{Z}^\top (\mathbf{Z}\mathbf{Z}^\top)^{-1} \mathbf{Z} - \mathbf{X}\|_F^2 \\
 = & \min_{\mathbf{Z} \in \mathcal{M}_{K,N}(\mathbb{R})} \text{Tr}(\mathbf{Z}^\top (\mathbf{Z}\mathbf{Z}^\top)^{-1} \mathbf{Z}\mathbf{Y}^\top \mathbf{Y}\mathbf{Z}^\top (\mathbf{Z}\mathbf{Z}^\top)^{-1} \mathbf{Z}) - 2 \text{Tr}(\mathbf{Y}^\top \mathbf{Y}\mathbf{Z}^\top (\mathbf{Z}\mathbf{Z}^\top)^{-1} \mathbf{Z}) + \|\mathbf{Y}\|_F^2 \\
 & + \lambda \text{Tr}(\mathbf{Z}^\top (\mathbf{Z}\mathbf{Z}^\top)^{-1} \mathbf{Z}\mathbf{X}^\top \mathbf{X}\mathbf{Z}^\top (\mathbf{Z}\mathbf{Z}^\top)^{-1} \mathbf{Z}) - 2\lambda \text{Tr}(\mathbf{X}^\top \mathbf{X}\mathbf{Z}^\top (\mathbf{Z}\mathbf{Z}^\top)^{-1} \mathbf{Z}) + \lambda \|\mathbf{X}\|_F^2 \\
 = & \min_{\mathbf{Z} \in \mathcal{M}_{K,N}(\mathbb{R})} \text{Tr}((\mathbf{Z}\mathbf{Z}^\top)^{-1} \mathbf{Z}\mathbf{Y}^\top \mathbf{Y}\mathbf{Z}^\top) - 2 \text{Tr}(\mathbf{Y}^\top \mathbf{Y}\mathbf{Z}^\top (\mathbf{Z}\mathbf{Z}^\top)^{-1} \mathbf{Z}) + \|\mathbf{Y}\|_F^2 \\
 & + \text{Tr}((\mathbf{Z}\mathbf{Z}^\top)^{-1} \mathbf{Z}\mathbf{X}^\top \mathbf{X}\mathbf{Z}^\top) - 2\lambda \text{Tr}(\mathbf{X}^\top \mathbf{X}\mathbf{Z}^\top (\mathbf{Z}\mathbf{Z}^\top)^{-1} \mathbf{Z}) + \lambda \|\mathbf{X}\|_F^2 \\
 = & \min_{\mathbf{Z} \in \mathcal{M}_{K,N}(\mathbb{R})} -\text{Tr}(\mathbf{Y}^\top \mathbf{Y}\mathbf{Z}^\top (\mathbf{Z}\mathbf{Z}^\top)^{-1} \mathbf{Z}) + \|\mathbf{Y}\|_F^2 - \lambda \text{Tr}(\mathbf{X}^\top \mathbf{X}\mathbf{Z}^\top (\mathbf{Z}\mathbf{Z}^\top)^{-1} \mathbf{Z}) + \lambda \|\mathbf{X}\|_F^2 \\
 = & \min_{\mathbf{Z} \in \mathcal{M}_{K,N}(\mathbb{R})} -\text{Tr}((\mathbf{Y}^\top \mathbf{Y} + \lambda \mathbf{X}^\top \mathbf{X}) \mathbf{Z}^\top (\mathbf{Z}\mathbf{Z}^\top)^{-1} \mathbf{Z}) + \|\mathbf{Y}\|_F^2 + \lambda \|\mathbf{X}\|_F^2 \\
 = & \min_{\mathbf{Z} \in \mathcal{M}_{K,N}(\mathbb{R}): \mathbf{Z}^\top \mathbf{Z} = \mathbf{I}} -\text{Tr}((\mathbf{Y}^\top \mathbf{Y} + \lambda \mathbf{X}^\top \mathbf{X}) \mathbf{Z}^\top \mathbf{Z}) + \|\mathbf{Y}\|_F^2 + \lambda \|\mathbf{X}\|_F^2
 \end{aligned}$$

which is solved by \mathbf{Z} being any orthogonal matrix in the subspace of the top- K eigenvectors of $(\mathbf{Y}^\top \mathbf{Y} + \lambda \mathbf{X}^\top \mathbf{X})$. Now as $\lambda \rightarrow \infty$ as the encoder will become more and more linear, ultimately converging to $f_\theta(\mathbf{x}) = \mathbf{U}\mathbf{x}$ with $\mathbf{U} \in \text{span}\{\text{eigvec}(\mathbf{X}^\top \mathbf{X})_1, \dots, \text{eigvec}(\mathbf{X}^\top \mathbf{X})_K\}$

Nonlinear decoder setting. To prove the result in the nonlinear decoder setting, our argument will look at the contribution of the data principal subspace to the gradient norm used to learn the encoder network parameters. In short, we will see that most of the gradient energy comes from the principal subspace of the data, effectively forcing the autoencoder model to focus on that subspace first. Let's first derive the gradient of the loss with respect to the encoder output

$$\begin{aligned}
 \nabla_{\mathcal{L}}(f(\mathbf{x}_n)) &= \sum_{n=1}^N \mathbf{J}_g(f(\mathbf{x}_n))^\top (g(f(\mathbf{x}_n)) - \mathbf{x}_n) \\
 &= \sum_{n=1}^N \mathbf{J}_g(f(\mathbf{x}_n))^\top g(f(\mathbf{x}_n)) - \sum_{n=1}^N \mathbf{J}_g(f(\mathbf{x}_n))^\top \mathbf{x}_n \\
 &= \sum_{n=1}^N \mathbf{J}_g(f(\mathbf{x}_n))^\top g(f(\mathbf{x}_n)) - \sum_{n=1}^N \sum_{k=1}^K \langle \mathbf{x}_n, \mathbf{v}_k(\mathbf{X}) \rangle \mathbf{J}_g(f(\mathbf{x}_n))^\top \mathbf{v}_k(\mathbf{X}) \\
 &= \sum_{n=1}^N \mathbf{J}_g(f(\mathbf{x}_n))^\top g(f(\mathbf{x}_n)) - \sum_{k=1}^K \sigma_k(\mathbf{X}) \left(\sum_{n=1}^N \mathbf{J}_g(f(\mathbf{x}_n))^\top \mathbf{U}_{n,k}(\mathbf{X}) \right) \mathbf{v}_k(\mathbf{X}),
 \end{aligned}$$

where we recall that $\mathbf{v}_k(\mathbf{X})$ is a unit-norm vector (right singular vector of \mathbf{X}), and we have singular values in descending order, $\sigma_i \geq \sigma_j, \forall i \leq j$. As a result, we first notice that the gradient will involve the sum of backpropagated right singular vector of \mathbf{X} weighted by the corresponding singular values. Since we know that

$$\sigma_{\min}^2 \left(\sum_{n=1}^N \mathbf{J}_g(f(\mathbf{x}_n)) \mathbf{U}_{n,k}(\mathbf{X}) \right) \leq \left\| \left(\sum_{n=1}^N \mathbf{J}_g(f(\mathbf{x}_n))^\top \mathbf{U}_{n,k}(\mathbf{X}) \right) \mathbf{v}_k(\mathbf{X}) \right\|_2^2 \leq \sigma_{\max}^2 \left(\sum_{n=1}^N \mathbf{J}_g(f(\mathbf{x}_n)) \mathbf{U}_{n,k}(\mathbf{X}) \right)$$

we can guarantee that the amount of information coming from the first principal singular vector of \mathbf{X} is greater than the information coming from the last singular vector as soon as

$$\frac{\sigma_1^2(\mathbf{X})}{\sigma_K^2(\mathbf{X})} \geq \frac{\sigma_{\max}^2 \left(\sum_{n=1}^N \mathbf{J}_g(f(\mathbf{x}_n)) \mathbf{U}_{n,K}(\mathbf{X}) \right)}{\sigma_{\min}^2 \left(\sum_{n=1}^N \mathbf{J}_g(f(\mathbf{x}_n)) \mathbf{U}_{n,1}(\mathbf{X}) \right)}.$$

Given that U is a unitary matrix, and that (at initialization) $J_g(f(\mathbf{x}_n))$ is well-behaved, we conclude that the above is always true at initialization in most practical settings where $\sigma_1^2(\mathbf{X}) \gg \sigma_K^2(\mathbf{X})$. Lastly, the fact that the principal subspace of \mathbf{X} impacts the gradient of the loss with respect to $f(\mathbf{x})$ is sufficient since training is done through backpropagation. Hence, any bias that $\nabla_{\mathcal{L}}(f(\mathbf{x}_n))$ has towards the principal subspace of the data will naturally flow through the backpropagation steps to all the internal parameters of f .

□

5.7. Eigendecomposition

Given a matrix $\mathbf{X} \in \mathcal{M}_{D,N}(\mathbb{R})$ with $D > N$, computing the eigendecomposition of $\mathbf{X}\mathbf{X}^\top$, a $D \times D$ matrix is $\mathcal{O}(D^3)$ which instead can be obtained in $\mathcal{O}(N^3 + DN^2)$ as

```
def fast_gram_eigh(X, major="C", unit_test=False):
    """
    compute the eigendecomposition of the Gram matrix:
    - XX.T using column (C) major notation
    - X.T@X using row (R) major notation
    """
    if major == "C":
        X_view = X.T
    else:
        X_view = X

    if X_view.shape[1] < X_view.shape[0]:
        # this case is the usual formula
        U, S = np.linalg.eigh(X_view.T @ X_view)
    else:
        # in this case we work in the tranpose domain
        U, S = np.linalg.eigh(X_view @ X_view.T)
        S = X_view.T @ S
        S[U>0] /= np.sqrt(U[U>0])
        # ensuring that we have the correct values
        if unit_test:
            Uslow, Sslow = np.linalg.eigh(X_view.T @ X_vew)
            assert np.allclose(U, Uslow)
            assert np.allclose(S, Sslow)

    return U, S
```

since we have the relation

$$\mathbf{X}\mathbf{X}^\top \mathbf{v} = \lambda \mathbf{v} \iff \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{v}) = \frac{\lambda}{\|\mathbf{X}^\top \mathbf{v}\|_2} (\mathbf{X}^\top \mathbf{v}),$$

and thus we can simply compute the eigenvectors of the $K \times K$ matrix $\mathbf{X}^\top \mathbf{X}$ and get the eigenvectors of the $N \times N$ matrix $\mathbf{X}\mathbf{X}^\top$ by left-multiplying them by \mathbf{X}^\top , and their corresponding eigenvalues are rescaled by $\frac{1}{\|\mathbf{X}^\top \mathbf{v}\|_2}$.

5.8. Fast Implementation of Alignment Metric (Eq. (7))

We want to sweep over the latent dimension K . As such, we can avoid recomputing the metric for each value and get them all at once as below. We again use the column-major notations as per Section 2:

```
def alignment_sweep(X, Y, major="C"):
    U, S, Vh = np.linalg.svd(X, full_matrices=False)
    if major == "C":
        denom = np.square(np.linalg.norm(Y @ Y.T))
        numer = np.linalg.multi_dot([Y.T, Y, Vh.T])
    else:
        denom = np.square(np.linalg.norm(Y.T @ Y))
        numer = np.linalg.multi_dot([Y, Y.T, U])
    numer = np.linalg.norm(numer, axis=0)**2
    return np.cumsum(numer) / denom
```

5.9. Additional figures

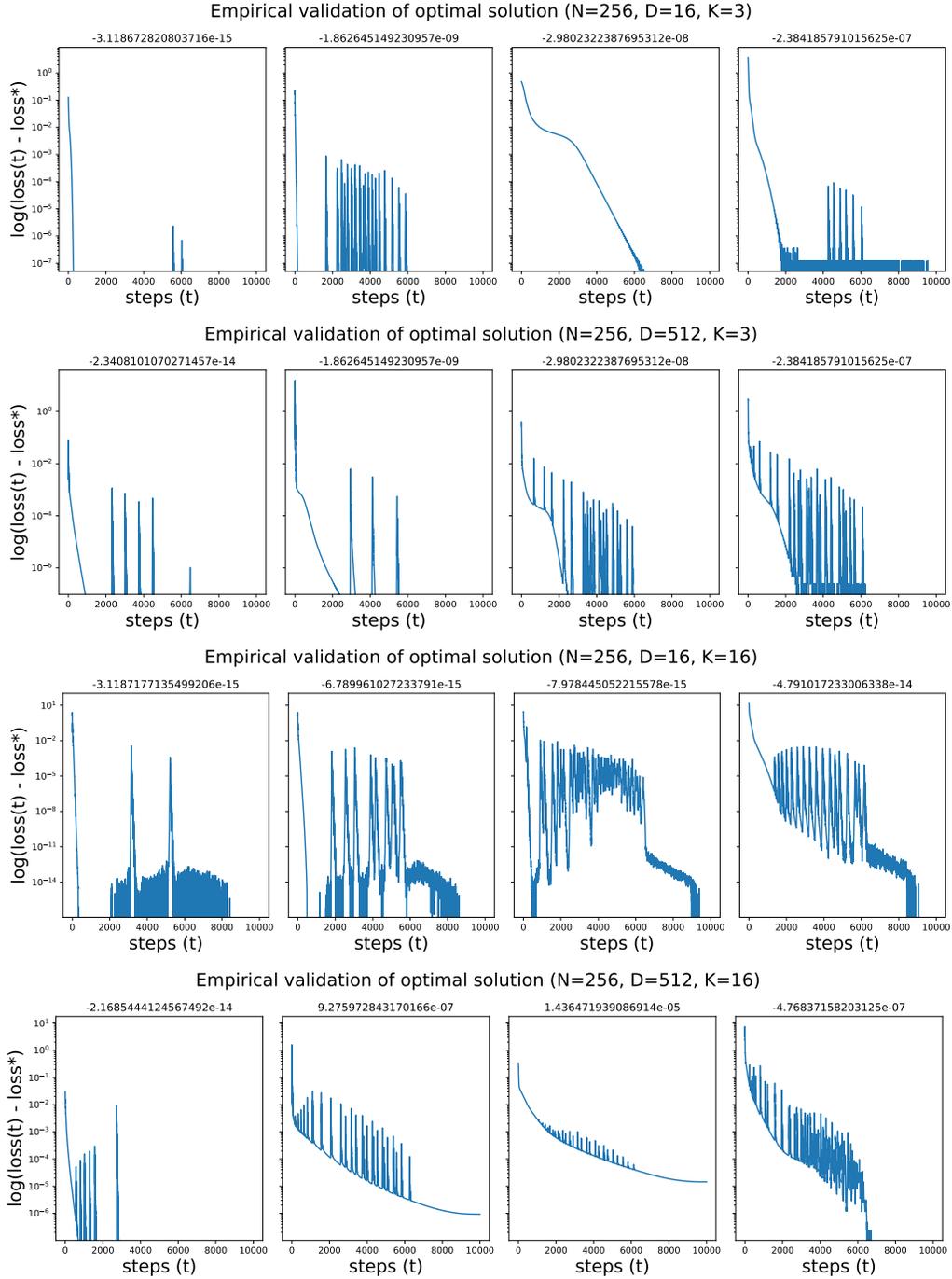


Figure 9. Empirical validation of Theorem 1 comparing the loss value at the optimum (from Eqs. (4) to (6)) against the one minimized with gradient descent (Adam optimizer) (y-axis) during gradient steps (t, x-axis). We expect that different to get close to 0 as the gradient updates converge to the minimum value of the loss. Although that quantity (loss(optimum) - loss(t)) is nonnegative in theory, we observe that its minimum value (reported in the title of each subplot) is sometimes negative with negligible value due to round off error. We compare numerous values of K, D, N as given in the titles of each row, and different values of $\lambda \in \{0.0, 0.1, 1, 10\}$ (column).

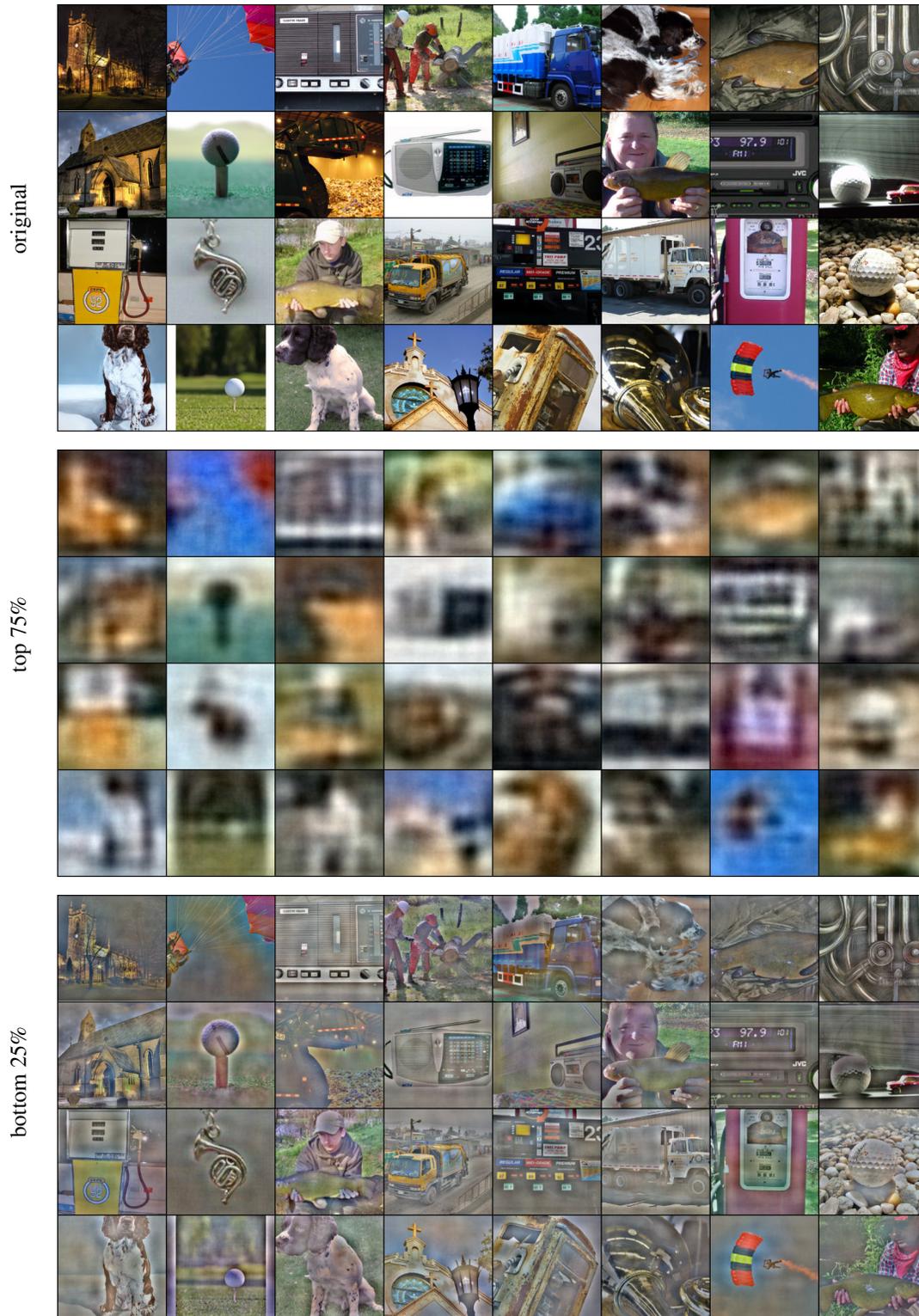


Figure 10. Depiction of Imagenet images (**top**) projected onto different subspaces obtained from Principal Component Analysis corresponding to the subspace explaining the top 75% of pixel variance (**middle**) and bottom 25% of pixel variance (**bottom**). We clearly observe that the image representation preserved after projection onto the bottom subspace makes the perception tasks (classification) easier to solve than if projected onto the top subspace, where the lower frequency information is insufficient to classify what is the object depicted (recall the classification performances of DNs applied onto those different projections from Figs. 1 and 4).