



Aligning Reconstruction-Based and Supervised Representations Through Masking

A discussion on: How Learning by Reconstruction Produces Uninformative Features For Perception (Balestrieri & LeCun, ICML 2024)

Panagiotis Koromilas

University of Athens

CV & Robotics reading group

7 October 2024

Q: Why does reconstruction-based learning generate **highly compelling reconstructed samples**, yet **struggle to produce informative latent representations** for downstream tasks, often **requiring extensive fine-tuning** to be effective?

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

- encoder $\mathbf{V} \in \mathcal{M}_{K,D}(\mathbb{R})$
- decoder $\mathbf{Z} \in \mathcal{M}_{D,K}(\mathbb{R})$
- predictor head $\mathbf{W} \in \mathcal{M}_{C,K}(\mathbb{R})$, C is the number of classes

Theorem 1. The combined loss function is minimized for

$$\mathbf{V}^* \text{ spans } \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} \mathbf{D}_{\mathbf{X}\mathbf{X}}^{-\frac{1}{2}} (\mathbf{P}_H)_{1:K},$$

$$\mathbf{W}^* = \left(\mathbf{V}^{*\top} \mathbf{X}\mathbf{X}^\top \mathbf{V}^* \right)^{-1} \mathbf{V}^{*\top} \mathbf{X}\mathbf{Y}^\top,$$

$$\mathbf{Z}^* = \left(\mathbf{V}^{*\top} \mathbf{X}\mathbf{X}^\top \mathbf{V}^* \right)^{-1} \mathbf{V}^{*\top} \mathbf{X}\mathbf{X}^\top,$$

$$\text{where } \mathbf{X}\mathbf{X}^\top = \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top}^\top, \text{ and } \mathbf{H} \triangleq \mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{-\frac{1}{2}} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top}^\top \mathbf{A} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{-\frac{1}{2}}.$$

Theorem 1. The combined loss function is minimized for

$$\mathbf{V}^* \text{ spans } \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} \mathbf{D}_{\mathbf{X}\mathbf{X}}^{-\frac{1}{2}} (\mathbf{P}_H)_{1:K},$$

$$\mathbf{W}^* = \left(\mathbf{V}^{*\top} \mathbf{X}\mathbf{X}^\top \mathbf{V}^* \right)^{-1} \mathbf{V}^{*\top} \mathbf{X}\mathbf{Y}^\top,$$

$$\mathbf{Z}^* = \left(\mathbf{V}^{*\top} \mathbf{X}\mathbf{X}^\top \mathbf{V}^* \right)^{-1} \mathbf{V}^{*\top} \mathbf{X}\mathbf{X}^\top,$$

where $\mathbf{X}\mathbf{X}^\top = \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top}^\top$, and $\mathbf{H} \triangleq \mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{-\frac{1}{2}} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top}^\top \mathbf{A} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{-\frac{1}{2}}$.

Proof steps:

- find optimal \mathbb{W}^* and \mathbb{Z}^* as a function of \mathbb{V}
- find optimal \mathbf{V} as the solution of a generalized eigenvalue problem

Theorem 1. The combined loss function is minimized for

\mathbf{V}^* spans $\mathbf{P}_{\mathbf{X}\mathbf{X}^\top} \mathbf{D}_{\mathbf{X}\mathbf{X}}^{-\frac{1}{2}} (\mathbf{P}_H)_{1:K}$,

$$\mathbf{W}^* = \left(\mathbf{V}^{*\top} \mathbf{X}\mathbf{X}^\top \mathbf{V}^* \right)^{-1} \mathbf{V}^{*\top} \mathbf{X}\mathbf{Y}^\top,$$

$$\mathbf{Z}^* = \left(\mathbf{V}^{*\top} \mathbf{X}\mathbf{X}^\top \mathbf{V}^* \right)^{-1} \mathbf{V}^{*\top} \mathbf{X}\mathbf{X}^\top,$$

where $\mathbf{X}\mathbf{X}^\top = \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top}^\top$, and $\mathbf{H} \triangleq \mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{-\frac{1}{2}} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top}^\top \mathbf{A} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} \mathbf{D}_{\mathbf{X}\mathbf{X}^\top}^{-\frac{1}{2}}$.

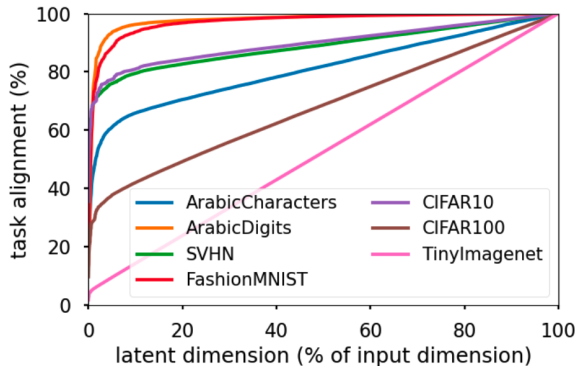
Corrolary: The solution from Theorem 1 recovers the OLS solution for $\mathbf{W}^{*\top} \mathbf{V}^{*\top}$ as $\lambda \rightarrow 0$, and the PCA solution for $\mathbf{Z}^{*\top} \mathbf{V}^{*\top}$ as $\lambda \rightarrow \infty$.

Q: under which condition \mathbb{V}^* is not impacted by λ ?

$$\text{alignment}(k) \triangleq \frac{\| \mathbf{Y}^\top \mathbf{Y} (\mathbf{P}_{\mathbf{X}\mathbf{X}^\top})_{1:k} \|_F^2}{\| \mathbf{Y}^\top \mathbf{Y} \mathbf{P}_{\mathbf{X}\mathbf{X}^\top} \|_F^2}$$

is the minimum supervised error that can be achieved given the $(\mathbf{V}^\top \mathbf{X})$ minimizes reconstruction which is measured by how much of the matrix $\mathbf{Y}^\top \mathbf{Y}$ can be reconstructed from the top- k subspa of $\mathbf{X}^\top \mathbf{X}$

Corollary 1.2. $\text{alignment}(k)$ increases with k , has value 0 iff the two losses are misaligned, and has value 1 iff the two losses are aligned.



- alignment for images without background
- alignment decreases for more classes (CIFAR10 vs CIFAR100)
- alignment decreases for higher resolution

Figure: Figure from Balestrieri & LeCun (ICML 2024)

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

Theorem 2. For any high-capacity encoder f_θ , studying the linear and the non-linear equation is equivalent at initialization for any decoder, and is always equivalent when the decoder is linear.

That is, linear and non-linear encoders
learn the principal subspace early in the training

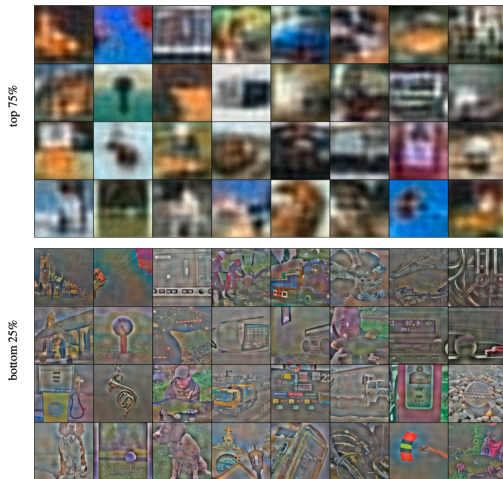


Figure: Figure from Balestrierio & LeCun (ICML 2024)

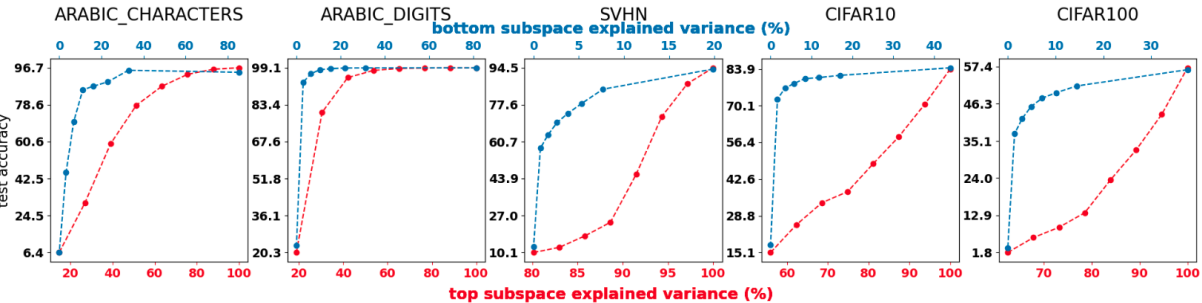
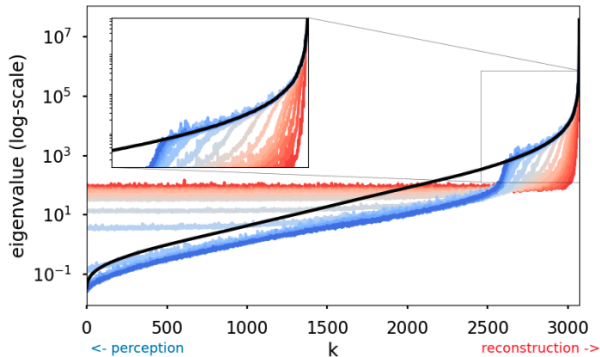


Figure: Figure from Balestrierio & LeCun (ICML 2024)



bottom subspace is learned exponentially slower than the top subspace

Figure: Figure from Balestrierio & LeCun (ICML 2024)

F-principle: DNNs fit target functions from low to high frequencies

Is this paper just a specific case of the F-principle?

How do other SSL methods perform well on downstream tasks while following the F-principle?

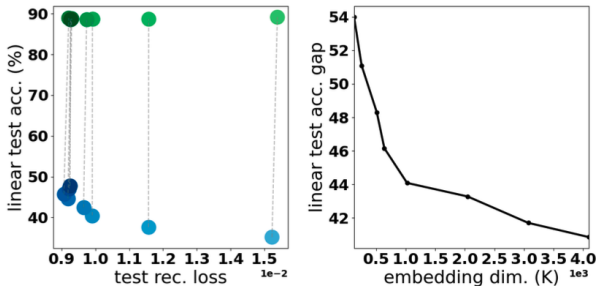


Figure: Figure from Balestrieri & LeCun (ICML 2024)

Enforcing a DN to use the "informative" subspace has minimal impact on the reconstruction loss

$$\text{alignment}(k) \triangleq \min_{\mathbf{W}} \left\| \mathbf{W}^{\top} \mathbf{V}^{*\top} \mathbf{X} - \mathbf{Y} \right\|_F^2,$$

$$\mathbf{V}^* = \arg \min_{\mathbf{V}} \min_{\mathbf{Z}} \mathbb{E}_{\mathbf{X}'|\mathbf{X}} \left[\left\| \mathbf{Z}^{\top} \mathbf{V}^{\top} \mathbf{X}' - \mathbf{X} \right\|_F^2 \right],$$

$$\text{alignment}(k) \triangleq \min_{\mathbf{W}} \|\mathbf{W}^\top \mathbf{V}^{*\top} \mathbf{X} - \mathbf{Y}\|_F^2,$$

$$\mathbf{V}^* = \arg \min_{\mathbf{V}} \min_{\mathbf{Z}} \mathbb{E}_{\mathbf{X}'|\mathbf{X}} \left[\|\mathbf{Z}^\top \mathbf{V}^\top \mathbf{X}' - \mathbf{X}\|_F^2 \right],$$

Theorem 3. The closed form solution for \mathbf{V}^* for this problem is given by \mathbf{V}^* spans $\mathbf{P}_G \mathbf{D}_G^{-\frac{1}{2}} (\mathbf{P}_H)_{:,1:K}$,

where $\mathbf{H} \triangleq \mathbf{D}_G^{-\frac{1}{2}} \mathbf{P}_G^\top \mathbf{S} \mathbf{X}^\top \mathbf{X} \mathbf{S}^\top \mathbf{P}_G \mathbf{D}_G^{-\frac{1}{2}}$ and $\mathbf{G} \triangleq \mathbb{E}_{\mathbf{X}'|\mathbf{X}} [\mathbf{X}' \mathbf{X}'^\top]$ and $\mathbf{S} \triangleq \mathbb{E}_{\mathbf{X}'|\mathbf{X}} [\mathbf{X}']$

$$\text{alignment}(k) \triangleq \min_{\mathbf{W}} \left\| \mathbf{W}^\top \mathbf{V}^{*\top} \mathbf{X} - \mathbf{Y} \right\|_F^2,$$

$$\mathbf{V}^* = \arg \min_{\mathbf{V}} \min_{\mathbf{Z}} \mathbb{E}_{\mathbf{X}'|\mathbf{X}} \left[\left\| \mathbf{Z}^\top \mathbf{V}^\top \mathbf{X}' - \mathbf{X} \right\|_F^2 \right],$$

Theorem 3. The closed form solution for \mathbf{V}^* for this problem is given by \mathbf{V}^* spans $\mathbf{P}_G \mathbf{D}_G^{-\frac{1}{2}} (\mathbf{P}_H)_{:,1:K}$,

where $\mathbf{H} \triangleq \mathbf{D}_G^{-\frac{1}{2}} \mathbf{P}_G^\top \mathbf{S} \mathbf{X}^\top \mathbf{X} \mathbf{S}^\top \mathbf{P}_G \mathbf{D}_G^{-\frac{1}{2}}$ and $\mathbf{G} \triangleq \mathbb{E}_{\mathbf{X}'|\mathbf{X}} [\mathbf{X}' \mathbf{X}'^\top]$ and $\mathbf{S} \triangleq \mathbb{E}_{\mathbf{X}'|\mathbf{X}} [\mathbf{X}']$

Corollary. Under these settings, **additive Gaussian noise has no impact in the downstream 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.

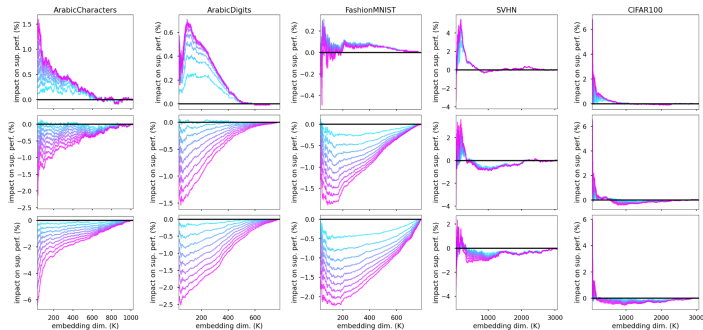


Figure 7. Depiction of the relative alignment difference when employing denoising tasks (recall Eq. (9)) with masking noise, with probability of dropping ranging from 0% to 99% (cyan to pink) for patch size of (1, 1) recovering multiplicative dropout (**top**), (2, 2) (**middle**), and (4, 4) (**bottom**) 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.

Figure: Figure from Balestrieri & LeCun (ICML 2024)

- **Misalignment** between learning by reconstruction and learning for downstream tasks

- **Misalignment** between learning by reconstruction and learning for downstream tasks
- **Ill-conditioned**: Perception features are learned last, requiring long training time

- **Misalignment** between learning by reconstruction and learning for downstream tasks
- **Ill-conditioned**: Perception features are learned last, requiring long training time
- **Ill-posed**: different model parameters can produce same reconstruction error but vastly different perception performance

- **Misalignment** between learning by reconstruction and learning for downstream tasks
- **Ill-conditioned**: Perception features are learned last, requiring long training time
- **Ill-posed**: different model parameters can produce same reconstruction error but vastly different perception performance
- **Masking helps**

- **Misalignment** between learning by reconstruction and learning for downstream tasks
- **Ill-conditioned**: Perception features are learned last, requiring long training time
- **Ill-posed**: different model parameters can produce same reconstruction error but vastly different perception performance
- **Masking helps**
- MAEs need:
 - large training times and architectures depending on image resolution, background etc
 - fine tuning

- [1] Balestrieri, R., & LeCun, Y. How Learning by Reconstruction Produces Uninformative Features For Perception. In Forty-first International Conference on Machine Learning.
- [2] Xu, Z. Q. J., Zhang, Y., Luo, T., Xiao, Y., & Ma, Z. Frequency Principle: Fourier Analysis Sheds Light on Deep Neural Networks. Communications in Computational Physics
- [3] Cao, Y., Fang, Z., Wu, Y., Zhou, D.-X., & Gu, Q.. Towards Understanding the Spectral Bias of Deep Learning. In Z.-H. Zhou, Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence.

