

Unsupervised and Self-Supervised Learning

Ju Sun

Computer Science & Engineering

University of Minnesota, Twin Cities

April 23, 2026

Our roadmap

Covered: Fundamentals

Fundamental belief: universal approximation theorem

Basics of numerical optimization

Training DNNs: basic methods and tricks

Covered: Structured data: images, sequences, graphs

Work with images: convolutional neural networks & applications

Work with sequences: recurrent neural networks & applications

Working with graphs: graph neural networks & applications

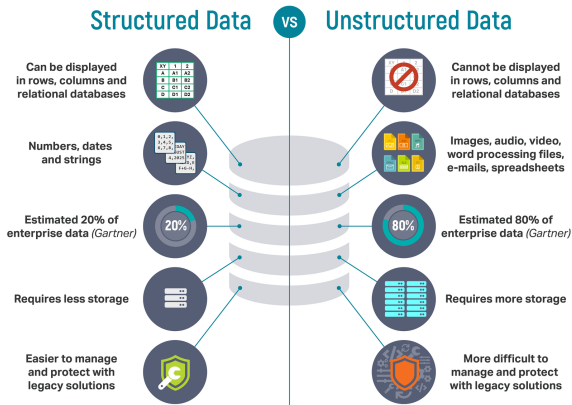
Transformers, large-language models, and foundation models

Now: Unsupervised/self-supervised/generative learning

Learning representation without labels: dictionary learning, autoencoders, self-supervised learning

Learning probability distributions: generative models

Structured vs. unstructured data

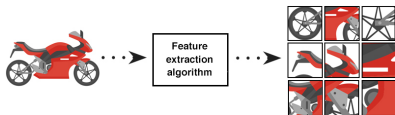


Credit: <https://lawtomated.com/>

[structured-data-vs-unstructured-data-what-are-they-and-why-care/](#)

- structured data typically **tabular data**
- structured data often directly fed into classical ML tools
- the success of DL mostly lies at **learning useful features/patterns from unstructured data**, i.e., **representation learning**

Feature engineering for unstructured data: old and new



Feature engineering: derive features for **efficient** learning

Credit: [Elgendy, 2020]

Traditional learning pipeline



- feature extraction is “independent” of the learning models and tasks
- features are handcrafted and/or learned

Modern learning pipeline

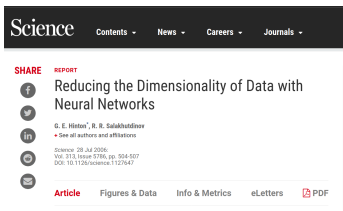


- end-to-end DNN learning

Unsupervised representation learning

Learning feature/representation **without task information** (e.g., labels)
(ICLR — International Conference on **Learning Representation**)

- **Historical:** Unsupervised representation learning key to the revival of deep learning (i.e., autoencoders, [[Hinton et al., 2006](#), [Hinton, 2006](#)])



Science Contents - News - Careers - Journals -

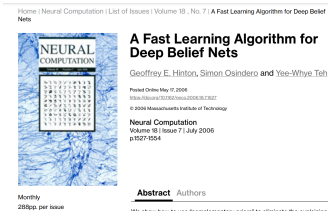
SHARE REPORT

Reducing the Dimensionality of Data with Neural Networks

G. E. Hinton¹, R. R. Salakhutdinov²
+ See all authors and affiliations

Science | 28 Jul 2006
Vol. 313, Issue 5746, pp. 504-507
DOI: 10.1126/science.1127647

Article Figures & Data Info & Metrics eLetters PDF



Home | Neural Computation | List of Issues | Volume 18, No. 7 | A Fast Learning Algorithm for Deep Belief Nets

NEURAL COMPUTATION

A Fast Learning Algorithm for Deep Belief Nets

Geoffrey E. Hinton, Simon Osindero and Yee-Whye Teh

Printed Online May 17, 2006
S0899-7667(06)00048-0/00000
© 2006 Massachusetts Institute of Technology

Neural Computation
Volume 18 | Issue 7, July 2006
p.1527-1554

Monthly
288pp. per issue

Abstract Authors

We show how to use "overparameterized networks" to alleviate the vanishing...

- **Practical:** Numerous advanced models built on top of the ideas in unsupervised representation learning (e.g., encoder-decoder networks, Transformers, U-Net in segmentation)

PCA for linear data

Autoencoder: extensions of PCA for nonlinear data

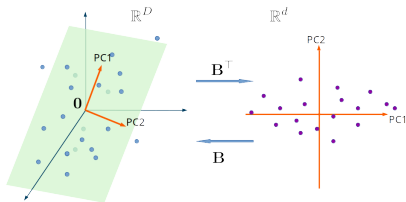
Applications of autoencoder

Self-supervised learning (SSL)

PCA: the geometric picture

Principal component analysis (PCA)

- $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^D$ zero-centered and write $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m]^\top \in \mathbb{R}^{m \times D}$
- Compact SVD $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$, where $\mathbf{V} \in \mathbb{R}^{D \times r}$ spans the row space of \mathbf{X}
- Take top right singular vectors \mathbf{B} from \mathbf{V} , and obtain $\mathbf{X}\mathbf{B}$



PCA is effectively to identify the best-fit subspace to $\mathbf{x}_1, \dots, \mathbf{x}_m$

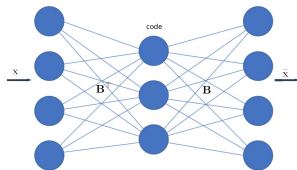
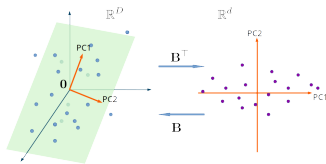
- \mathbf{B} has orthonormal columns, i.e., $\mathbf{B}^\top \mathbf{B} = \mathbf{I}$ ($\mathbf{B}\mathbf{B}^\top \neq \mathbf{I}$ when $D \neq d$)
- sample to representation:
 $\mathbf{x} \mapsto \mathbf{x}' \doteq \mathbf{B}^\top \mathbf{x}$ ($\mathbb{R}^D \rightarrow \mathbb{R}^d$, dimension reduction)
- representation to sample:
 $\mathbf{x}' \mapsto \hat{\mathbf{x}} \doteq \mathbf{B}\mathbf{x}'$ ($\mathbb{R}^d \rightarrow \mathbb{R}^D$)
- $\hat{\mathbf{x}} = \mathbf{B}\mathbf{B}^\top \mathbf{x} \approx \mathbf{x}$

Autoencoders

story in digital communications ...



autoencoder: [Bourlard and Kamp, 1988, Hinton and Zemel, 1994]



To find the basis B , solve ($d \leq D$)

$$\min_{B \in \mathbb{R}^{D \times d}} \sum_{i=1}^m \|x_i - BB^T x_i\|_2^2$$

or:

$$\min_{B \in \mathbb{R}^{D \times d}} \|X - XBB^T\|_F^2$$

– **Encoding:**

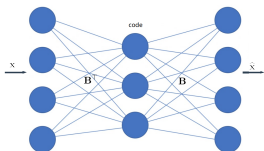
$$x \mapsto x' = B^T x$$

– **Decoding:**

$$x' \mapsto BB^T x = \hat{x}$$

Autoencoders

autoencoder:



To find the basis B , solve

$$\min_{B \in \mathbb{R}^{D \times d}} \sum_{i=1}^m \|x_i - BB^T x_i\|_2^2$$

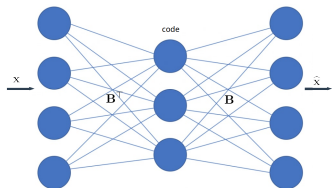
So the autoencoder is performing PCA!

One can even relax the weight tying:

$$\min_{B \in \mathbb{R}^{D \times d}, A \in \mathbb{R}^{D \times d}} \sum_{i=1}^m \|x_i - BA^T x_i\|_2^2,$$

which finds a basis (**not necessarily orthonormal**) B that spans the top singular space also [Baldi and Hornik, 1989], [Kawaguchi, 2016], [Lu and Kawaguchi, 2017].

Factorization



To perform PCA,

$$\min_{B \in \mathbb{R}^{D \times d}} \sum_{i=1}^m \|x_i - BB^T x_i\|_2^2$$
$$\min_{B \in \mathbb{R}^{D \times d}, A \in \mathbb{R}^{D \times d}} \sum_{i=1}^m \|x_i - BA^T x_i\|_2^2,$$

But: the basis B and the representations/codes z_i 's are all we care about

Factorization: (or autoencoder without encoder)

$$\min_{B \in \mathbb{R}^{D \times d}, z'_i s \in \mathbb{R}^d} \sum_{i=1}^m \|x_i - Bz_i\|_2^2.$$

All three formulations will find three **different** B 's that span the **same** principal subspace [Tan and Mayrovouniotis, 1995, Li et al., 2020b, Li et al., 2020a, Valavi et al., 2020]. They're all doing PCA!

Sparse coding

Factorization: (or autoencoder without encoder)

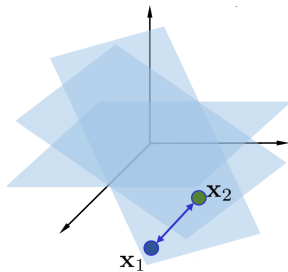
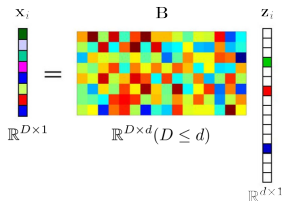
$$\min_{B \in \mathbb{R}^{D \times d}, z'_i s \in \mathbb{R}^d} \sum_{i=1}^m \|x_i - Bz_i\|_2^2.$$

What happens when we allow $d \geq D$? Underdetermined even if B is known.

Sparse coding (i.e., dictionary learning): assuming z_i 's are sparse and $d \geq D$

$$\min_{B \in \mathbb{R}^{D \times d}, z'_i s \in \mathbb{R}^d} \sum_{i=1}^m \|x_i - Bz_i\|_2^2 + \lambda \sum_{i=1}^m \Omega(z_i)$$

where Ω promotes sparsity, e.g., $\Omega = \|\cdot\|_1$.



More on sparse coding (dictionary learning)

MENU **nature**

Letter | Published: 13 June 1996

Emergence of simple-cell receptive field properties by learning a sparse code for natural images

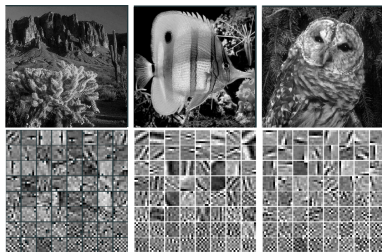
Bruno A. Olshausen & David J. Field

Nature 381, 607–609(1996) | [Cite this article](#)

5409 Accesses | 2901 Citations | 29 Altmetric | [Metrics](#)

Abstract

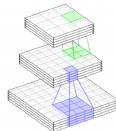
THE receptive fields of simple cells in mammalian primary visual cortex can be characterized as being spatially localized, oriented¹⁻⁴ and bandpass (selective to structure at different spatial scales), comparable to



denoising



super resol.



recognition

References: [Olshausen and Field, 1996, Mairal, 2014, Sun et al., 2017, Bai et al., 2018, Qu et al., 2019]

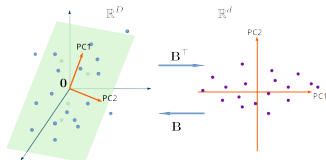
PCA for linear data

Autoencoder: extensions of PCA for nonlinear data

Applications of autoencoder

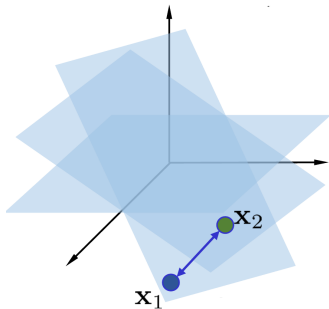
Self-supervised learning (SSL)

Quick summary of the linear models



PCA is effectively to identify
the best-fit subspace to

$$\mathbf{x}_1, \dots, \mathbf{x}_m$$



– B from V of $X = USV^T$

– autoencoder:

$$\min_{B \in \mathbb{R}^{D \times d}} \sum_{i=1}^m \|\mathbf{x}_i - BB^T \mathbf{x}_i\|_2^2$$

– autoencoder:

$$\min_{B \in \mathbb{R}^{D \times d}, A \in \mathbb{R}^{D \times d}} \sum_{i=1}^m \|\mathbf{x}_i - BA^T \mathbf{x}_i\|_2^2$$

– factorization:

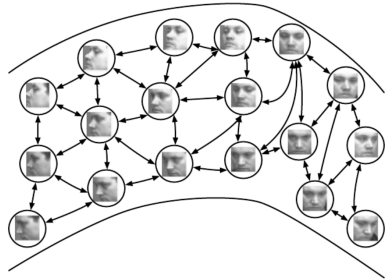
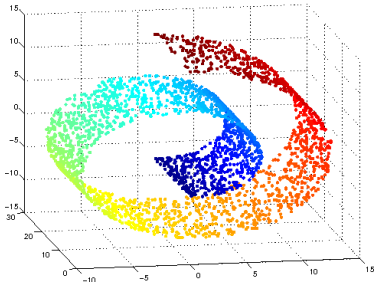
$$\min_{B \in \mathbb{R}^{D \times d}, \mathbf{z}'_i \in \mathbb{R}^d} \sum_{i=1}^m \|\mathbf{x}_i - B\mathbf{z}'_i\|_2^2$$

– when $d \geq D$, sparse coding/dictionary learning

$$\min_{B \in \mathbb{R}^{D \times d}, \mathbf{z}'_i \in \mathbb{R}^d} \sum_{i=1}^m \|\mathbf{x}_i - B\mathbf{z}'_i\|_2^2 + \lambda \sum_{i=1}^m \Omega(\mathbf{z}'_i)$$

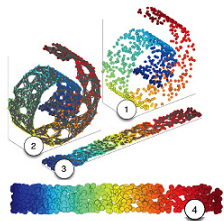
e.g., $\Omega = \|\cdot\|_1$

What about nonlinear data?



- Manifold, but not mathematically (i.e., differential geometry sense) rigorous
- **(No. 1?) (Union-of-)Manifold hypothesis for high-dimensional data:** practical data lie (approximately) on union of **low-dimensional** “manifolds”.
Why?
 - * data generating processes often controlled by very few parameters

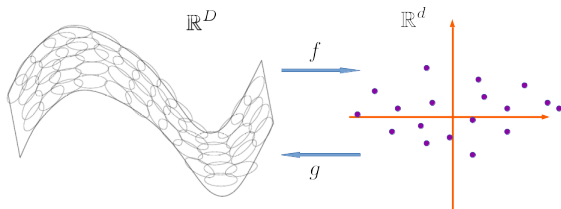
Manifold learning



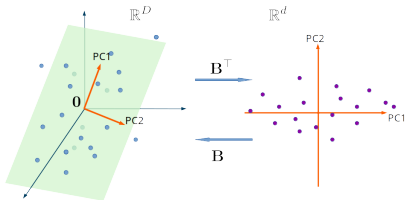
Classic methods (mostly for visualization): e.g.,

- ISOMAP [Tenenbaum, 2000]
- Locally-linear embedding [Roweis, 2000]
- Laplacian eigenmap [Belkin and Niyogi, 2001]
- t-distributed stochastic neighbor embedding (t-SNE) [van der Maaten and Hinton, 2008]

Nonlinear dimension reduction and representation learning

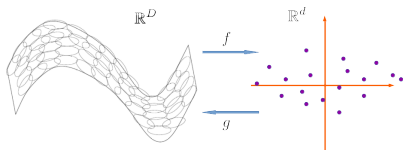


From autoencoders to deep autoencoders



$$\min_{B \in \mathbb{R}^{D \times d}} \sum_{i=1}^m \|\mathbf{x}_i - BB^T \mathbf{x}_i\|_2^2$$
$$\min_{B \in \mathbb{R}^{D \times d}, A \in \mathbb{R}^{D \times d}} \sum_{i=1}^m \|\mathbf{x}_i - BA^T \mathbf{x}_i\|_2^2$$

nonlinear generalization of the linear mappings:



deep autoencoders

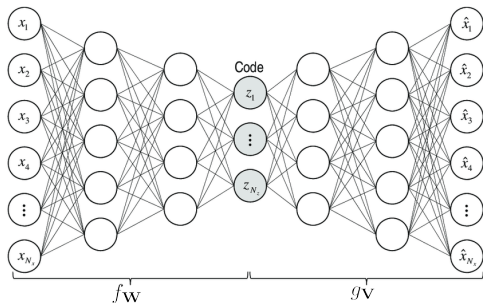
$$\min_{V, W} \sum_{i=1}^m \|\mathbf{x}_i - gV \circ fW(\mathbf{x}_i)\|_2^2$$

simply $A^T \rightarrow f_W$ and $B \rightarrow g_V$

A side question: why not calculate “nonlinear basis”?

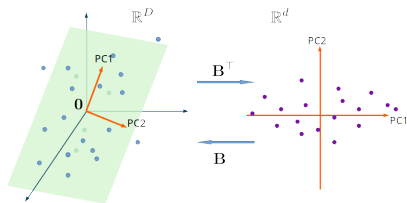
Deep autoencoders

$$\min_{\mathbf{V}, \mathbf{W}} \sum_{i=1}^m \|\mathbf{x}_i - g\mathbf{v} \circ f\mathbf{w}(\mathbf{x}_i)\|_2^2$$



the landmark paper [Hinton, 2006] ... that introduced **pretraining**

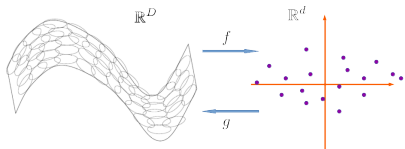
From factorization to deep factorization



factorization

$$\min_{B \in \mathbb{R}^{D \times d}, z'_i \in \mathbb{R}^d} \sum_{i=1}^m \|\mathbf{x}_i - Bz'_i\|_2^2$$

nonlinear generalization of the linear mappings:



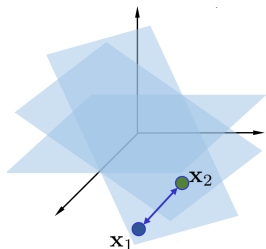
deep factorization

$$\min_{V, z'_i \in \mathbb{R}^d} \sum_{i=1}^m \|\mathbf{x}_i - gV(z'_i)\|_2^2$$

simply $B \rightarrow gV$

[Tan and Mayrovouniotis, 1995, Fan and Cheng, 2018, Bojanowski et al., 2017, Park et al., 2019, Li et al., 2020b], also known as **deep decoder**.

From sparse coding to deep sparse coding



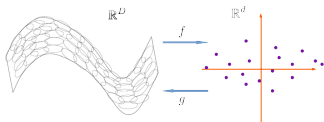
- when $d \geq D$, sparse coding/dictionary learning

$$\min_{B \in \mathbb{R}^{D \times d}, z'_i s \in \mathbb{R}^d} \sum_{i=1}^m \|\mathbf{x}_i - Bz_i\|_2^2 + \lambda \sum_{i=1}^m \Omega(z_i)$$

e.g., $\Omega = \|\cdot\|_1$

nonlinear generalization of the linear mappings: ($d \geq D$)

deep sparse coding/dictionary learning

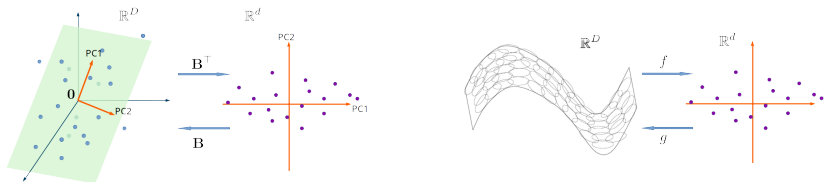


$$\min_{V, z'_i s \in \mathbb{R}^d} \sum_{i=1}^m \|\mathbf{x}_i - gV(z_i)\|_2^2 + \lambda \sum_{i=1}^m \Omega(z_i)$$

$$\min_{V, W} \sum_{i=1}^m \|\mathbf{x}_i - gV \circ f_W(\mathbf{x}_i)\|_2^2 + \sum_{i=1}^m \Omega(f_W(\mathbf{x}_i))$$

the 2nd also called **sparse autoencoder** [Ranzato et al., 2006].

Quick summary of linear vs nonlinear models



	linear models	nonlinear models
autoencoder	$\min_B \sum_{i=1}^m \ell(\mathbf{x}_i, \mathbf{B}\mathbf{B}^T \mathbf{x}_i)$ $\min_{B,A} \sum_{i=1}^m \ell(\mathbf{x}_i, \mathbf{B}\mathbf{A}^T \mathbf{x}_i)$	$\min_{V,W} \sum_{i=1}^m \ell(\mathbf{x}_i, g_V \circ f_W(\mathbf{x}_i))$
factorization	$\min_{B,Z} \sum_{i=1}^m \ell(\mathbf{x}_i, \mathbf{B}\mathbf{z}_i)$	$\min_{V,Z} \sum_{i=1}^m \ell(\mathbf{x}_i, g_V(\mathbf{z}_i))$
sparse coding	$\min_{B,Z} \sum_{i=1}^m \ell(\mathbf{x}_i, \mathbf{B}\mathbf{z}_i)$ $+\lambda \sum_{i=1}^m \Omega(\mathbf{z}_i)$	$\min_{V,Z} \sum_{i=1}^m \ell(\mathbf{x}_i, g_V(\mathbf{z}_i))$ $+\lambda \sum_{i=1}^m \Omega(\mathbf{z}_i)$ $\min_{V,W} \sum_{i=1}^m \ell(\mathbf{x}_i, g_V \circ f_W(\mathbf{x}_i))$ $+\lambda \sum_{i=1}^m \Omega(f_W(\mathbf{x}_i))$

ℓ can be general loss functions other than $\|\cdot\|_2$

Ω promotes sparsity, e.g., $\Omega = \|\cdot\|_1$

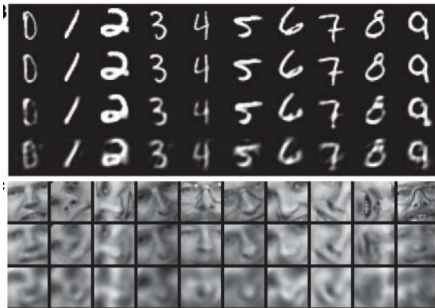
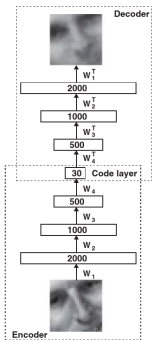
PCA for linear data

Autoencoder: extensions of PCA for nonlinear data

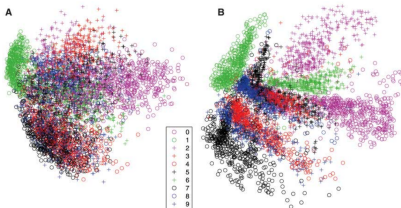
Applications of autoencoder

Self-supervised learning (SSL)

Nonlinear dimension reduction



autoencoder vs. PCA vs. logistic PCA



[Hinton, 2006]

Traditional learning pipeline

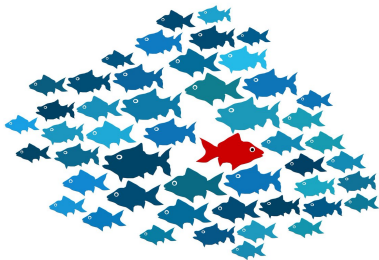


- feature extraction is “independent” of the learning models and tasks
- features are handcrafted and/or learned

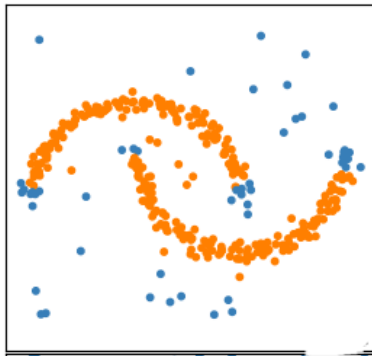
Use the low-dimensional codes as features/representations

- task agnostic
- less overfitting
- semi-supervised (rich unlabeled data + little labeled data) learning

Outlier detection

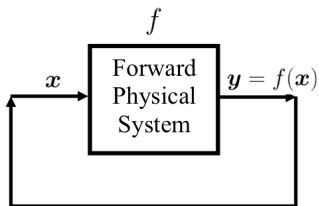


(Credit: towardsdatascience.com)



- idea: outliers don't obey the manifold assumption — the reconstruction error $\ell(\mathbf{x}_i, g_V \circ f_W(\mathbf{x}_i))$ is large after autoencoder training
- for effective detection, better use ℓ that penalizes large errors less harshly than $\|\cdot\|_2^2$, e.g., $\ell(\mathbf{x}_i, g_V \circ f_W(\mathbf{x}_i)) = \|\mathbf{x}_i - g_V \circ f_W(\mathbf{x}_i)\|_2$ [Lai et al., 2019]

Deep generative prior



Inverse Problem f^{-1}

- **inverse problems**: given f and $\mathbf{y} \approx f(\mathbf{x})$, estimate \mathbf{x}
- often ill-posed, i.e., \mathbf{y} doesn't contain enough info for recovery
- **regularized data-fitting** formulation:

$$\min_{\mathbf{x}} \ell(\mathbf{y}, f(\mathbf{x})) + \lambda \Omega(\mathbf{x})$$

where Ω contains extra info about \mathbf{x}

Suppose $\mathbf{x}_1, \dots, \mathbf{x}_m$ come from the same manifold as \mathbf{x}

- train a deep factorization model on $\mathbf{x}_1, \dots, \mathbf{x}_m$:

$$\min_{\mathbf{V}, \mathbf{Z}} \sum_{i=1}^m \ell(\mathbf{x}_i, g_{\mathbf{V}}(\mathbf{z}_i))$$

- $\mathbf{x} \approx g_{\mathbf{V}}(\mathbf{z})$ for a certain \mathbf{z} so: $\min_{\mathbf{z}} \ell(\mathbf{y}, f \circ g_{\mathbf{V}}(\mathbf{z}))$. Some recent work even uses random \mathbf{V} , i.e., without training

See: [Pan et al., 2020, Ulyanov et al., 2018, Bora et al., 2017,

Wang et al., 2021, Zhuang et al., 2022]

PCA for linear data

Autoencoder: extensions of PCA for nonlinear data

Applications of autoencoder

Self-supervised learning (SSL)

SSL: marriage of supervised and unsupervised learning

Why not supervised learning?

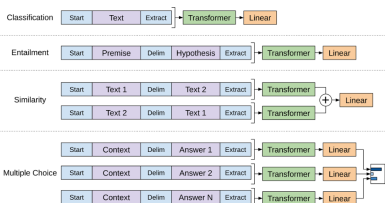
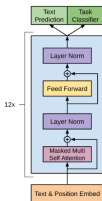
- labeling is expensive
- unlabeled data can be abundant
- supervised learning are task-specific (despite transfer learning)

What's self-supervised learning?

- like unsupervised learning: no task-specific labels
- like supervised learning: trained on tasks defined on the unlabeled data—**pretext tasks**

Pretraining: (transformer-based)
decoder-only architectures pretrained
on language model $p(x^{(t+1)} | x^{(0)}, \dots, x^{(t)})$

Finetuning: on task-specific
supervised data

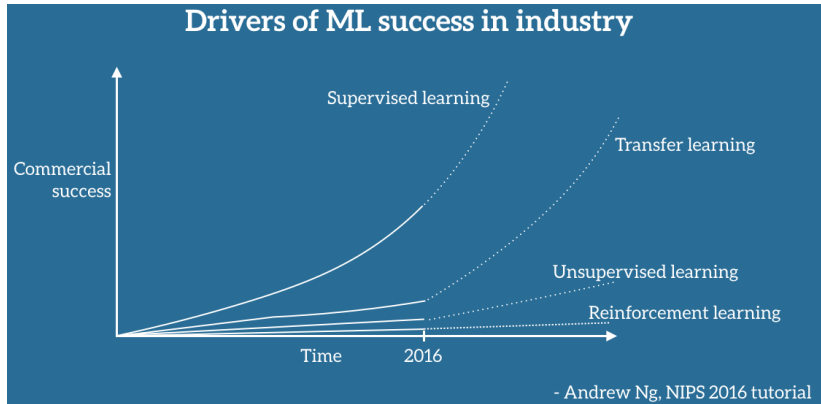


How Much Information is the Machine Given during Learning?

- ▶ **“Pure” Reinforcement Learning (cherry)**
 - ▶ The machine predicts a scalar reward given once in a while.
 - ▶ **A few bits for some samples**
- ▶ **Supervised Learning (icing)**
 - ▶ The machine predicts a category or a few numbers for each input
 - ▶ Predicting human-supplied data
 - ▶ **10→10,000 bits per sample**
- ▶ **Self-Supervised Learning (cake génoise)**
 - ▶ The machine predicts any part of its input for any observed part.
 - ▶ Predicts future frames in videos
 - ▶ **Millions of bits per sample**



Success of different ML strategies



SSL in NLP: sequential prediction

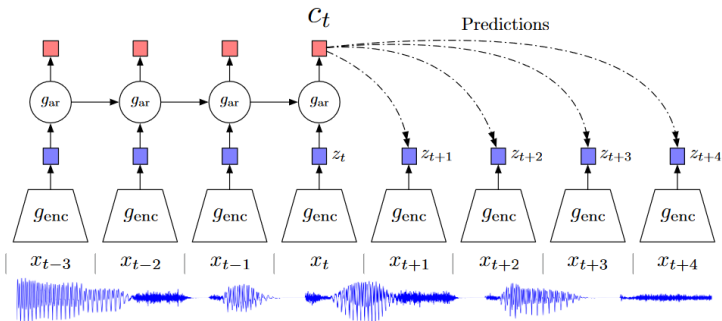
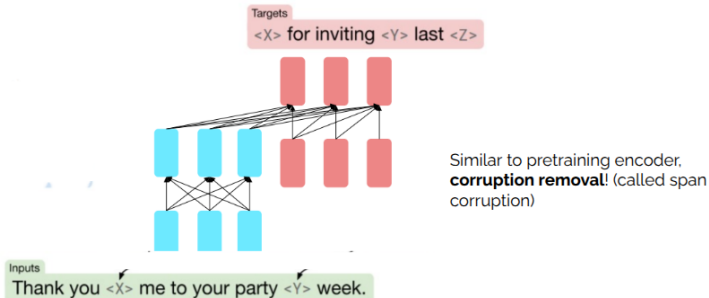


Image credit: <https://www.v7labs.com/blog/self-supervised-learning-guide>

Language modeling is a special case, e.g., causal language modeling. Used in e.g., GPT, LLaMA

SSL in NLP: masked token prediction



Force the model to understand the context both before and after. Used in e.g., BERT, RoBERTa

SSL: contrastive learning in CV

learning embedding/representation that respects certain predefined constraints/goals

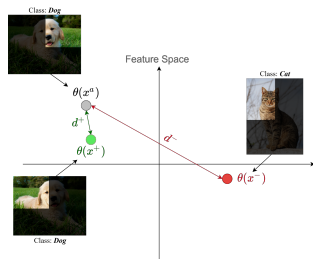
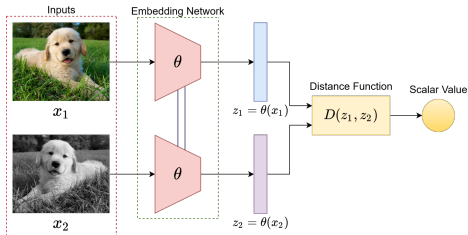
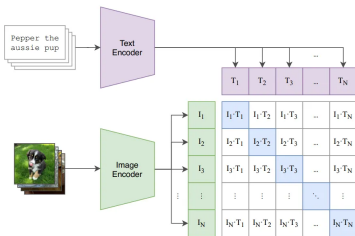


Image credit: <https://www.v7labs.com/blog/self-supervised-learning-guide> link to DINO:

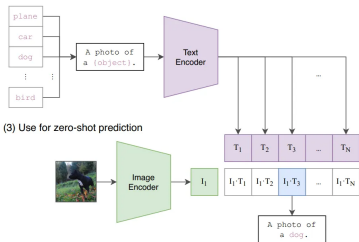
<https://github.com/facebookresearch/dino>

SSL: contrastive learning in CV/NLP

(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

Align/contrast different modalities—better called *weakly supervised learning*.

Used in e.g., CLIP <https://arxiv.org/abs/2103.00020>

SSL: contrastive learning + missing prediction in CV

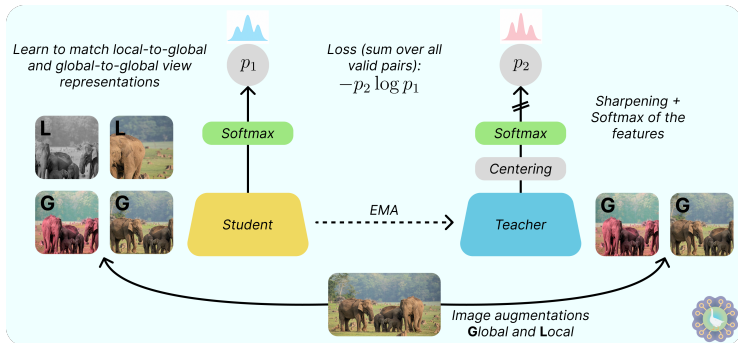


Image credit: <https://mlhonk.substack.com/p/39-from-dino-to-dinov3>

Meta's DINO models <https://ai.meta.com/research/dinov3/>

More about self-supervised learning

- Awesome Self-Supervised Learning
<https://github.com/jason718/awesome-self-supervised-learning>
- A Cookbook of Self-Supervised Learning
<https://arxiv.org/abs/2304.12210>
- Know Your Self-supervised Learning: A Survey on Image-based Generative and Discriminative Training <https://arxiv.org/abs/2305.13689>
- https://cs229.stanford.edu/notes2021spring/notes2021spring/cs229_lecture_selfsupervision_final.pdf
- Self-Supervised Representation Learning <https://lilianweng.github.io/posts/2019-11-10-self-supervised/>

- [Bai et al., 2018] Bai, Y., Jiang, Q., and Sun, J. (2018). **Subgradient descent learns orthogonal dictionaries.** *arXiv:1810.10702*.
- [Baldi and Hornik, 1989] Baldi, P. and Hornik, K. (1989). **Neural networks and principal component analysis: Learning from examples without local minima.** *Neural Networks*, 2(1):53–58.
- [Belkin and Niyogi, 2001] Belkin, M. and Niyogi, P. (2001). **Laplacian eigenmaps and spectral techniques for embedding and clustering.** In Dietterich, T. G., Becker, S., and Ghahramani, Z., editors, *Advances in Neural Information Processing Systems 14 [Neural Information Processing Systems: Natural and Synthetic, NIPS 2001, December 3-8, 2001, Vancouver, British Columbia, Canada]*, pages 585–591. MIT Press.
- [Bojanowski et al., 2017] Bojanowski, P., Joulin, A., Lopez-Paz, D., and Szlam, A. (2017). **Optimizing the latent space of generative networks.** *arXiv:1707.05776*.
- [Bora et al., 2017] Bora, A., Jalal, A., Price, E., and Dimakis, A. G. (2017). **Compressed sensing using generative models.** In *Proceedings of the 34th International Conference on Machine Learning*, volume 70.

- [Bourlard and Kamp, 1988] Bourlard, H. and Kamp, Y. (1988). **Auto-association by multilayer perceptrons and singular value decomposition.** *Biological Cybernetics*, 59(4-5):291–294.
- [Elgendy, 2020] Elgendy, M. (2020). **Deep Learning for Vision Systems.** MANNING PUBN.
- [Fan and Cheng, 2018] Fan, J. and Cheng, J. (2018). **Matrix completion by deep matrix factorization.** *Neural Networks*, 98:34–41.
- [Hinton, 2006] Hinton, G. E. (2006). **Reducing the dimensionality of data with neural networks.** *Science*, 313(5786):504–507.
- [Hinton et al., 2006] Hinton, G. E., Osindero, S., and Teh, Y.-W. (2006). **A fast learning algorithm for deep belief nets.** *Neural Computation*, 18(7):1527–1554.
- [Hinton and Zemel, 1994] Hinton, G. E. and Zemel, R. S. (1994). **Autoencoders, minimum description length and helmholtz free energy.** In *Advances in neural information processing systems*, pages 3–10.
- [Kawaguchi, 2016] Kawaguchi, K. (2016). **Deep learning without poor local minima.** *arXiv:1605.07110*.

- [Lai et al., 2019] Lai, C.-H., Zou, D., and Lerman, G. (2019). **Robust subspace recovery layer for unsupervised anomaly detection.** *arXiv:1904.00152*.
- [Li et al., 2020a] Li, S., Li, Q., Zhu, Z., Tang, G., and Wakin, M. B. (2020a). **The global geometry of centralized and distributed low-rank matrix recovery without regularization.** *arXiv:2003.10981*.
- [Li et al., 2020b] Li, T., Mehta, R., Qian, Z., and Sun, J. (2020b). **Rethink autoencoders: Robust manifold learning.** *ICML workshop on Uncertainty and Robustness in Deep Learning*.
- [Lu and Kawaguchi, 2017] Lu, H. and Kawaguchi, K. (2017). **Depth creates no bad local minima.** *arXiv:1702.08580*.
- [Mairal, 2014] Mairal, J. (2014). **Sparse modeling for image and vision processing.** *Foundations and Trends® in Computer Graphics and Vision*, 8(2-3):85–283.
- [Olshausen and Field, 1996] Olshausen, B. A. and Field, D. J. (1996). **Emergence of simple-cell receptive field properties by learning a sparse code for natural images.** *Nature*, 381(6583):607–609.

- [Pan et al., 2020] Pan, X., Zhan, X., Dai, B., Lin, D., Loy, C. C., and Luo, P. (2020). **Exploiting deep generative prior for versatile image restoration and manipulation.** In *Computer Vision – ECCV 2020*, pages 262–277. Springer International Publishing.
- [Park et al., 2019] Park, J. J., Florence, P., Straub, J., Newcombe, R., and Lovegrove, S. (2019). **DeepSDF: Learning continuous signed distance functions for shape representation.** In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 165–174. IEEE.
- [Qu et al., 2019] Qu, Q., Zhai, Y., Li, X., Zhang, Y., and Zhu, Z. (2019). **Analysis of the optimization landscapes for overcomplete representation learning.** *arXiv:1912.02427*.
- [Ranzato et al., 2006] Ranzato, M., Poultney, C. S., Chopra, S., and LeCun, Y. (2006). **Efficient learning of sparse representations with an energy-based model.** In *Advances in Neural Information Processing Systems*.
- [Roweis, 2000] Roweis, S. T. (2000). **Nonlinear dimensionality reduction by locally linear embedding.** *Science*, 290(5500):2323–2326.

- [Sun et al., 2017] Sun, J., Qu, Q., and Wright, J. (2017). **Complete dictionary recovery over the sphere i: Overview and the geometric picture.** *IEEE Transactions on Information Theory*, 63(2):853–884.
- [Tan and Mayrovouniotis, 1995] Tan, S. and Mayrovouniotis, M. L. (1995). **Reducing data dimensionality through optimizing neural network inputs.** *AIChE Journal*, 41(6):1471–1480.
- [Tenenbaum, 2000] Tenenbaum, J. B. (2000). **A global geometric framework for nonlinear dimensionality reduction.** *Science*, 290(5500):2319–2323.
- [Ulyanov et al., 2018] Ulyanov, D., Vedaldi, A., and Lempitsky, V. (2018). **Deep image prior.** In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9446–9454.
- [Valavi et al., 2020] Valavi, H., Liu, S., and Ramadge, P. J. (2020). **The landscape of matrix factorization revisited.** *arXiv:2002.12795*.
- [van der Maaten and Hinton, 2008] van der Maaten, L. and Hinton, G. (2008). **Visualizing data using t-sne.** *Journal of Machine Learning Research*, 9:2579–2605.
- [Wang et al., 2021] Wang, H., Li, T., Zhuang, Z., Chen, T., Liang, H., and Sun, J. (2021). **Early stopping for deep image prior.** *arXiv:2112.06074*.

[Zhuang et al., 2022] Zhuang, Z., Li, T., Wang, H., and Sun, J. (2022). **Blind image deblurring with unknown kernel size and substantial noise.** *arXiv:2208.09483*.