# EECS 182/282A

**Today:** Self-supervision
State-space models

Reading: Prince through Ch 11 + Ch 13
Note: Little in this lecture is in
the Prince textbook. (Just 9.3.7)

## Architecture Order In Class:

MLPs $\longrightarrow$ CNNs $\longrightarrow$ Graph NN $\longrightarrow$ RNN/state-space $\longrightarrow$ Transformers

## General Principle: Self-supervision

"I need labels to train, I don't have labels, so make my own labels from data"

Lessons from example of : 1) We can learn a partial pattern that can be useful
trying to learn what you can       2) Might need scaffolding parts of my NN.
in the Kalman Filtering context    3) Generic idea of "next-thing" prediction
                                       in causal sequence modeling

## Step Back: Connect to unsupervised learning in classic ML.

Two approaches: 1) Dimensionality Reduction
2) Clustering

## Recall Dimensionality Reduction: Think about PCA

All we have are $\{\vec{x}_i\}_{i=1}^N \leftarrow$ d-dim    Unlabeled Data From Interesting Distribution

Classic Recipe (Neglecting Means):
1) Construct
$$X = \begin{bmatrix} \vec{x}_1^T \\ \vec{x}_2^T \\ \vdots \\ \vec{x}_N^T \end{bmatrix}$$

2) Compute SVD $X = U \Sigma V^T = \sum_i \sigma_i \vec{u}_i \vec{v}_i^T$

3) Keep top k    $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_k$    singular vectors to use for dim-reduction

4) Given some other problem $\vec{x} \longrightarrow \begin{bmatrix} \vec{v}_1^T \vec{x} \\ \vdots \\ \vec{v}_k^T \vec{x} \end{bmatrix}$ } k-dim Features

Classic Perspective has no loss, no labels, no gradients, no optimizer.
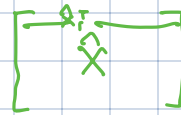No mini-batches, ...

Why was this reasonable?
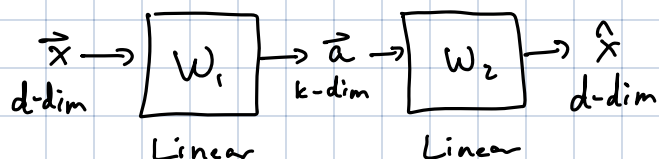
Recall Eckhart-Young-Mirsky Theorem for Frobenius Norm

Given $X$, $\hat{X} = \sum_{i=1}^{k} \sigma_i \vec{u}_i \vec{v}_i^T$ is the rank-$k$ matrix that minimizes

$$\| X - \hat{X} \|_F^2$$

$$\left[ \overbrace{\phantom{--}}^{x_i^T} \; X \right] \qquad \left[ \overbrace{\phantom{--}}^{\hat{x}_i^T} \; \hat{X} \right]$$

Let's interpret this in neural-net terms: (row by row perspective on $\hat{X}$)

$\vec{x} \longrightarrow$ | $W_1$ | $\to \vec{a} \to$ | $W_2$ | $\to \hat{x}$   to minimize $\| \vec{x} - \hat{x} \|_2^2$
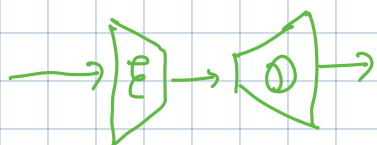d-dim       k-dim      d-dim
     Linear        Linear

Trained over the entire data set $\{\vec{x}_i\}_{i=1}^{N}$

If this reaches the minimizer, it has to find the same subspace to project into using $W_1$ as the first $k$ s.v. $\vec{v}_1, \ldots, \vec{v}_k$.

This approach has a name: <u>Auto-encoders</u>

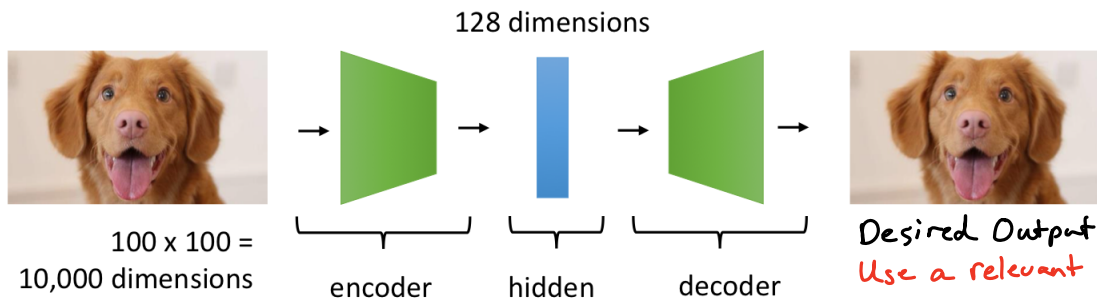Core Ingreddents: Labels are $\vec{x}_i$ itself.
Architecture has an encoder followed by a decoder
Bottleneck in the middle.

$\longrightarrow$ E $\to$ D $\longrightarrow$

Traditional Perspective: Decoder is scaffolding.

<u>Some Pictures...</u>

128 dimensions



100 x 100 =
10,000 dimensions

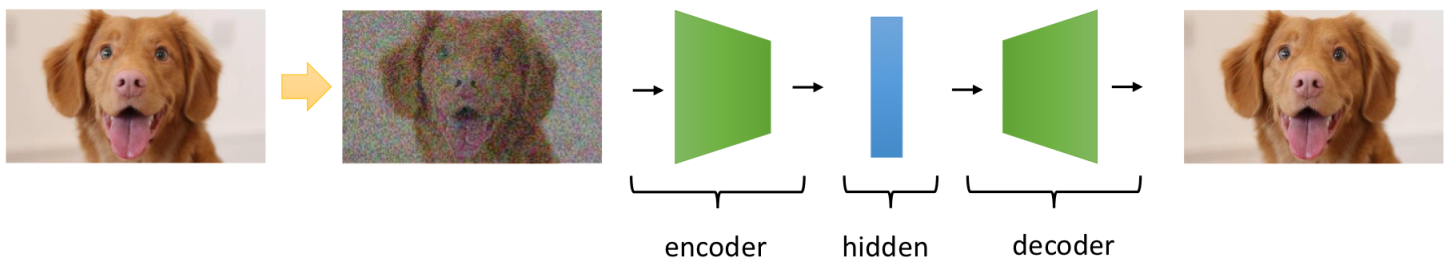encoder        hidden        decoder

Desired Output
Use a relevant loss function to train

Classic Autoencoder... Into the encoder & decoder boxes goes a Neural Net Architecture choice tied to domain. E.g. CNN for image.
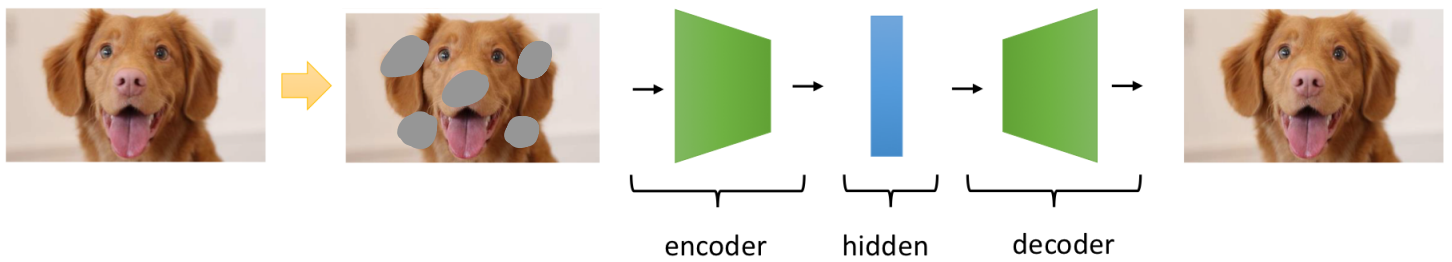
Aux Loss that is sparsity-seeking e.g $\|\cdot\|$, Penalty



100 x 100 =
10,000 dimensions

encoder        hidden        decoder

Desired Output
Standard end-to-end loss

Sparse Autoencoder : Can make hidden "bottleneck" even bigger than the original but impose an auxilliary loss.



encoder        hidden        decoder

Denoising Autoencoder: Instead of regularizing with a bottleneck or aux loss, use data augmentation (e.g. adding noise)

encoder     hidden     decoder

Masked Autoencoders: Same spirit as denoising autoencoders, but the data augmentation is random masking

Note: Can involve learned $\vec{?}$ vectors if this is appropriate.

Note from O.H. In our default perspective, reconstruction loss is on everything. But there's a variant when it is only on masked places.

Recall Perspective on Traditional Unsupervised learning in standard ML
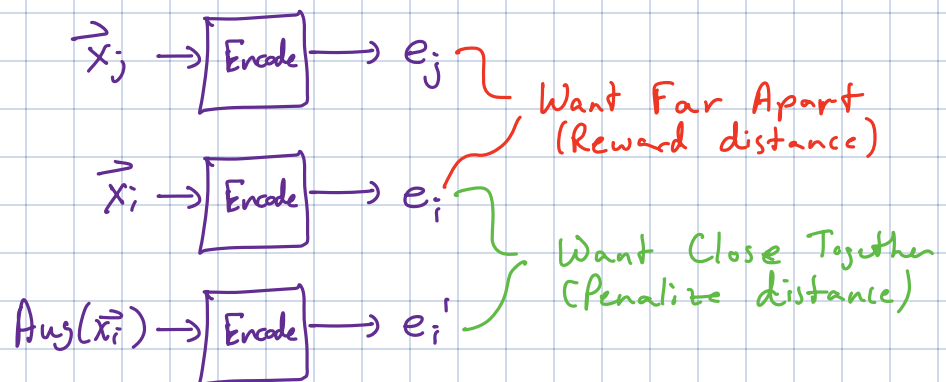Two approaches: 1) Dimensionality Reduction $\Longleftrightarrow$ Autoencoder-style Self-Supervision
                   2) Clustering $\Longleftrightarrow$ ???

## Contrastive Self-Supervision

Core Idea: An example and its own augmentations should be in the same cluster, while Fundamentally different examples are in different clusters.

(First formalized in "Provable Guarantees for Self-Supervised Deep Learning with Spectral Contrastive Loss" by HaoChen, et.al. in 2021 but builds on latent intuition in "A theoretical analysis of Contrastive Unsupervised Representation Learning" by Arora, et.al. in 2019)

Most basic structure: Embed/Encode 3 things

$$\vec{x}_j \rightarrow \boxed{\text{Encode}} \rightarrow e_j$$

Want Far Apart
(Reward distance)

$$\vec{x}_i \rightarrow \boxed{\text{Encode}} \rightarrow e_i$$

Want Close Together
(Penalize distance)

$$Aug(\vec{x}_i) \rightarrow \boxed{\text{Encode}} \rightarrow e_i{}'$$

Fun going beyond.... It's possible to avoid negative examples entirely, but this requires other tricks to prevent collapse...