

## UNIT-II

### Questions:

- 1) Principal Component Analysis / Diff b/w PCA & Autoencoder
- 2) Autoencoder Architecture.
- 3) Denoising of Autoencoder for
- 4) Sparsity in Autoencoder
- 5) Lower-Dimensional Representation
- 6) Diff b/w AE & PCA.
- 7)
- 8)
- 9)
- 10)
- 11)
- 12)
- 13)
- 14)
- 15)

## Chapter 2:-

### 1) Learning Lower Dimensional Representation

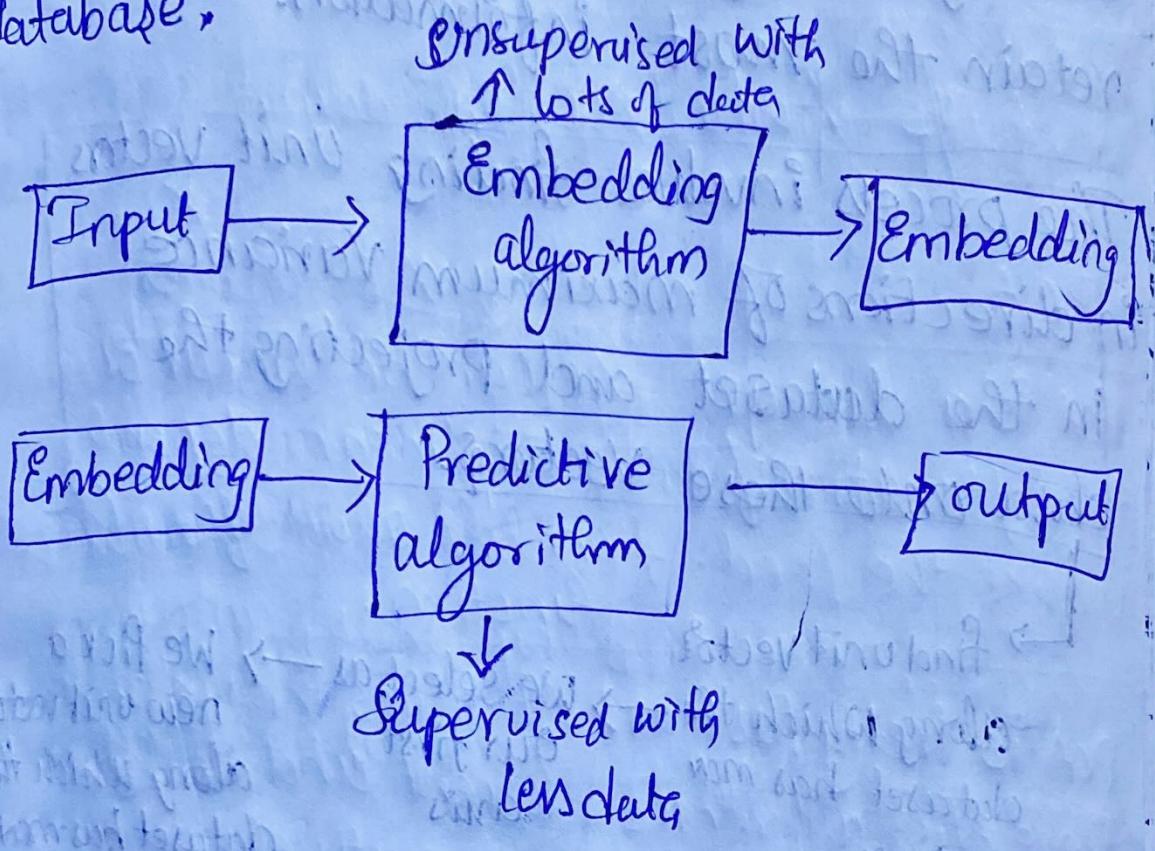
Large input vectors leads to larger models with many parameters, which can be Expressive but are also data-hungry and Prone to Overfitting.

The Convolutional architectures address the curse of dimensionality by Reducing the number of Parameters, maintaining Expressiveness without need for Extensive training data.

There is often a Scarcity of labeled data, which is necessary for training Convolutional networks, leading to the needs for methods that can work effectively with less labeled data.

One solution is to develop models that learn from unlabeled data by creating lower-dimensional embeddings in an unsupervised manner.

These embeddings capture essential features and can be used with smaller models, reducing the requirements for large labeled database.



## Principal Component Analysis:

Principal Component Analysis is a statistical procedure that aims to transform a set of possibly correlated variables into a set of linearly uncorrelated variables known as principal components.

The goal of PCA is to retain as much seek to reduce the dimensionality of data by finding the new ones that retain the most information.

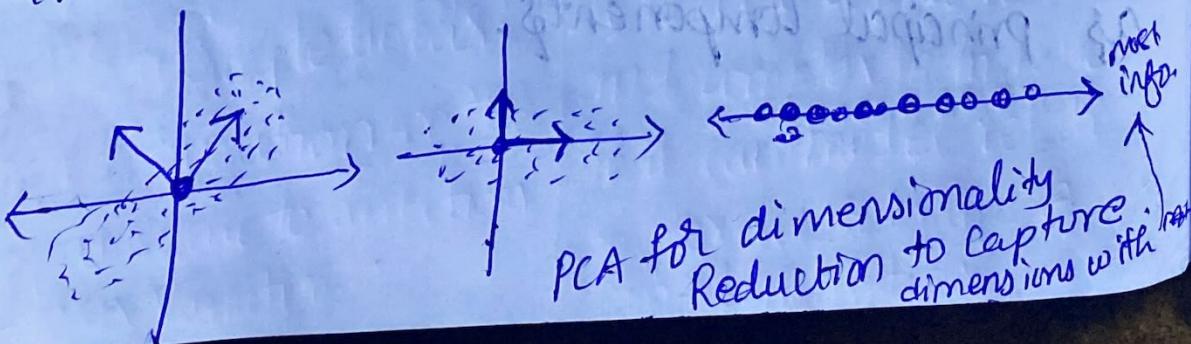
The process involves finding unit vectors in directions of maximum variance in the dataset and projecting the data onto these vectors.

→ find unit vector along which the dataset has max variance → we select as our first axis → we pick a new unit vector along which the dataset has max variance.

Continue this process until we have found a total of d new vectors. This is our second axis.

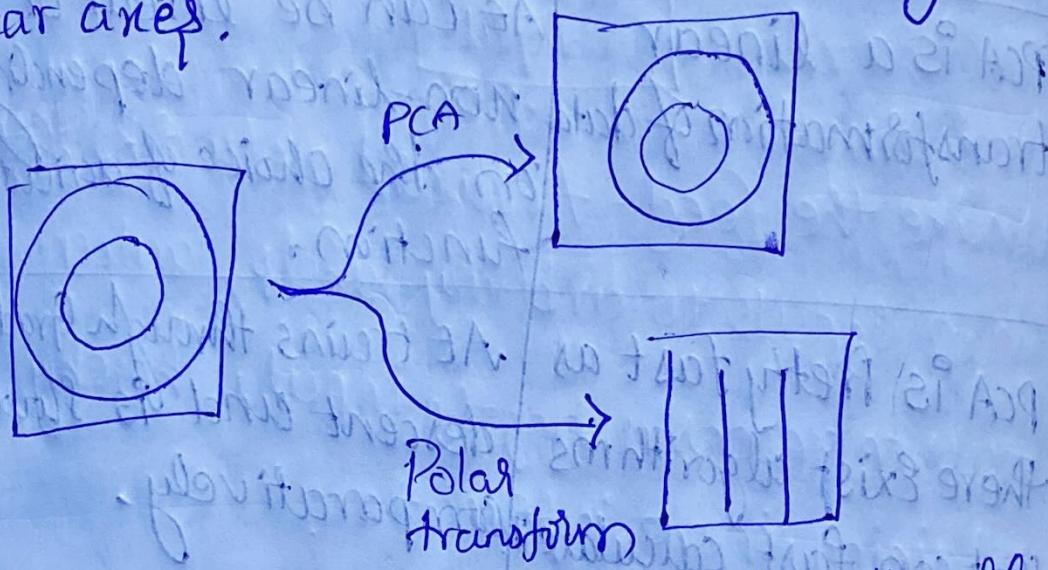
Project our data onto this new set of axes.

This result is shown in below fig.



The dataset is represented as a matrix, and an embedding matrix is created that reduces dimensions while preserving as much info as possible.

PCA can fail with complex datasets, particularly where important patterns are nonlinear, as it primarily capture variance along linear axes.



A situation in which PCA fails to optimally transform the data for dimensionality Reducing

The non-linear Relationships in data, such as the above concentric circles, may require diff approaches, like polar transformation, to Separate features Effectively.

The limitations of PCA in dealing with non linear relationships highlight the needs

for nonlinear dimensionality reduction techniques, which are addressed by neural models in deep learning.

## PCA vs Auto Encoders

### PCA

PCA is a linear transformation of data

PCA is pretty fast as there exist algorithms that can fast calculate.

PCA Projects data into dimensions that are orthogonal to each other resulting in very low or close to zero correlation in the projected data.

PCA is a simple linear transformation on the input space to directions of max variations

### AE

AE can be linear or non-linear depending on the choice of activation function.

AE trains through Gradient descent and is slower comparatively.

AE transformed data doesn't guarantee that because the way it's trained is merely to minimize the reconstruction loss.

AE is more sophisticated and complex technique that can model relatively complex relationships & non-linearities.

one rule to thumb could be  
the size of Data. PCA for  
small datasets.

AE for larger  
dataset.

PCA hyperparameter is  
"K"

AE it is the architecture  
of the neural network.

AE with single layer.  
& linear activation  
has similar performance  
as PCA.

AE with multiple layers  
and non-activation function  
Prone to Overfitting.

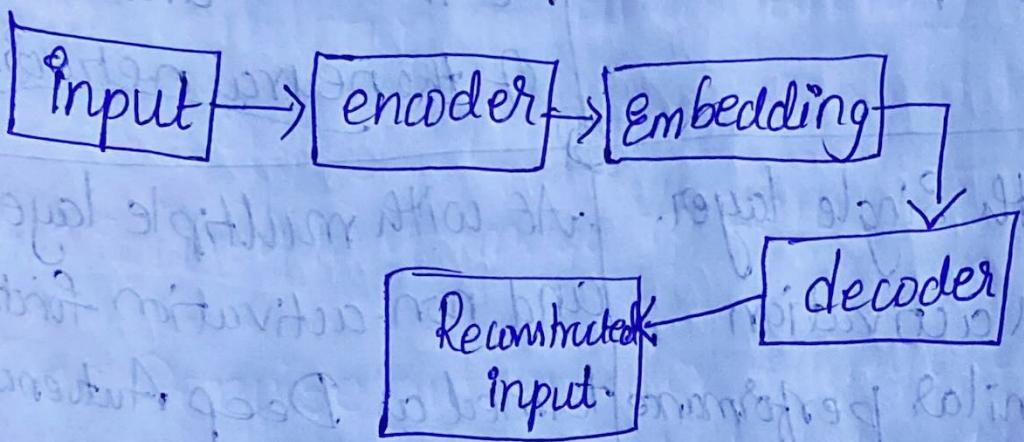
## Autoencoder Architecture:

The feed-forward networks learn increasingly refined representation of input data, with the final convolutional layer output serving as a compressed representation.

Autoencoders are introduced as an architecture designed to address the loss of information by capturing and compressing input data into a lower-dimensional space.

An autoencoder consists of two main components:

an encoder that compresses the input into a code, and a decoder that reconstructs the original input from this code.

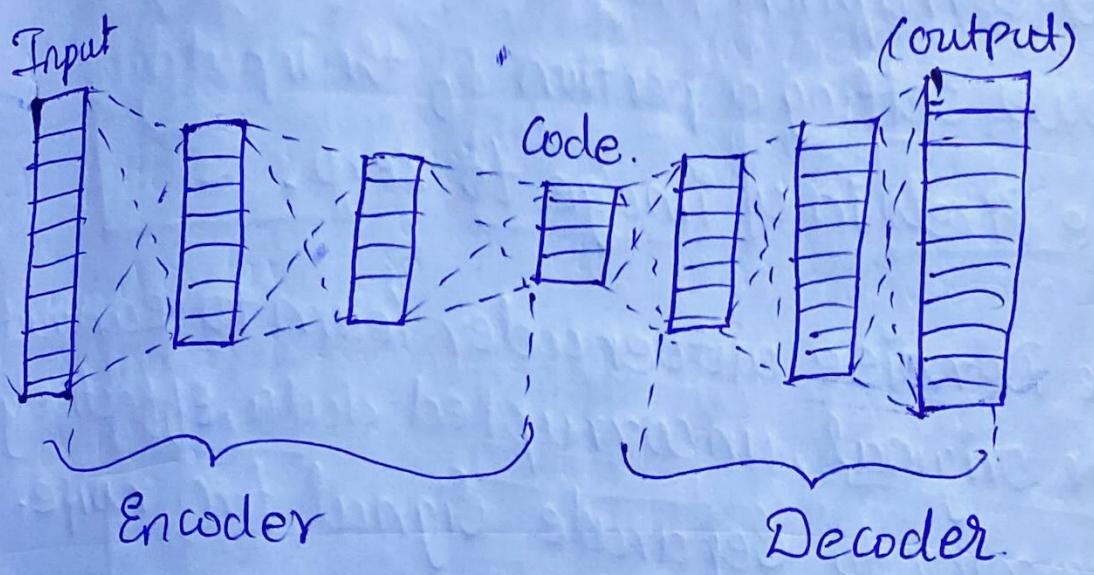


Unlike the feed-forward networks that map inputs to arbitrary labels, autoencoders focus on reconstructing the input, thereby retaining more information.

The encoder part of the autoencoder learns to identify and encode the most critical features of the input data.

The decoder part learns to reverse the encoding process to reconstruct the input data from the encoded representation.

Autoencoders can be used to generate lower dimensional data representation without the supervision required in traditional methods.



## Denoising of autoencoder

### Denoising:

Denoising is used to enhance autoencoders by generating noise-resistant Embedding.

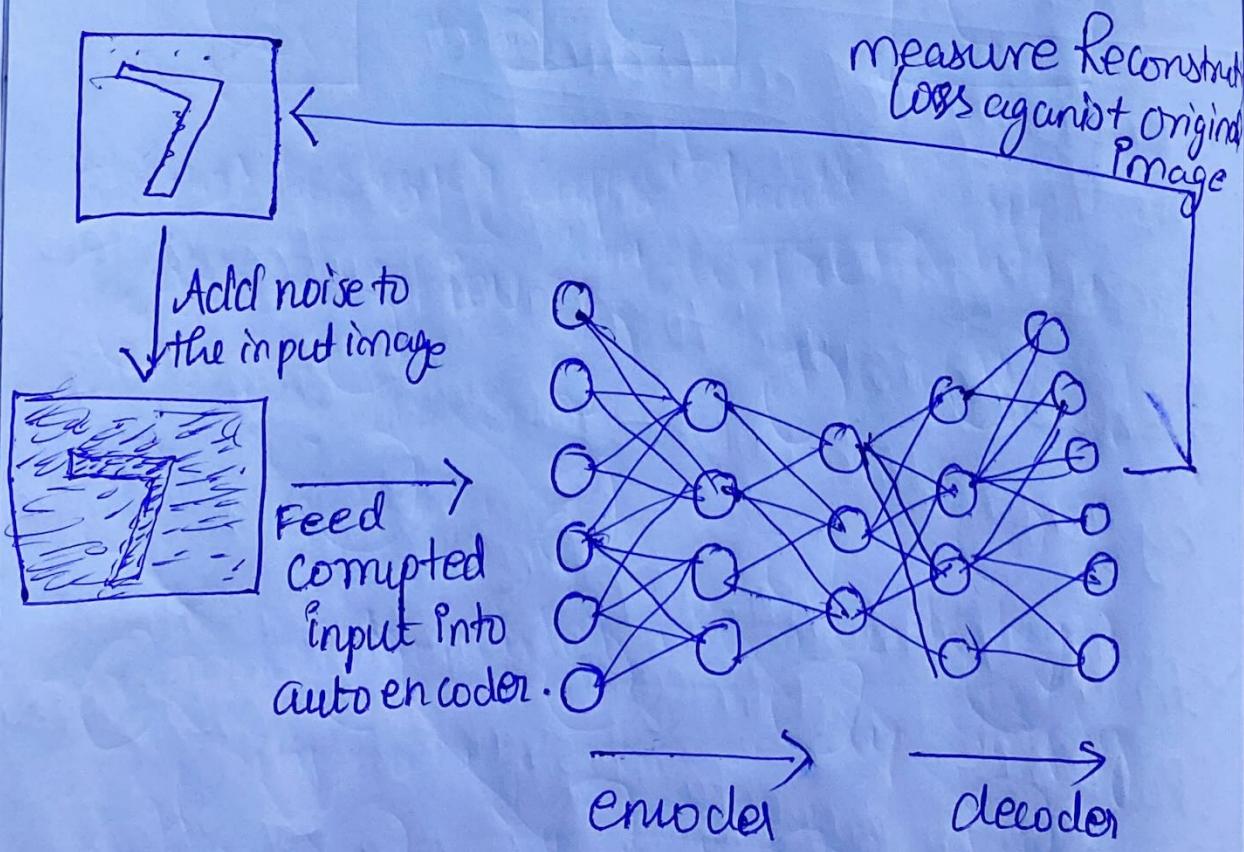
In the human Perception we can recognize Patterns and Images despite significant noise, a property desirable for ML models.

We can even distinguish  $\alpha$  and  $\beta$  if all the pixels got corrupted easily.

The Denoising encoders are trained using corrupted inputs, forcing them to learn robust features.

The corruption process in denoising involves setting a portion of the input data, like pixels in an image to zero.

The Denoise autoencoder learns to reconstruct the original, uncorrupted data, effectively learning to separate signal from noise.



The geometric interpretation involves the concept of a manifold, the shape within the data that the autoencoder learns to reconstruct.

An autoencoder must identify whether a data point belongs to one manifold or another during the reconstruction process.

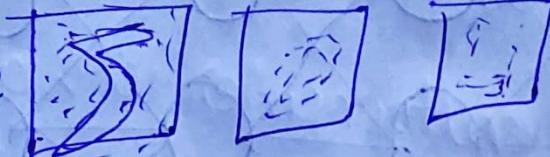
The denoising process conceptually expands the data space to include points around the manifold, with the autoencoder then collapsing these points back onto the manifold.

The denoising objective helps the autoencoder learn to distinguish b/w generalizable features and noise.

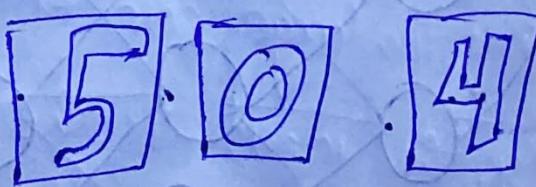
Original Image



Corrupted Image

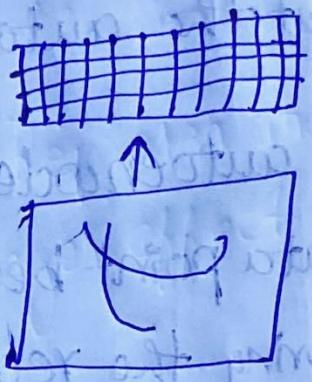
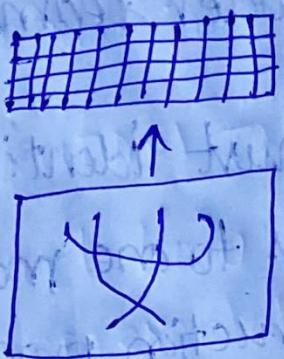
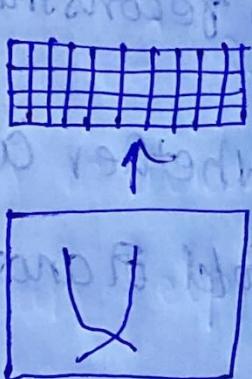


Reconstructed Image

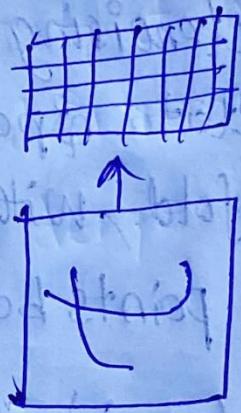
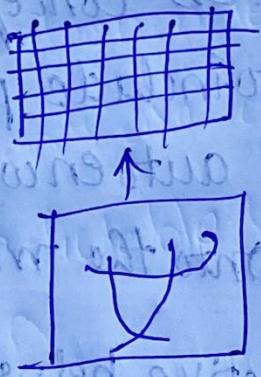
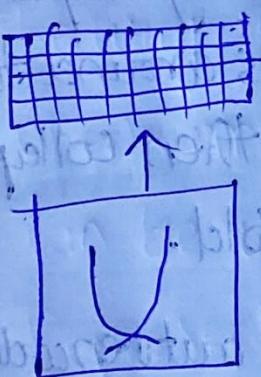


# Sparsity in Autoencoders

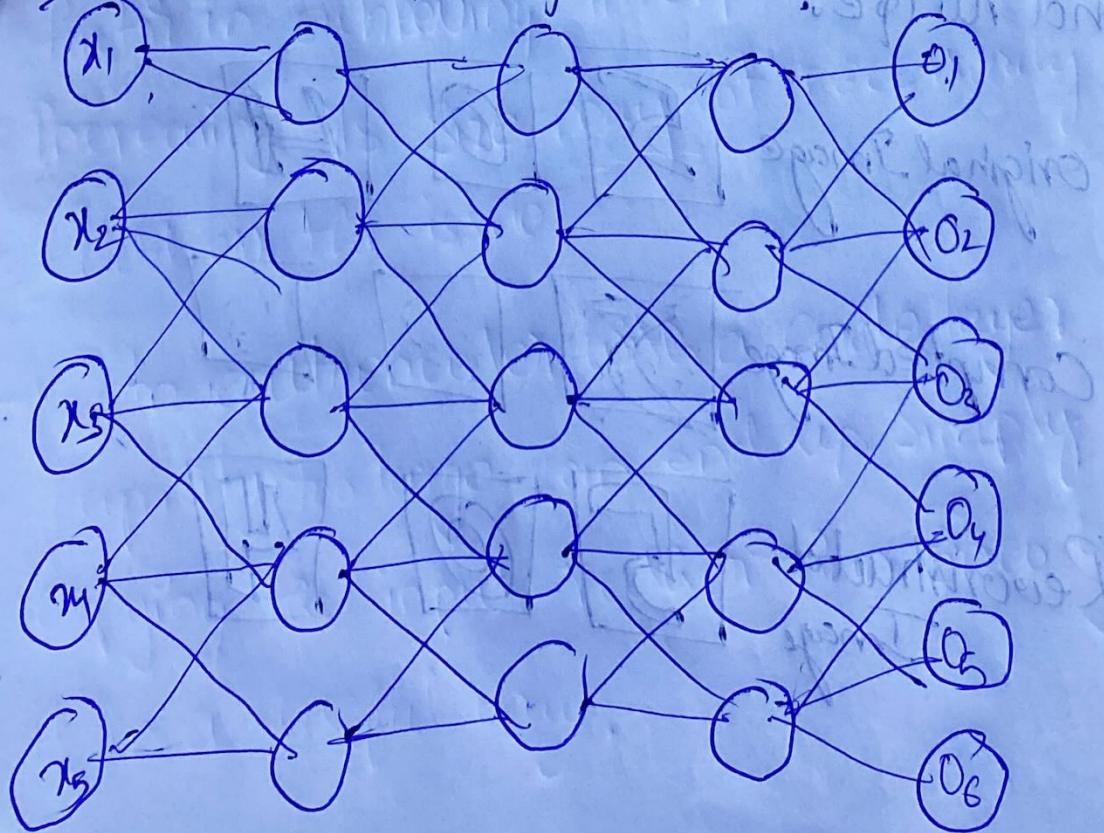
A



B



Input layer  $\xrightarrow{\text{encoder}}$  hidden layer  $\xrightarrow{\text{decoder}}$  output layer



Deep learning models struggle with interpretability due to their complexity and nonlinearity. Deep learning models are generally very difficult to interpret because of the nonlinearities and massive number of parameters that make up a model.

The interpretability can be solved by autoencoder. In general, an autoencoder representations are dense, which means that many features are combined in ways that are difficult to separate or interpret.

In fig A. The activation of a dense representation combine and overlay information from multiple features in ways that are difficult to interpret.

Introducing Sparsity into autoencoder representation can help make the model output more interpretable.

The goal is to achieve a one-to-one correspondence b/w the high level features in the input

data and the components in the autoencoder code.

→ The Fig B. colour - coding the contribution of individual features to the representation can help understand how changes effect the overall representation.

→ The way an autoencoder's output changes when components are added or removed can be unpredictable with dense representation.

→ The code's layer capacity can be limiting factor in achieving interpretability, increasing its size may not necessarily address the issue.

→ An autoencoder with a code layer capacity too large can end up simply copying the input rather than encoding it in a useful way.

→ The Autoencoder can effectively summarize data points when all relevant info is present in the individual data points themselves.

---