

# **Introduction to Autoencoders**

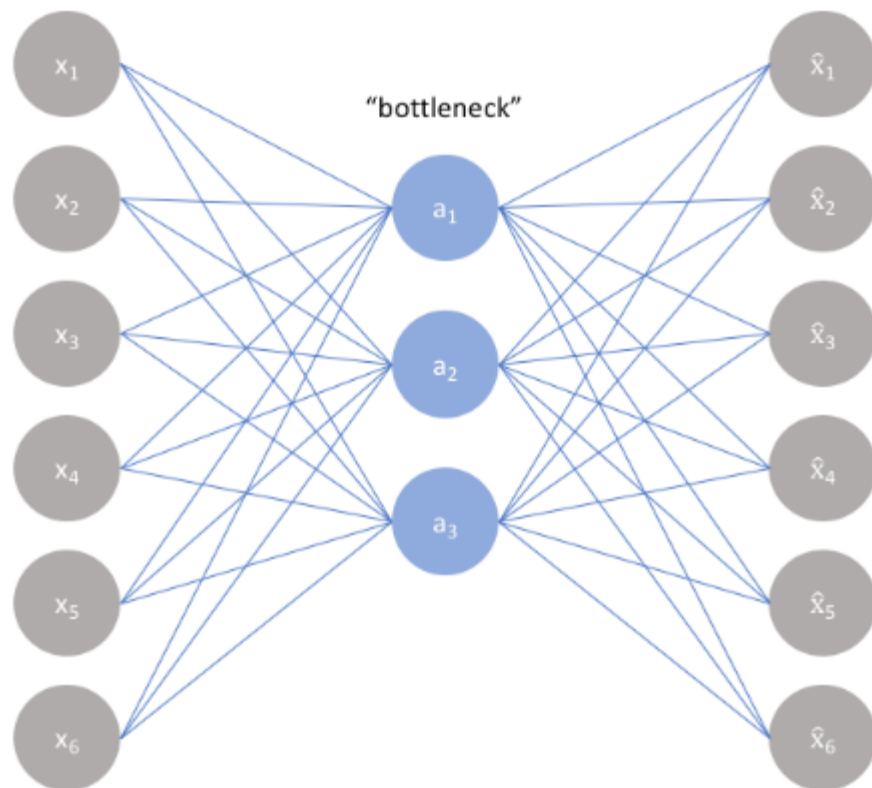
- Autoencoders are an unsupervised learning technique
- Autoencoders leverage neural networks for representation learning.
- Autoencoders introduce a bottleneck which forces a compressed knowledge representation of the original input.
- The idea is to capture a structure in the data if one exists in the bottleneck. Once identified, the underlying structure can be used for other things like anomaly detection.

- First, we take an unlabeled dataset and attempt to reconstruct the original ( $x$ ) by obtaining its reconstruction( $\hat{x}$ ). This network can be trained by minimizing the **reconstruction error**.
- As the name implies, the reconstruction error represents the difference between the original and the reconstruction.
- Note that the reconstruction occurs through a **bottleneck**.
- A bottleneck constrains the amount of information that can traverse the full network, **forcing a learned compression of the input data**.

Input layer

Hidden layer

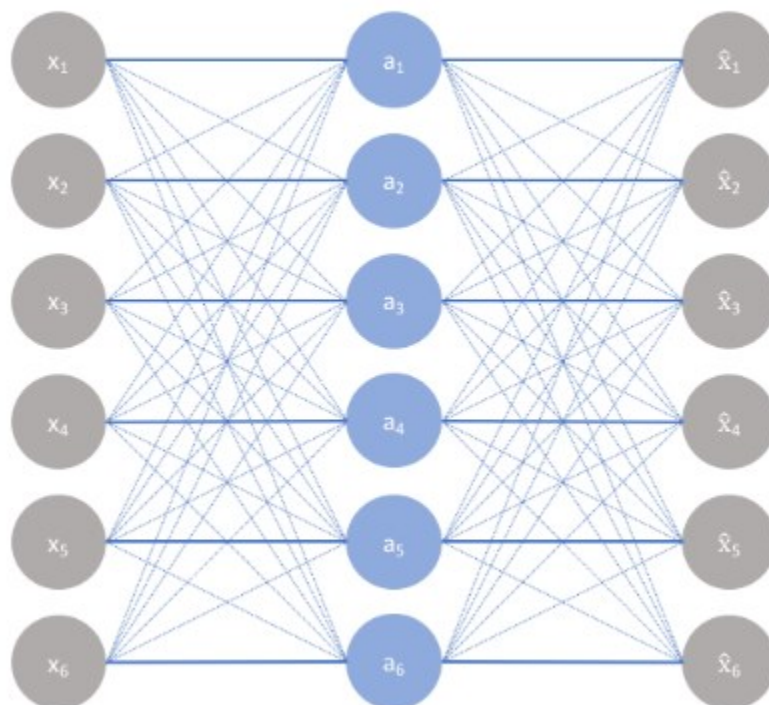
Output layer



Input layer

Hidden layer

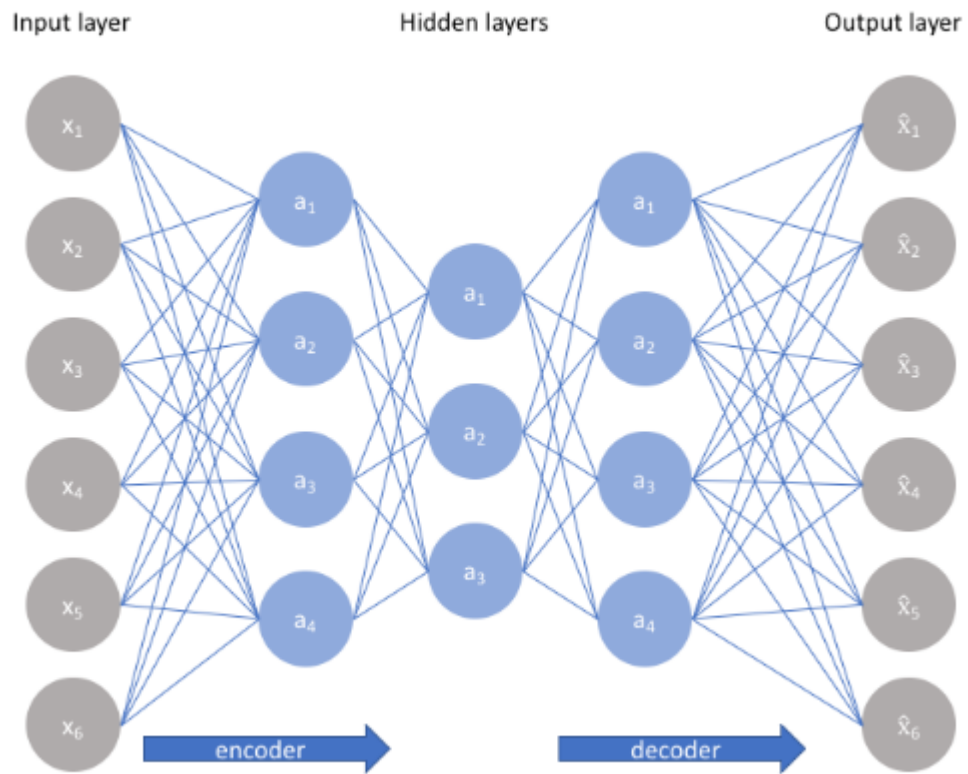
Output layer



The ideal autoencoder manages a trade-off between two things

- Sensitivity to inputs to be able to create a reconstruction
- Insensitivity to inputs to avoid memorizing the training data.
- This involves constructing a loss function which is a representation error and a regularizer

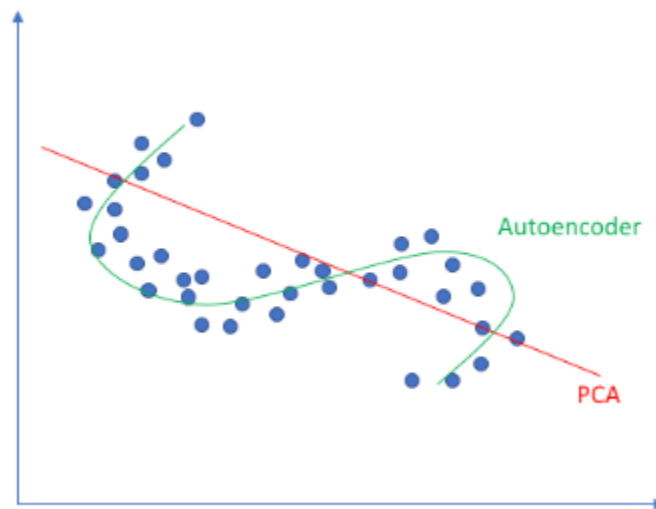
$$\mathcal{L}(x, \hat{x}) + \text{regularizer}$$



- Because neural networks are capable of learning nonlinear relationships, this can be thought of as a more powerful (nonlinear) generalization of PCA.
- Whereas PCA attempts to discover a lower dimensional hyperplane which describes the original data, autoencoders are capable of learning nonlinear manifolds (a manifold is defined in simple terms as a continuous, non-intersecting surface). The difference between these two approaches is visualized below.



## Linear vs nonlinear dimensionality reduction



## Sparse autoencoders

- Sparse autoencoders offer us an alternative method for introducing an information bottleneck without requiring a reduction in the number of nodes at our hidden layers.
- Instead of the bottleneck layer, we will introduce an artificial scarcity by learning encoding and decoding from a small number of neurons(spare autoencoding).
- In the diagram of a generic sparse autoencoder, the opacity of a node corresponds with the level of activation.

