



Applications of Autoencoders

1. Encoder

The encoder is the part of the network that takes the input data and compresses it into a smaller lower-dimensional representation.

- **Input Layer:** This is where the original data enters the network e.g., an image or a set of features.
- **Hidden Layers:** These layers apply transformations to the input data. The encoder's goal is to extract essential features and reduce the data's dimensionality.
- **Output of Encoder (Latent Space):** The encoder outputs a compressed version of the data often called the latent representation or encoding. This is a condensed version of the input retaining only the important features.

2. Bottleneck (Latent Space)

The bottleneck is the smallest layer of the network where the data is represented in its most compressed form. It's often referred to as the latent space or code.

- This layer contains a reduced set of features representing the most important information from the input.
- The idea is that through this compression the network learns the key patterns and structures of the input data.

3. Decoder

The decoder is responsible for taking the compressed representation from the latent space and reconstructing it back into the original data form.

- **Hidden Layers:** The decoder uses a series of layers to gradually expand the compressed data back to the original input's dimensions.
- **Output Layer:** This layer produces the reconstructed data and aim to closely resemble the input data.

Loss Function in Autoencoder Training

During training an autoencoder aims to minimize the reconstruction loss which measures the difference between the original input and the reconstructed output. The choice of loss function depends on the type of data:

- **Mean Squared Error (MSE):** This is commonly used for continuous data. It measures the average squared differences between the input and the reconstructed data.
- **Binary Cross-Entropy:** Used for binary data (0 or 1 values). It calculates the difference in probability between the original and reconstructed output.

The network adjusts its weights to minimize this reconstruction loss learning to extract and retain only the most important features of the input data which are then stored in the latent space (bottleneck layer).



DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

Autoencoder Efficient Representations

Constraining an autoencoders makes it learn and show efficient representation. Constraining an autoencoder means that the network learns meaningful, compact and useful features from the input data. After the network is trained only the encoder part is used to encode similar data for future tasks. Several techniques can be employed to achieve this goal:

- **Keep Small Hidden Layers:** By making each hidden layer as small as possible the network is forced to learn only the most important features of the data. A smaller layer size means less redundancy and efficient encoding.
- **Regularization:** Regularization methods like L1 or L2 regularization add a penalty terms to the loss function, preventing the network from overfitting. This ensures that the learned representations are more general and useful by discouraging excessively large weights.
- **Denoising:** In denoising autoencoders random noise is added to the input data during training. The network is then forced to remove this noise during reconstruction which helps it focus on the core, noise-free features of the data and makes the model more robust.
- **Tuning the Activation Functions:** Modifying the activation functions can encourage sparsity in the hidden layer, where only a few neurons are active at a time. This sparsity forces the network to learn only the most relevant features, reducing the complexity of the model and improving its efficiency.

Types of Autoencoders

Autoencoders come in several types each suited for different tasks and with unique features. Let's understand the main types of autoencoders.

1. Denoising Autoencoder

Denoising Autoencoder is trained to work with corrupted or noisy input and learns to remove the noise and reconstruct the original, clean data. It helps the network avoid memorizing the input and forces it to learn the core features of the data instead.

2. Sparse Autoencoder

Sparse Autoencoder contains more hidden units than the input but only a few are allowed to be active at once. This property is called the sparsity of the network. The sparsity of the network can be controlled by either manually zeroing the required hidden units, tuning the activation functions or by adding a loss term to the cost function.

3. Variational Autoencoder

Variational autoencoder (VAE) makes assumptions about the probability distribution of the data and tries to learn a better approximation of it. It uses stochastic gradient descent to optimize and learn the distribution of latent variables. VAEs are often used for generating new data such as creating realistic images or text.

It assumes that the data is generated by a Directed Graphical Model and tries to learn an approximation to $q\phi(z|x)$ to the conditional property $q\theta(z|x)$ where ϕ and θ are the parameters of the encoder and the decoder respectively.



4. Convolutional Autoencoder

Convolutional autoencoder uses convolutional neural networks (CNNs) which are specifically designed for processing images. In this type of autoencoder the encoder uses convolutional layers to extract features from an image and the decoder applies deconvolution also called upsampling to reconstruct the image.

Implementation of Autoencoders

We've created an autoencoder comprising two Dense layers: an encoder responsible for condensing the images into a 64-dimensional latent vector and a decoder tasked with reconstructing the initial image based on this latent space.