Autoencoder Overview

An **autoencoder** is a type of neural network architecture designed for unsupervised learning, primarily used for dimensionality reduction, feature learning, and data compression. It learns to encode input data into a lower-dimensional representation (called the **latent space** or **code**) and then decode it back to the original input. The network is trained to minimize the reconstruction error — the difference between the original input and the output generated by the decoder.

Autoencoders are composed of two main parts:

- 1. **Encoder**: This part of the network takes the input data and compresses it into a lowerdimensional latent space representation. The encoder is typically made of one or more layers of neurons that gradually reduce the dimensionality of the input data.
- 2. **Decoder**: The decoder takes the compressed representation from the encoder and attempts to reconstruct the original input from this representation. The decoder typically mirrors the structure of the encoder but in reverse, gradually expanding the representation back to the original data's dimensionality.

Basic Structure of an Autoencoder:

- Input layer: The original data (e.g., an image or a vector).
- **Encoder**: Compresses the input into a latent space representation (usually of lower dimensionality).
- Latent space: A compressed encoding of the input data.
- **Decoder**: Reconstructs the input data from the latent space representation.
- **Output layer**: The reconstructed data, ideally as close as possible to the original input.

Mathematical Formulation

Let the input data be $x \in Rnx \in Rnx \in Rn \in Rn \in Rn$ (e.g., an nnn-dimensional vector representing an image or some data point). The autoencoder is trained to minimize the reconstruction error:

Where:

- xxx is the original input.
- $x^{h} \{x\} x^{h}$ is the reconstructed output from the decoder.
- ||·||| \cdot ||||·|| typically represents the L2 norm (Euclidean distance) between the input and the reconstructed output.

The encoder function f(x)f(x)f(x) maps the input xxx to the latent representation, and the decoder function g(z)g(z)g(z) reconstructs the data from the latent representation zzz.

z=f(x) and $x^{2}=g(z)z = f(x) \quad \text{(and)} \quad \text{(and)}$

Training the Autoencoder

Training an autoencoder is done by minimizing the reconstruction error across a dataset. This is typically done using gradient-based optimization methods (like stochastic gradient descent, or variants such as Adam) and backpropagation. The goal is to adjust the weights of the encoder and decoder networks so that the output is as close as possible to the input data.

Applications of Autoencoders:

- 1. **Dimensionality Reduction**: Autoencoders can be used for reducing the dimensionality of data in an unsupervised manner (similar to PCA), by taking the compressed representation from the encoder.
- 2. **Data Compression**: The encoder can learn to represent data in a more compact form, which can be used for data compression.
- 3. Anomaly Detection: Since autoencoders learn to reconstruct the input data, they can identify anomalies by measuring the reconstruction error. If the model cannot reconstruct certain data points well (due to being different from the majority of data), those points are flagged as anomalies.
- 4. **Denoising**: A **denoising autoencoder** can be trained to learn a cleaner representation of noisy data. In this case, the autoencoder is trained to reconstruct the original (clean) input from a noisy version of the data.

Regularized Autoencoder

A **regularized autoencoder** refers to an autoencoder that has additional regularization terms added to the loss function to encourage certain desired properties in the learned representations. Regularization techniques help prevent the model from overfitting the training data and improve the generalization of the learned features.

Why Regularize Autoencoders?

Autoencoders are prone to overfitting, especially when the network is large and has too many parameters relative to the available training data. Regularization helps control the capacity of the network, leading to better generalization to unseen data. It can also encourage the autoencoder to learn meaningful, sparse, or structured representations that are useful for downstream tasks.

Types of Regularized Autoencoders

- 1. Sparse Autoencoder:
 - The sparse autoencoder enforces **sparsity** in the hidden layer activations. This means that, for any given input, only a small fraction of the hidden units should be active. To enforce this, a penalty term is added to the loss function based on the average activation of the hidden units.

Loss function with sparsity regularization:

 $\label{eq:listic_list$

Where:

- ο pj\rho_jpj is the average activation of hidden unit jjj.
- \circ λ lambda λ is a regularization parameter that controls the sparsity level.

This regularization encourages the network to learn a compact representation of the data with only a few active neurons for each input.

2. Denoising Autoencoder (DAE):

- A denoising autoencoder is a variant where the input is deliberately corrupted (for example, by adding noise or removing some pixels in an image) and the autoencoder is trained to reconstruct the original, clean input.
- The denoising process forces the autoencoder to learn more robust and general features of the data, making it less likely to memorize noise and more likely to extract meaningful patterns.

Loss function for denoising:

$$\label{eq:loss_loss} \begin{split} Ldenoising=||x-x^{\wedge}||2 \\ mathcal_{L}_{\det}(denoising) &= ||x - hat_{x}||^{2}\\ Ldenoising=||x-x^{\wedge}||^{2} \end{split}$$

Where xxx is the original input, and $x^{hat}{x}x^{h}$ is the reconstruction from the corrupted input.

3. Contractive Autoencoder:

- The contractive autoencoder adds a penalty to the loss function that encourages the learned representation to be **smooth** and **robust** to small changes in the input data.
- This is achieved by adding a **Jacobian** regularization term, which measures how much the activations of the hidden units change with respect to small changes in the input.

Loss function with contractive regularization:

 $Lcontractive=||x-x^{||}2+\beta\sum_{i,j}(\partial hi\partial x_{j})2 \\ \mbox{ + beta } sum_{i,j} \left(frac_{partial h_i} (partial x_{j}) right)^2 \\ \mbox{ Lcontractive } = ||x-x^{||}2+\beta_{i,j}\sum_{j}(\partial x_{j}\partial h_{j})2 \\ \mbox{ + beta } sum_{i,j} \left(frac_{i,j} \left(frac_{i,j}) \\ \mbox{ Lcontractive } = ||x-x^{||}2+\beta_{i,j}\sum_{j}(\partial x_{j}\partial h_{j})2 \\ \mbox{ + beta } sum_{i,j} \left(frac_{i,j} \left(frac_{i,j}) \\ \mbox{ Lcontractive } = ||x-x^{||}2+\beta_{i,j}\sum_{j}(\partial x_{j}\partial h_{j})2 \\ \mbox{ Homology } sum_{i,j} \left(frac_{i,j}) \\ \mbox{ Lcontractive } sum_{i,j} \left(frac_{i,j}) \left(frac_{i,j}) \\ \mbox{ Lcontractive } sum_{i,j} \left(frac_{i,j}) \left(frac_{i,j}) \\ \mbox{ Lcontractive } sum_{i,j} \left(frac_{i,j}) \left(frac_{i,j})$

Where $\partial hi\partial xj \frac{\lambda_i}{\beta}$ (partial h_i) { $\rho x_j \partial x_j \partial h_i$ is the partial derivative of the iii-th hidden unit with respect to the jjj-th input feature, and $\beta beta\beta$ is the regularization coefficient.

This regularization helps prevent the model from overfitting to small fluctuations in the data by ensuring that small changes in the input do not lead to large changes in the hidden representation.

4. Variational Autoencoder (VAE):

- A Variational Autoencoder (VAE) introduces a probabilistic approach to autoencoding. In a VAE, the encoder outputs parameters of a distribution (typically a Gaussian distribution), and the latent representation is sampled from this distribution.
- The decoder then reconstructs the input from this latent representation, but unlike regular autoencoders, a VAE also introduces a **KL-divergence** term in the loss function to ensure that the learned latent space distribution is close to a standard normal distribution.

Loss function for VAE:

$$\begin{split} LVAE = &||x-x^{\wedge}||^{2} + \beta \cdot DKL(q(z|x)||p(z)) \\ &| attacher ||x-x^{\wedge}||^{2} + \beta \cdot DKL(q(z|x)||p(z)) \\ &D_{KL}(q(z|x) || p(z)) \\ LVAE = &||x-x^{\wedge}||^{2} + \beta \cdot DKL(q(z|x)||p(z)) \end{split}$$

Where:

- \circ q(z|x)q(z|x)q(z|x) is the learned distribution over the latent space.
- \circ p(z)p(z)p(z) is the prior distribution (usually standard normal).
- DKLD_{KL}DKL is the **Kullback-Leibler divergence** between the learned and prior distributions.

The VAE regularizes the latent space by forcing it to follow a known distribution, which makes the representation more useful for generating new data (sampling from the latent space).

Conclusion

- Autoencoders are neural networks designed for unsupervised learning of efficient representations of input data by encoding it into a lower-dimensional space and then decoding it back. They are useful for dimensionality reduction, denoising, and anomaly detection.
- **Regularized autoencoders** introduce additional constraints to the autoencoder's architecture to improve generalization, enforce sparsity, smoothness, or create more structured latent spaces. Regularization techniques such as sparse autoencoders, denoising autoencoders, and contractive autoencoders help ensure that the learned representations capture meaningful features while avoiding overfitting.