

# Table of Contents

Chapter 0 Introduction

Chapter 1 Basics and Mechanism

Chapter 2 Applications and Types



# Chapter 0 Introduction

---

## Autoencoders: A Comprehensive Definition

### What is an Autoencoder?

An **autoencoder** is a specialized type of neural network architecture designed to learn efficient data representations in an **unsupervised manner**. It's a self-supervised learning model that learns to compress data into a compact form and then reconstruct it back to closely match the original input. The key principle is learning to encode the essential features of input data while filtering out noise and irrelevant information.

### Formal Definition

Formally, an autoencoder consists of two key functions:

1. **Encoder Function:**  $g : \mathbb{R}^d \rightarrow \mathbb{R}^k$
2. Maps input data from the original space (dimension  $d$ ) to a representation space (dimension  $k$ )
3. Creates a compressed representation  $a \in \mathbb{R}^k$  where typically  $k < d$
4. Learns to extract the most important features
5. **Decoder Function:**  $h : \mathbb{R}^k \rightarrow \mathbb{R}^d$
6. Maps the compressed representation back to the original data space
7. Attempts to reconstruct the input as accurately as possible
8. Acts as the inverse operation of the encoder

# Key Characteristics

## Architecture

- **Unsupervised Learning:** No labeled data is required for training
- **Encoder-Decoder Structure:** Symmetrical architecture with a bottleneck
- **Latent Space Representation:** The compressed representation in the middle layer
- **Reconstruction Loss:** Measures how well the reconstructed output matches the input

## Core Mechanism

The autoencoder is trained to minimize the difference between the input and its reconstruction:

$$L = ||\mathbf{x} - \text{Decoder}(\text{Encoder}(\mathbf{x}))||^2$$

Where:

- $\mathbf{x}$  is the original input
- The loss encourages the network to learn the most compact and informative representation

## What Autoencoders Learn

Autoencoders learn to identify and extract:

- **Key Patterns:** Underlying structure in the data
- **Essential Features:** The most discriminative attributes
- **Data Distribution:** How data points are organized in the feature space
- **Dimensionality Reduction:** A lower-dimensional representation without losing critical information

## Why Use Autoencoders?

### 1. Dimensionality Reduction

- Compress high-dimensional data into lower-dimensional representations
- Useful when dealing with large datasets or images

### 2. Feature Extraction

- Automatically learn meaningful features without manual engineering
- The encoder can be repurposed for other tasks

### 3. Noise Handling

- Denoising autoencoders can remove corruption from data
- Useful for data cleaning and preprocessing

### 4. Anomaly Detection

- Data that deviates from the normal distribution will have higher reconstruction error
- Useful for fraud detection, system monitoring, etc.

### 5. Generative Modeling

- Variational autoencoders (VAEs) can generate new data samples
- Learn the distribution of the training data

### 6. Data Augmentation

- Generate synthetic variations of training data
- Helps improve model generalization with limited data

## Historical Context

Autoencoders have evolved from simple linear autoencoders (related to Principal Component Analysis) to:

- **Deep Autoencoders:** Multiple layers enabling learning of hierarchical representations
- **Variational Autoencoders (VAEs):** Probabilistic approach for generative modeling
- **Denoising Autoencoders (DAEs):** Specifically designed for noise removal
- **Convolutional Autoencoders (CAEs):** For image and spatial data
- **Adversarial Autoencoders (AAEs):** Combine autoencoders with adversarial training

## Comparison with Related Techniques

Aspect	Autoencoder	PCA	VAE
Learning Type	Unsupervised	Unsupervised	Unsupervised
Architecture	Neural Network	Linear Transformation	Neural Network + Probabilistic
Non-linearity	Yes (can be)	No	Yes

Aspect	Autoencoder	PCA	VAE
Generation	Limited	No	Yes
Complexity	High	Low	Very High

## When to Use Autoencoders

### ✓ Good for:

- Compressing images or structured data
- Learning data representations
- Detecting anomalies
- Handling unlabeled data
- Complex non-linear relationships

### ✗ Not ideal for:

- Very small datasets (limited training data)
- Simple linear relationships (use PCA)
- When interpretability is critical
- Real-time applications requiring fast inference

## Summary

Autoencoders are powerful, versatile neural networks that learn to compress and reconstruct data, enabling unsupervised learning of meaningful data representations. They form the foundation for many advanced techniques in deep learning and have applications ranging from data compression to generative modeling.

---

**Next:** Chapter 1 will explore the detailed architecture and mechanism of autoencoders.

---

# Chapter 1 Basics and Mechanism

---

## Chapter 1: Basics and Mechanism of Autoencoders

### 1.1 Architecture Overview

An autoencoder consists of three main components:

#### 1. Encoder Network

- Takes high-dimensional input data
- Progressively reduces dimensionality through layers
- Output is a compressed representation (latent vector)
- Maps: Input Space (high-dim)  $\rightarrow$  Latent Space (low-dim)

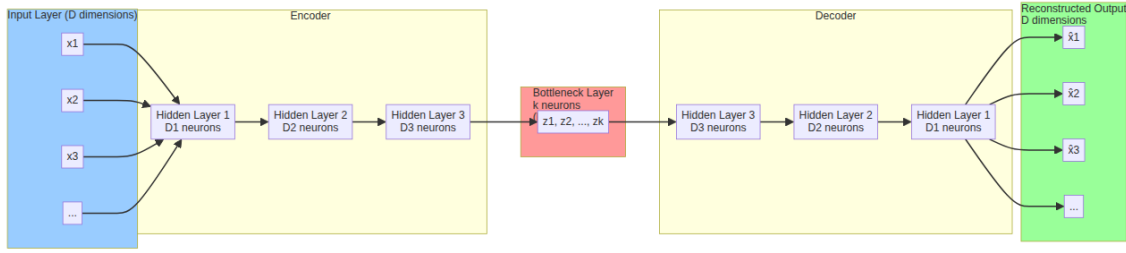
#### 2. Bottleneck (Latent Layer)

- The compressed representation of the input
- Has the smallest number of neurons
- Serves as the "information bottleneck"
- Forces the network to learn only the essential features

#### 3. Decoder Network

- Takes the compressed representation
- Progressively increases dimensionality back to original size
- Attempts to reconstruct the original input
- Maps: Latent Space (low-dim)  $\rightarrow$  Output Space (high-dim)

## 1.2 Architecture Diagram



## 1.3 Data Flow and Mathematical Formulation

### Forward Pass

#### Encoding Phase:

$$\mathbf{a} = g(\mathbf{W}_e \mathbf{x} + \mathbf{b}_e)$$

Where:

- $\mathbf{x}$  = input vector
- $\mathbf{W}_e$  = encoder weight matrix
- $\mathbf{b}_e$  = encoder bias
- $g$  = activation function (ReLU, sigmoid, tanh, etc.)
- $\mathbf{a}$  = encoded representation (latent vector)

#### Decoding Phase:

$$\hat{\mathbf{x}} = h(\mathbf{a}) = \sigma(\mathbf{W}_d \mathbf{a} + \mathbf{b}_d)$$

Where:

- $\mathbf{W}_d$  = decoder weight matrix
- $\mathbf{b}_d$  = decoder bias
- $\hat{\mathbf{x}}$  = reconstructed output

### Complete Forward Pass

$$\hat{\mathbf{x}} = h(g(\mathbf{x})) = \text{Decoder}(\text{Encoder}(\mathbf{x}))$$

## 1.4 Loss Function

The autoencoder is trained by minimizing the **reconstruction loss**:

$$L(\mathbf{x}, \hat{\mathbf{x}}) = \|\mathbf{x} - \hat{\mathbf{x}}\|^2$$



This is typically the **Mean Squared Error (MSE)**:

$$L = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mathbf{x}}_i)^2$$

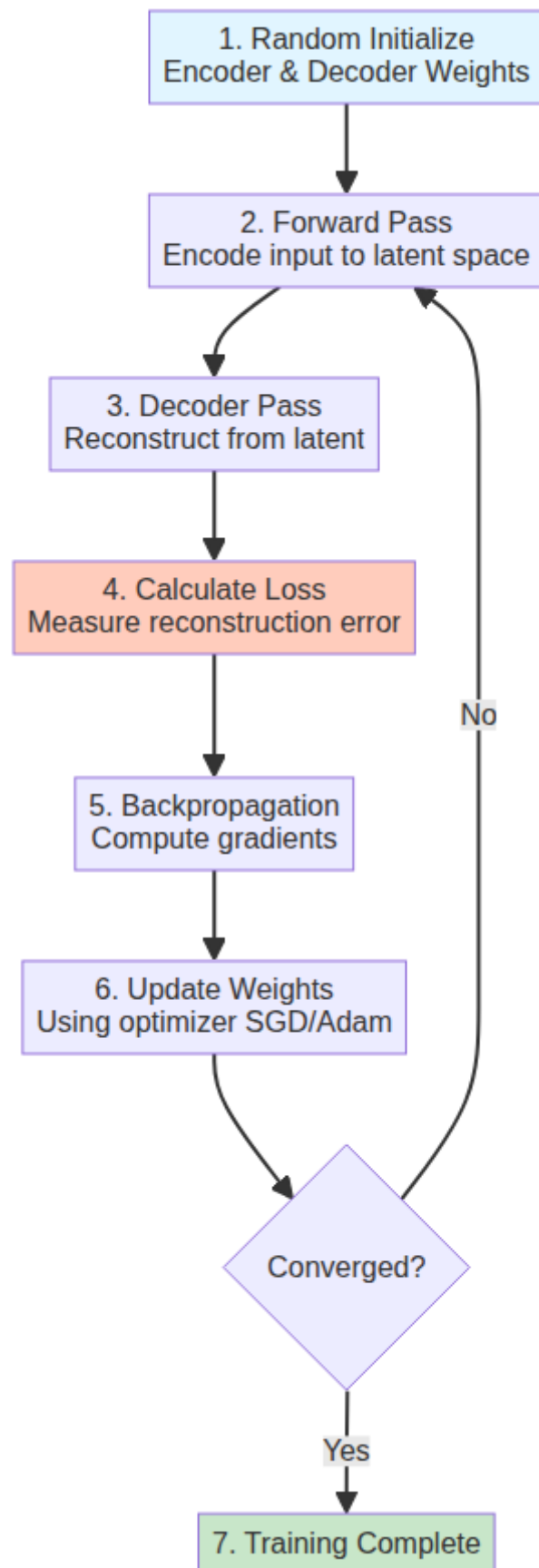
For binary data, **Binary Cross-Entropy** might be used:

$$L = -\frac{1}{n} \sum_{i=1}^n [\mathbf{x}_i \log(\hat{\mathbf{x}}_i) + (1 - \mathbf{x}_i) \log(1 - \hat{\mathbf{x}}_i)]$$

## 1.5 Training Process

The autoencoder learns through standard backpropagation:

## Step-by-Step Training



## Pseudo-code for Training Loop

```
for epoch in range(num_epochs):
    total_loss = 0
    for batch in training_data:
        # 1. Forward pass
        latent = encoder(batch)
        reconstruction = decoder(latent)

        # 2. Calculate loss
        loss = mse_loss(batch, reconstruction)

        # 3. Backward pass
        loss.backward()

        # 4. Update weights
        optimizer.step()
        optimizer.zero_grad()

    total_loss += loss.item()

print(f"Epoch {epoch}, Loss: {total_loss/len(training_data)}")
```

## 1.6 Activation Functions in Autoencoders

Different layers use different activation functions:

### Encoder Layers

Activation	Formula	Use Case
ReLU	$\max(0, x)$	Default choice, fast computation
Leaky ReLU	$\max(0.01x, x)$	Avoid dead neurons
ELU	$x$ if $x > 0$ else $\alpha(e^x - 1)$	Smooth gradients

Activation	Formula	Use Case
<b>Sigmoid</b>	$\frac{1}{1+e^{-x}}$	Squash to [0,1]
<b>Tanh</b>	$\frac{e^x - e^{-x}}{e^x + e^{-x}}$	Squash to [-1,1]

## Latent Layer

- Usually **no activation** (linear)
- Allows full range of values for representation

## Decoder Output Layer

- **Sigmoid**: For images (output in [0,1])
- **Tanh**: For centered data (output in [-1,1])
- **Linear**: For unbounded reconstruction
- **Softmax**: For probability distributions

# 1.7 Latent Space Representation

The latent space is the key to understanding autoencoders:

## Properties of Latent Space

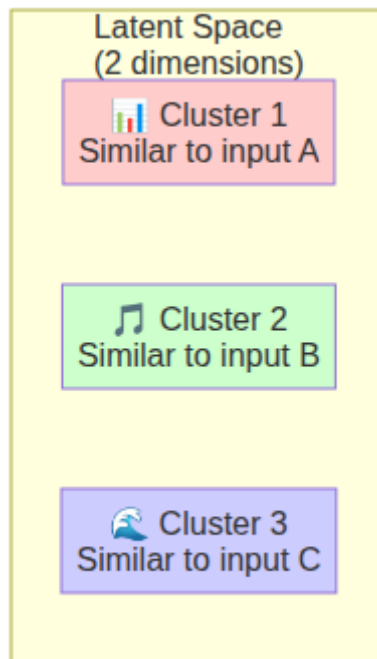
### Compressed Dimension:

- Original:  $28 \times 28 = 784$  dimensions (MNIST image)
- Latent: 32 dimensions (typical choice)
- **Compression ratio**:  $784/32 \approx 24.5x$

### Learned Representation:

- Captures underlying factors of variation
- Similar inputs map to nearby latent points
- Smooth interpolation is possible

### Visualization Example (2D Latent Space):



## 1.8 Encoder-Decoder Symmetry

Most autoencoders have **symmetric architecture**:

### Typical Configuration

**Encoder:**  $784 \rightarrow 512 \rightarrow 256 \rightarrow 128 \rightarrow 32$

- Reduces by factor of 2 at each layer

**Decoder:**  $32 \rightarrow 128 \rightarrow 256 \rightarrow 512 \rightarrow 784$

- Increases by factor of 2 at each layer (mirror of encoder)

This symmetry is useful but not required. **Asymmetric architectures** are also valid.

## 1.9 Key Hyperparameters

Parameter	Impact	Typical Values
Latent Dimension	Smaller = more compression, larger = better reconstruction	32, 64, 128

Parameter	Impact	Typical Values
Encoder Depth	More layers = more capacity but harder to train	2-4 layers
Hidden Layer Size	Capacity of intermediate representations	256, 512, 1024
Activation Function	Non-linearity, training dynamics	ReLU, Tanh
Learning Rate	Convergence speed and stability	0.001, 0.0001
Batch Size	Training stability and speed	32, 64, 128
Regularization	Prevent overfitting	L1, L2, Dropout

## 1.10 Common Issues and Solutions

### Issue 1: Blurry Reconstructions

**Cause:** Latent dimension too small, network underfitting

**Solutions:**

- Increase latent dimension
- Add more layers
- Increase training epochs

### Issue 2: Overfitting

**Cause:** Network learning noise instead of general features

**Solutions:**

- Add dropout regularization
- Use L1/L2 regularization
- Increase latent dimension (forces compression)
- Use denoising autoencoder variant

## Issue 3: Training Divergence

**Cause:** Learning rate too high, gradient explosion

**Solutions:**

- Reduce learning rate
- Use batch normalization
- Implement gradient clipping
- Use Adam optimizer (adaptive learning rates)

## Issue 4: Dead Neurons

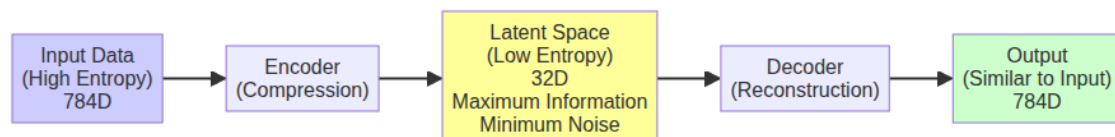
**Cause:** ReLU units can get stuck at 0

**Solutions:**

- Use Leaky ReLU or ELU
- Adjust learning rate
- Use batch normalization

## 1.11 Information Bottleneck Principle

The core principle of autoencoders:



**The Tradeoff:**

- **Smaller latent space:** Stronger compression, more information loss
- **Larger latent space:** Better reconstruction, less compression

The optimal size balances:

$\text{Compression Ratio}$  vs  $\text{Reconstruction Quality}$



## 1.12 Comparison: Shallow vs Deep Autoencoders

### Shallow Autoencoder

Input (784) → Hidden (256) → Latent (32) → Hidden (256) → Output (784)

- **Pros:** Fast training, simple
- **Cons:** Limited capacity, struggles with complex patterns

### Deep Autoencoder

Input (784) → 512 → 256 → 128 → Latent (32) → 128 → 256 → 512 → Output (784)

- **Pros:** Learns hierarchical representations, better compression
- **Cons:** Harder to train, requires careful tuning

## Summary

Autoencoders work through a simple yet powerful mechanism:

1. **Compress** data into a bottleneck layer
2. **Reconstruct** the original data from the compressed representation
3. **Learn** by minimizing reconstruction error
4. **Extract** meaningful features through the information bottleneck

This forces the network to discover the essential structure in the data without explicit supervision.

---

**Next:** Chapter 2 explores specific applications and different types of autoencoders suited for various tasks.

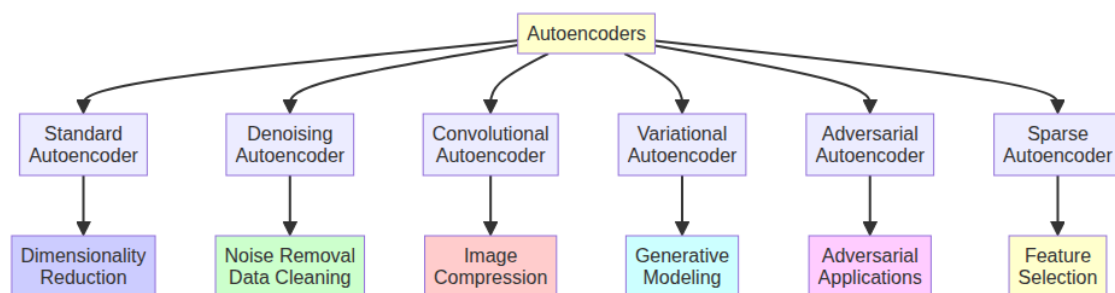
# Chapter 2 Applications and Types

---

## Chapter 2: Applications and Types of Autoencoders

### 2.1 Overview of Autoencoder Variants

Different types of autoencoders are designed for specific applications:



### 2.2 Standard (Vanilla) Autoencoder

#### Architecture

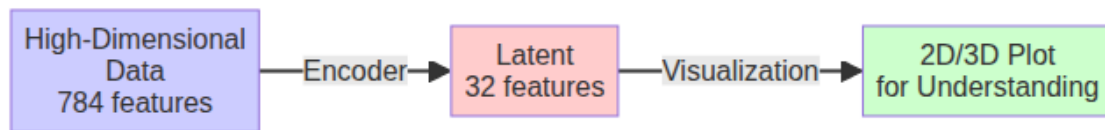
- Fully connected layers
- Symmetric encoder-decoder structure
- Minimizes MSE reconstruction loss

#### Applications

##### 1. Dimensionality Reduction

- **Problem:** High-dimensional data is hard to visualize and process
- **Solution:** Compress to 2-3 dimensions
- **Use Case:**

- MNIST digits: 784D → 32D
- Images: 65536D → 128D



## 2. Feature Learning

- **Problem:** Manual feature engineering is time-consuming
- **Solution:** Learn representations automatically
- **Use Case:** Transfer learning with pre-trained encoder

### Example Workflow:

1. Train autoencoder on unlabeled data
2. Extract encoder weights
3. Use encoder as feature extractor for supervised task
4. Fine-tune on labeled data

## 3. Data Compression

- **Problem:** Storage and transmission of large files
- **Solution:** Compress via latent representation
- **Use Case:**
  - Image compression (JPEG-like)
  - Video frame compression
  - Document compression

## Hyperparameter Tuning

Parameter	Impact	Recommendation
Latent Size	Compression ratio	Start with 5-10% of input size
Encoder Layers	Capacity	2-4 layers typically

Parameter	Impact	Recommendation
Learning Rate	Convergence	0.001 for Adam
Batch Size	Training stability	32-128

## 2.3 Denoising Autoencoder (DAE)

### Architecture

- Input: Corrupted data
- Output: Clean data
- Training: Add noise → Reconstruct clean version

### How It Works



### Mathematical Formulation

#### Training Process:

1. Start with clean input  $\mathbf{x}$
2. Add noise:  $\tilde{\mathbf{x}} = \mathbf{x} + \epsilon$
3. Encode corrupted input:  $\mathbf{a} = g(\tilde{\mathbf{x}})$
4. Decode:  $\hat{\mathbf{x}} = h(\mathbf{a})$
5. Minimize loss against **clean** input:  

$$\mathcal{L} = \|\mathbf{x} - \hat{\mathbf{x}}\|^2$$

#### At Test Time:

- Feed noisy input
- Network reconstructs clean version
- No need for explicit noise model

## Types of Noise

Noise Type	Formula	Application
Gaussian Noise	$\tilde{x} = x + \mathcal{N}(0, \sigma^2)$	General-purpose
Salt-and-Pepper	Random pixels to 0 or 1	Image corruption
Dropout Noise	Randomly zero neurons	Robustness training
Masking Noise	Randomly mask regions	Inpainting

## Applications

### 1. Image Denoising

- **Problem:** Photos corrupted by sensor noise or compression
- **Solution:** Train on noisy-clean pairs
- **Example:** Medical imaging, low-light photos

Medical Image (noisy) → DAE → Clean Medical Image

### 2. Missing Data Imputation

- **Problem:** Datasets with missing values
- **Solution:** Mask missing values as "noise"
- **Example:** User preferences in recommendation systems

User ratings (some missing) → DAE → Imputed ratings

### 3. Robust Feature Learning

- **Problem:** Features should be robust to small perturbations
- **Solution:** Train with noisy inputs
- **Example:** Speech recognition on noisy audio

## 4. Data Augmentation

- **Problem:** Limited training data
- **Solution:** Generate clean versions from noisy variants
- **Example:** Augment datasets for better generalization

## Key Advantage

The denoising autoencoder learns a **regularized representation** that's:

- Robust to input variations
- Captures essential features
- Better for downstream tasks than vanilla autoencoder

## 2.4 Convolutional Autoencoder (CAE)

### Architecture

- Uses convolutional layers in encoder
- Uses transpose/deconvolutional layers in decoder
- Preserves spatial structure of images

### Why Convolution?

#### Problems with Fully Connected:

- Treats images as 1D vectors
- Doesn't preserve spatial relationships
- Too many parameters for large images

#### Advantages of Convolution:

- Learns local features (edges, textures)
- Parameter sharing reduces model size
- Preserves 2D structure naturally

### Architecture Comparison



## Encoder Architecture Example

```
Input (28×28×1)
    ↓ Conv(3×3, 16 filters) + ReLU
Output (28×28×16)
    ↓ MaxPool(2×2)
Output (14×14×16)
    ↓ Conv(3×3, 32 filters) + ReLU
Output (14×14×32)
    ↓ MaxPool(2×2)
Latent (7×7×32) = 1,568 features
```

## Decoder Architecture Example

```
Latent (7×7×32)
    ↓ Reshape to proper format
Output (7×7×32)
    ↓ UpSample(2×2)
Output (14×14×32)
    ↓ ConvTranspose(3×3, 16 filters) + ReLU
Output (14×14×16)
    ↓ UpSample(2×2)
Output (28×28×16)
    ↓ ConvTranspose(3×3, 1 filter) + Sigmoid
Reconstructed (28×28×1)
```

## Applications

### 1. Image Compression

- **Problem:** Large file sizes for storage/transmission
- **Solution:** Compress via latent space
- **Example:** JPEG alternative with learned compression
- **Compression Ratio:** 10:1 to 100:1 possible

Original Image ( $28 \times 28 \times 3 = 2,352$  bytes)

↓ CAE Encoder

Latent ( $7 \times 7 \times 8 = 392$  features)

↓ CAE Decoder

Reconstructed Image ( $28 \times 28 \times 3$ )

## 2. Image Reconstruction

- **Problem:** Incomplete or damaged images
- **Solution:** Train on full images, inpaint at test time
- **Example:** Old photo restoration, satellite imagery

## 3. Feature Extraction

- **Problem:** Need discriminative features for classification
- **Solution:** Use learned representations from encoder
- **Example:** Pre-training for image classification

## 4. Anomaly Detection

- **Problem:** Identify unusual/defective items
- **Solution:** High reconstruction error = anomaly
- **Example:** Manufacturing quality control

Normal Product → CAE → Low Reconstruction Error

Defective Item → CAE → High Reconstruction Error

## Key Advantages

Advantage	Benefit
Fewer Parameters	Faster training, less memory
Spatial Awareness	Better for structured data like images
Better Features	Learns hierarchical visual features



Advantage	Benefit
Scalability	Can handle high-resolution images

## 2.5 Variational Autoencoder (VAE)

### Key Difference from Standard AE

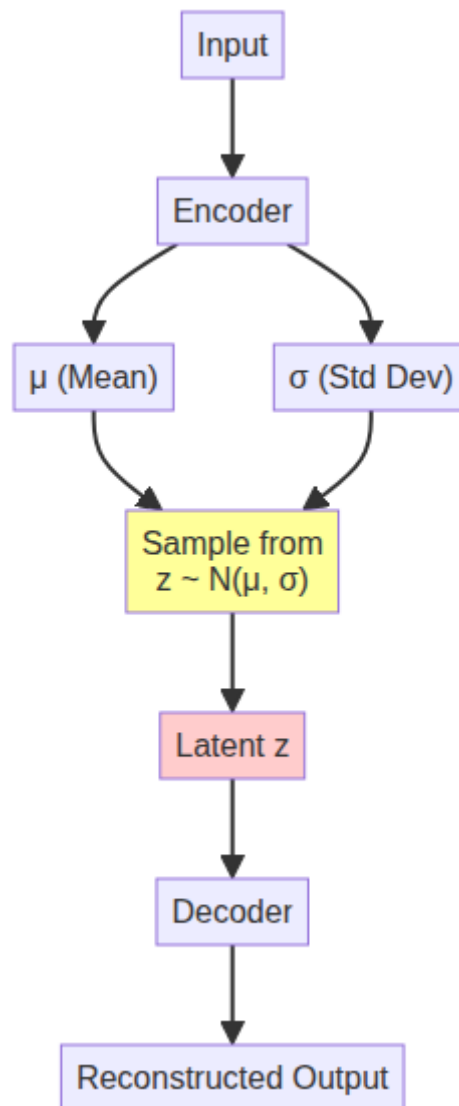
**Standard AE:** Learns a fixed mapping

$$\mathbf{x} \rightarrow z \rightarrow \hat{\mathbf{x}}$$

**VAE:** Learns a **probability distribution** over latent space

$$\mathbf{x} \rightarrow P(z|\mathbf{x}) \rightarrow \hat{\mathbf{x}}$$

## Architecture



## Loss Function

VAEs use a special loss combining:

1. **Reconstruction Loss:** Standard MSE/BCE
2. **KL Divergence:** Ensures latent space is close to prior  $\mathcal{N}(0,1)$

$$\mathcal{L}_{\text{VAE}} = \mathcal{L}_{\text{reconstruction}} + \beta \cdot D_{\text{KL}}(q(z|x) || p(z))$$

Where:

- $q(z|x)$  = learned distribution (encoder)
- $p(z)$  = prior distribution (usually standard normal)
- $\beta$  = weighting factor (controls generation vs reconstruction)

# Applications

## 1. Generative Modeling

- **Problem:** Generate new data similar to training data
- **Solution:** Sample from learned latent distribution
- **Example:** Generate new faces, handwritten digits

Sample  $z \sim N(0,1) \rightarrow$  VAE Decoder  $\rightarrow$  New Image

## 2. Interpolation and Morphing

- **Problem:** Create smooth transitions between data points
- **Solution:** Interpolate in latent space
- **Example:** Face morphing, image blending

Image A  $\rightarrow$  Encode  $\rightarrow z\_A$

Image B  $\rightarrow$  Encode  $\rightarrow z\_B$

$t \cdot z\_A + (1-t) \cdot z\_B \rightarrow$  Decode  $\rightarrow$  Intermediate Image (for  $t \in [0,1]$ )

## 3. Semi-Supervised Learning

- **Problem:** Limited labeled data
- **Solution:** Pre-train VAE on unlabeled data
- **Example:** Classification with mostly unlabeled examples

## 4. Disentangled Representations

- **Problem:** Separate factors of variation
- **Solution:** Use  $\beta$ -VAE or Factor-VAE
- **Example:** Separate rotation, scale, position in images

## When to Use VAE



**Good for:**

- Generating new samples
- Interpretable latent spaces

- Smooth interpolation between examples
- Semi-supervised learning

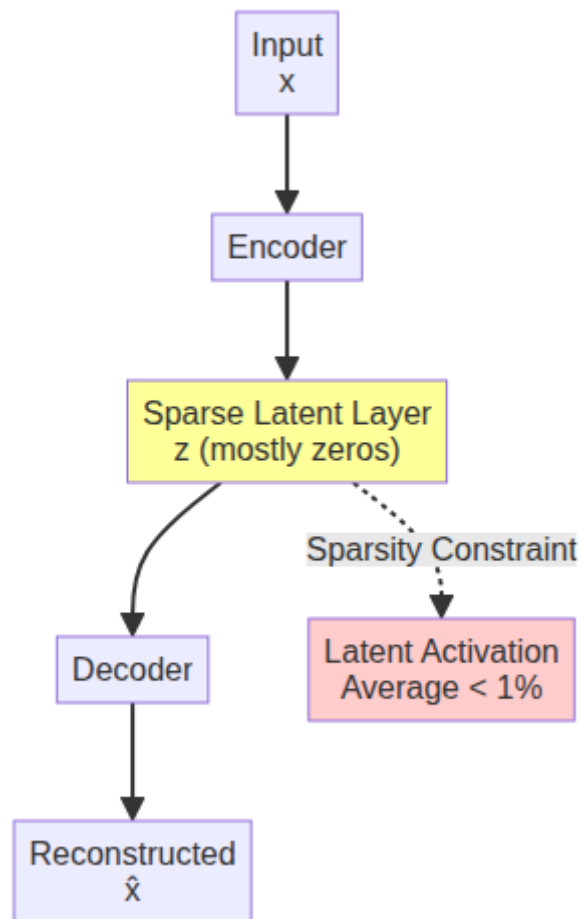
**✗ Not ideal for:**

- Pixel-perfect reconstruction (blurry outputs)
- When generation isn't needed
- Limited training data

## 2.6 Sparse Autoencoder

### Concept

Restrict the number of active neurons in the latent layer



### Loss Function

Adds sparsity penalty:

$$L = L_{\text{reconstruction}} + \lambda \cdot L_{\text{sparsity}}$$

Sparsity can be enforced with:

- **KL Divergence**: Prefer low activation rates
- **L1 Regularization**: Encourage zeros
- **Binary Masks**: Only activate top-k neurons

## Applications

### 1. Feature Selection

- **Problem**: Identify important features from high-dimensional data
- **Solution**: Active neurons = important features
- **Example**: Gene selection in genomics

### 2. Interpretability

- **Problem**: Understand what features matter
- **Solution**: Only few active units per output
- **Example**: Medical diagnosis explanation

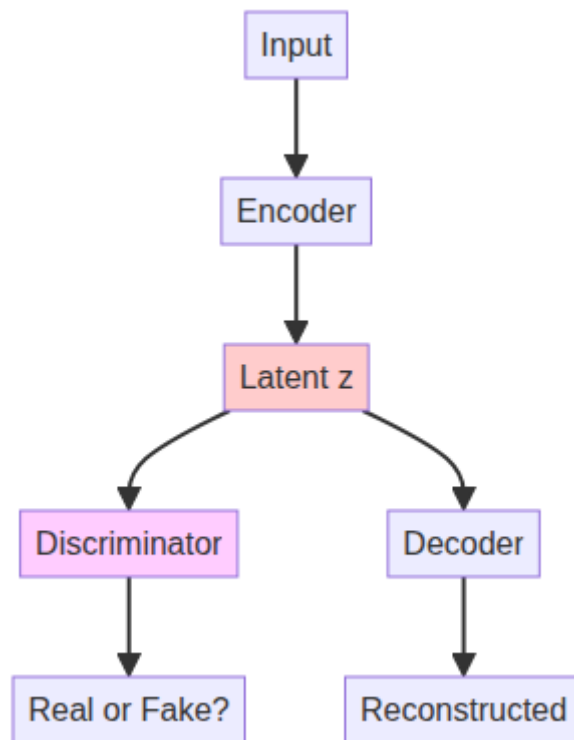
### 3. Efficient Storage

- **Problem**: Sparse latent codes compress better
- **Solution**: Use sparse representations
- **Example**: Semantic hashing, efficient indexing

## 2.7 Adversarial Autoencoder (AAE)

### Hybrid Approach

Combines autoencoders with GANs (Generative Adversarial Networks)



## Two Training Objectives

### 1. Reconstruction Objective

2. Standard autoencoder loss

3. Minimize:  $\|x - \text{Decoder}(\text{Encoder}(x))\|^2$

### 4. Adversarial Objective

5. Discriminator distinguishes encoded vs prior

6. Generator (encoder) fools discriminator

7. Ensures  $q(z) \approx p(z)$

## Applications

### 1. Flexible Generation

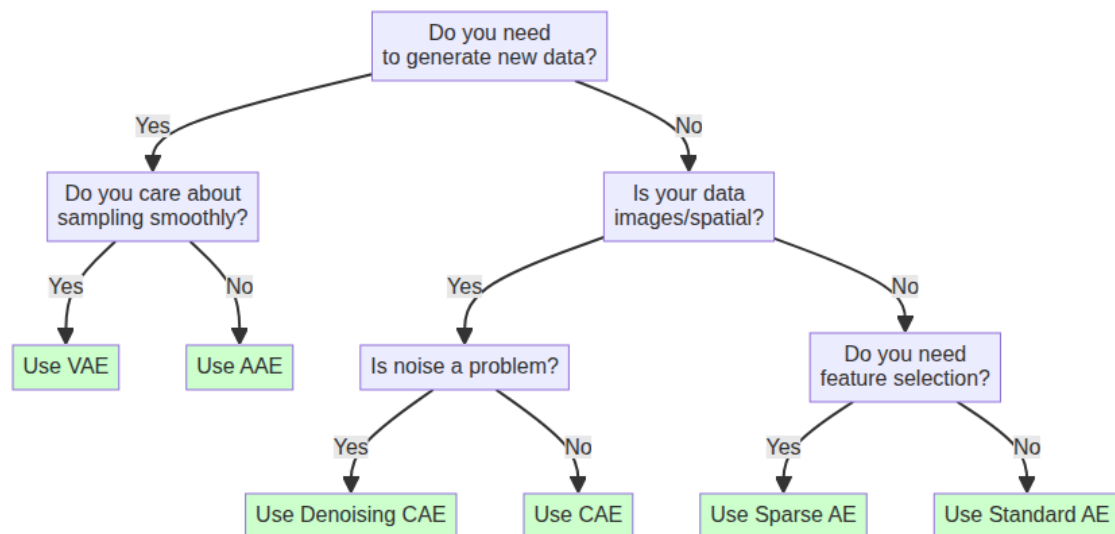
- More stable than VAE
- Higher quality than VAE
- Flexible latent space structure

## 2. Clustering

- Learn cluster structure in latent space
- Categorical variables in latent layer
- Example: Unsupervised classification

## 2.8 Application Selection Guide

Use this flowchart to choose the right autoencoder:



## 2.9 Comparison Table

Type	Input Type	Output	Generation	Noise Robust	Interpretable
Standard	Any	Reconstruction	✗	✗	Medium
Denoising	Noisy data	Clean version	✗	✓	Medium
Convolutional	Images	Image	✗	✗	Medium
Variational	Any	Reconstruction + New samples	✓	✗	High

Type	Input Type	Output	Generation	Noise Robust	Interpretable
Sparse	High-dim	Sparse features	✗	✗	✓
Adversarial	Any	Reconstruction + New samples	✓	✗	Low

## 2.10 Real-World Application Examples

### 1. Medical Imaging

- **Problem:** Noisy CT/MRI scans
- **Solution:** Denoising autoencoder
- **Benefit:** Improved diagnosis accuracy

### 2. Recommendation Systems

- **Problem:** Missing user ratings
- **Solution:** DAE for imputation
- **Benefit:** Better recommendations with sparse data

### 3. Facial Recognition

- **Problem:** Recognize faces from low-resolution images
- **Solution:** CAE pre-training + fine-tuning classifier
- **Benefit:** Better features with limited labeled data

### 4. Fraud Detection

- **Problem:** Detect unusual transaction patterns
- **Solution:** CAE with anomaly detection (high reconstruction error)
- **Benefit:** Identify fraudulent transactions

### 5. Drug Discovery

- **Problem:** Generate novel drug molecules
- **Solution:** VAE on chemical fingerprints



- **Benefit:** Accelerate drug development

## 6. Content Moderation

- **Problem:** Flag inappropriate content
- **Solution:** Sparse AE for feature selection
- **Benefit:** Understand which features trigger flags

## Summary

Each autoencoder type excels in different scenarios:

- **Standard:** General dimensionality reduction
- **Denoising:** Robust features and noise removal
- **Convolutional:** Image processing
- **Variational:** Generative tasks
- **Sparse:** Interpretability and feature selection
- **Adversarial:** High-quality generation with reconstruction

---

**Next:** Chapter 3 provides TensorFlow implementations of three key autoencoder types with complete working code.