

Generative Models: Autoencoders, Adversarial Networks, and Diffusion

October 22, 2025

Contents

1 Autoencoders (AE, VAE)

Autoencoders learn latent representations by reconstructing inputs. The encoder f_ϕ maps \mathbf{x} to latent code \mathbf{z} , and the decoder g_θ reconstructs $\hat{\mathbf{x}}$. Figure ?? illustrates the bottleneck structure.

1.1 Deterministic Autoencoders

The reconstruction loss (often mean squared error) is

$$\mathcal{L}_{\text{AE}}(\theta, \phi) = \frac{1}{N} \sum_{i=1}^N \|g_\theta(f_\phi(\mathbf{x}_i)) - \mathbf{x}_i\|_2^2. \quad (1)$$

While AEs excel at dimensionality reduction, latent spaces may be discontinuous. Regularized variants enforce structure:

- **Sparse AE:** Adds L_1 penalties on activations to encourage sparse latent representations.
- **Denoising AE:** Trains with corrupted inputs $\tilde{\mathbf{x}}$ but reconstructs the clean \mathbf{x} , enhancing robustness.
- **Contractive AE:** Penalizes the Jacobian $\|\nabla_{\mathbf{x}} f_\phi(\mathbf{x})\|_F^2$ to learn invariant features.

1.2 Variational Autoencoders

VAEs impose a probabilistic latent variable model. Given prior $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ and decoder likelihood $p_\theta(\mathbf{x} | \mathbf{z})$, the evidence lower bound (ELBO) for each sample is

$$\log p_\theta(\mathbf{x}) \geq \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})} [\log p_\theta(\mathbf{x} | \mathbf{z})] - \text{KL}(q_\phi(\mathbf{z} | \mathbf{x}) \| p(\mathbf{z})). \quad (2)$$

Parameterizing $q_\phi(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_\phi(\mathbf{x}), \text{diag}(\boldsymbol{\sigma}_\phi^2(\mathbf{x})))$ enables reparameterization:

$$\mathbf{z} = \boldsymbol{\mu}_\phi(\mathbf{x}) + \boldsymbol{\sigma}_\phi(\mathbf{x}) \odot \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (3)$$

The KL term has closed form:

$$\text{KL} = -\frac{1}{2} \sum_{j=1}^d (1 + \log \sigma_j^2 - \mu_j^2 - \sigma_j^2). \quad (4)$$

1.3 Beta-VAE and Disentanglement

β -VAE scales the KL divergence, $\mathcal{L} = \mathbb{E}[\log p_\theta(\mathbf{x} | \mathbf{z})] - \beta \text{KL}(q_\phi \| p)$, promoting disentangled latent factors when $\beta > 1$. Mutual information penalties and Total Correlation regularizers (TC-VAE) refine latent independence.

1.4 Training Procedure

Listing 1: Variational autoencoder training loop with KL annealing.

```
1 for epoch in range(num_epochs):
2     beta = min(1.0, epoch / kl_warmup_epochs)
3     for x in dataloader:
4         mu, logvar = encoder(x)
5         z = mu + torch.exp(0.5 * logvar) * torch.randn_like(mu)
6         x_recon = decoder(z)
7         recon_loss = reconstruction_criterion(x_recon, x)
8         kl = -0.5 * (1 + logvar - mu.pow(2) - logvar.exp()).sum(dim=1).mean()
9         loss = recon_loss + beta * kl
10        loss.backward()
11        optimizer.step()
12        optimizer.zero_grad()
```

Autoencoder Encoder-Decoder Architecture

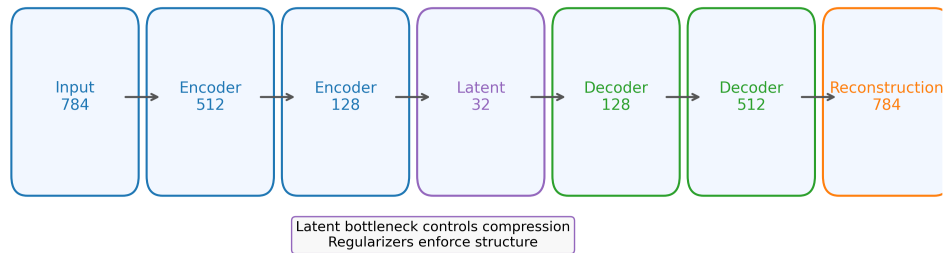


Figure 1: Autoencoder architecture with deterministic encoder/decoder. Bottleneck dimensionality controls compression.

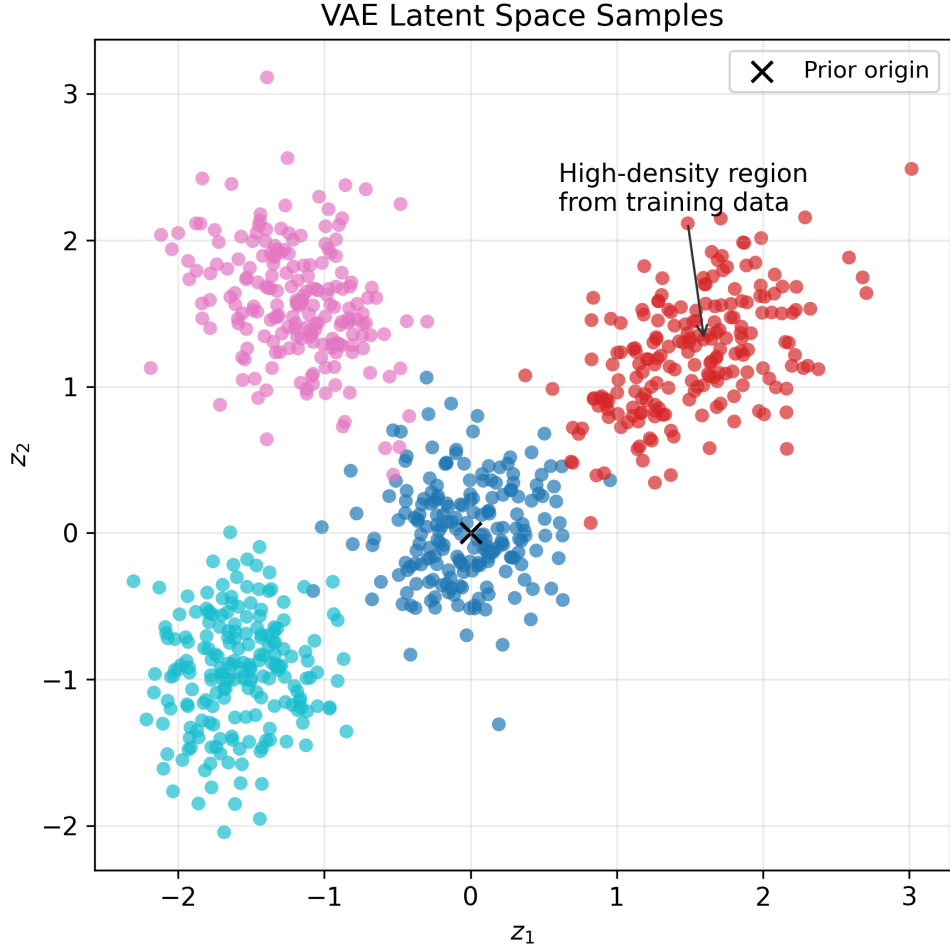


Figure 2: Latent space sampling in a VAE. Points near the origin correspond to plausible reconstructions.

2 Generative Adversarial Networks (GANs)

GANs pit a generator G against a discriminator D . The min-max objective,

$$\min_G \max_D \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] , \quad (5)$$

drives G to synthesize realistic samples. Figure ?? depicts loss curves and mode coverage.

2.1 GAN Variants

- **DCGAN:** Convolutional architectures with strided convolutions, batch normalization, and ReLU/Leaky ReLU activations for stable image synthesis.
- **WGAN and WGAN-GP:** Replace Jensen-Shannon divergence with Earth-Mover distance. The WGAN objective

$$\min_G \max_{D \in \mathcal{D}_1} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [D(\mathbf{x})] - \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [D(G(\mathbf{z}))], \quad (6)$$

constrains D to 1-Lipschitz. Gradient penalty adds $\lambda(\|\nabla_{\hat{\mathbf{x}}} D(\hat{\mathbf{x}})\|_2 - 1)^2$ to encourage Lipschitz continuity.

- **StyleGAN:** Introduces style-based modulation. Latent \mathbf{w} modulates convolution kernels via adaptive instance normalization, enabling control over coarse-to-fine features. Path-length regularization and noise injection improve fidelity.

2.2 Training Stabilization

Mode collapse, vanishing gradients, and discriminator overpowering require safeguards:

- Feature matching and minibatch discrimination regularize G .
- Spectral normalization enforces Lipschitz constraints by rescaling weight matrices.
- Two-time-scale update rule (TTUR) adjusts learning rates (η_D, η_G) to balance convergence.

2.3 Evaluation Metrics

Fréchet Inception Distance (FID) approximates the Wasserstein-2 distance between Inception features:

$$\text{FID} = \|\boldsymbol{\mu}_r - \boldsymbol{\mu}_g\|_2^2 + \text{Tr} \left(\boldsymbol{\Sigma}_r + \boldsymbol{\Sigma}_g - 2(\boldsymbol{\Sigma}_r \boldsymbol{\Sigma}_g)^{1/2} \right). \quad (7)$$

Precision/recall for generative models and Inception Score complement FID.

2.4 StyleGAN2 Generator Forward Pass

Listing 2: Simplified StyleGAN2 generator block with style modulation.

```

1 class StyledConv(nn.Module):
2     def __init__(self, in_channels, out_channels, style_dim, upsample):
3         super().__init__()
4         self.upsample = upsample
5         self.weight = nn.Parameter(torch.randn(1, out_channels, in_channels, 3, 3))
6         self.modulation = nn.Linear(style_dim, in_channels)
7         self.noise_weight = nn.Parameter(torch.zeros(1, out_channels, 1, 1))
8         self.activation = nn.LeakyReLU(0.2)
9
10    def forward(self, x, style, noise):
11        style = self.modulation(style).view(-1, 1, x.size(1), 1, 1)
12        weight = self.weight * (style + 1e-8)
13        if self.upsample:
14            x = upsample_2x(x)
15        x = conv2d_modulated(x, weight)
16        x = x + self.noise_weight * noise
17        return self.activation(x)

```

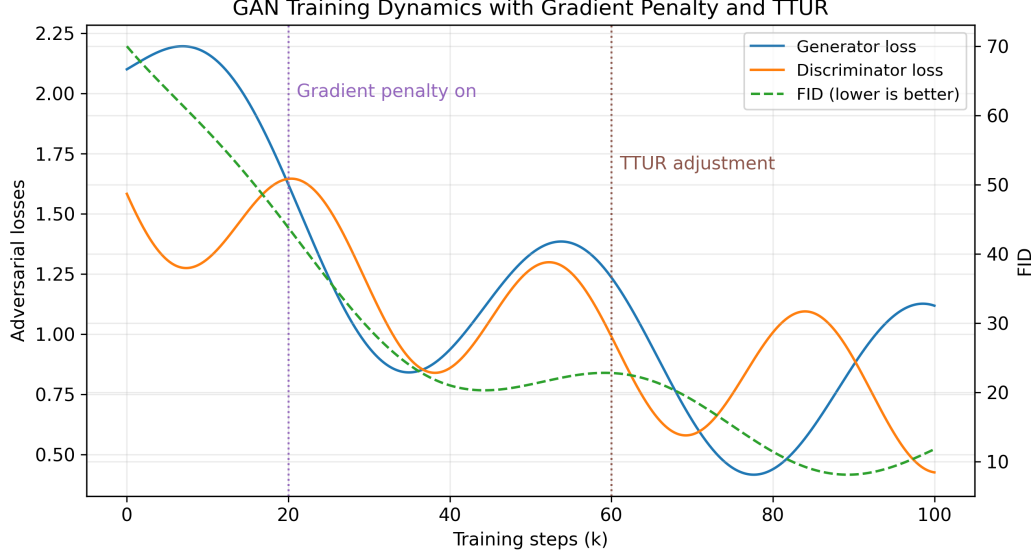


Figure 3: Generator and discriminator loss trajectories with gradient penalty and TTUR scheduling.

3 Diffusion Models Overview

Diffusion models generate data by reversing a gradual noising process. The forward process corrupts data \mathbf{x}_0 into \mathbf{x}_T using a variance schedule $\{\beta_t\}_{t=1}^T$:

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}). \quad (8)$$

Due to Gaussian composition,

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}), \quad \bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s). \quad (9)$$

3.1 Denoising Diffusion Probabilistic Models (DDPM)

The model learns $\epsilon_\theta(\mathbf{x}_t, t)$ to predict noise. The simplified training objective is

$$\mathcal{L}_{\text{simple}} = \mathbb{E}_{t, \mathbf{x}_0, \epsilon} [\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|_2^2]. \quad (10)$$

Sampling begins from $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and iteratively denoises:

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{1 - \beta_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (11)$$

Figure ?? visualizes forward and reverse trajectories.

3.2 Improvements and Variants

- **Guided diffusion:** Classifier guidance adjusts sampling drift, $\hat{\epsilon} = \epsilon_\theta - \sigma_t \nabla_{\mathbf{x}_t} \log p_\phi(y | \mathbf{x}_t)$, improving class-conditional fidelity.
- **Score-based generative models:** Learn $\nabla_{\mathbf{x}} \log q_t(\mathbf{x})$ with stochastic differential equations (SDEs) and integrate using predictor-corrector samplers.
- **Latent diffusion:** Compress data into latent space via VAE before diffusion (e.g., Stable Diffusion), reducing computational cost while leveraging cross-attention for conditioning.

3.3 Pseudo-code

Listing 3: Diffusion training step with cosine noise schedule.

```
1 def diffusion_training_step(model, scheduler, x0):
2     t = torch.randint(0, scheduler.num_steps, (x0.size(0),), device=x0.device)
3     noise = torch.randn_like(x0)
4     alpha_bar = scheduler.alpha_bar(t).view(-1, 1, 1, 1)
5     xt = torch.sqrt(alpha_bar) * x0 + torch.sqrt(1 - alpha_bar) * noise
6     noise_pred = model(xt, t)
7     loss = (noise - noise_pred).pow(2).mean()
8     loss.backward()
9     optimizer.step()
10    optimizer.zero_grad()
11    return loss
```

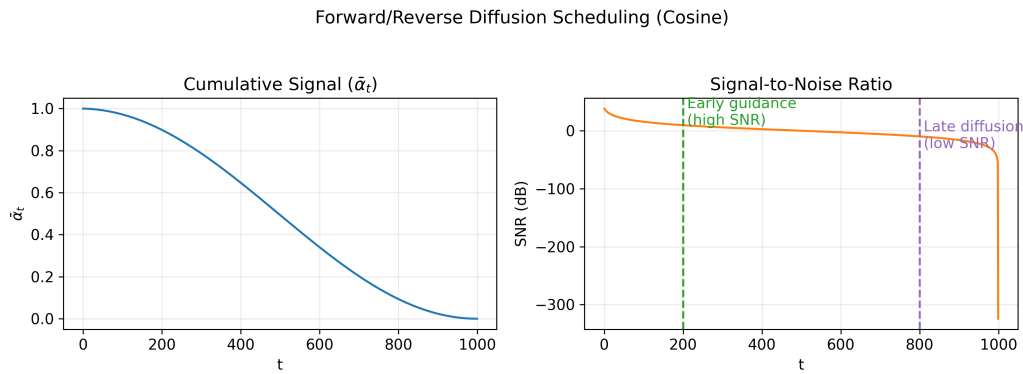


Figure 4: Forward noising and reverse denoising trajectories in a diffusion model with cosine schedule.

Further Reading

- Diederik P. Kingma and Max Welling. “Auto-Encoding Variational Bayes.” ICLR 2014.
- Ian Goodfellow et al. “Generative Adversarial Networks.” NIPS 2014.
- Tero Karras et al. “Analyzing and Improving the Image Quality of StyleGAN.” CVPR 2020.
- Jonathan Ho et al. “Denoising Diffusion Probabilistic Models.” NeurIPS 2020.
- Prafulla Dhariwal and Alexander Nichol. “Diffusion Models Beat GANs on Image Synthesis.” NeurIPS 2021.