

DS 4440

Auto-encoders & Diffusion Models

Many problems in ML amount to generative modeling.

Consider image generation

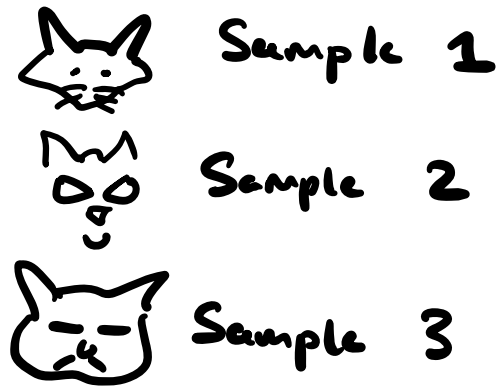
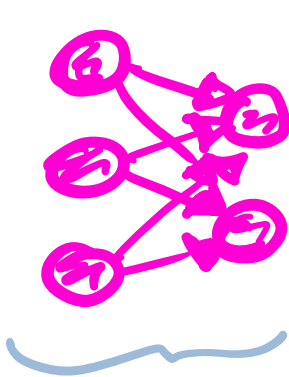


Image generation Model

Hopefully a bit better though

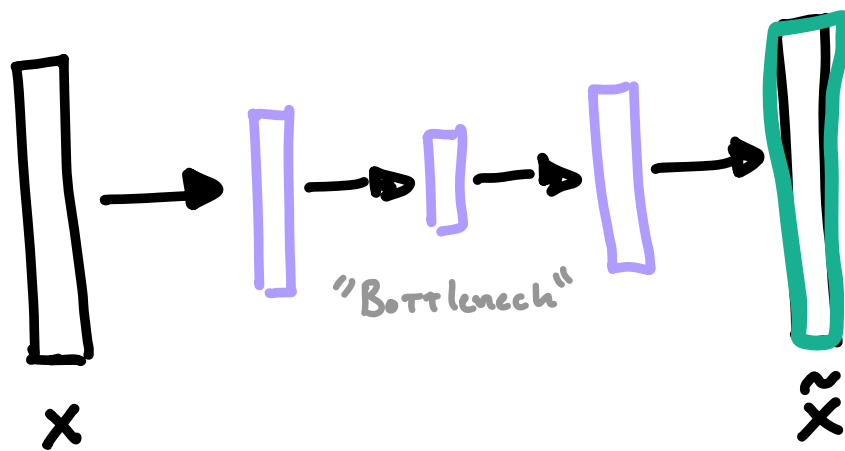
Want to be able to sample from some (unknown) distribution

$P^*(x)$

Cat image distribution

But this hard. (What would P look like for natural images?)

Auto-encoders!



Consume x , Compress it, Then Reconstruct it.

$$\begin{array}{ccc} \begin{array}{c} m \\ \text{[rectangle]} \\ x \end{array} & \xrightarrow{W} & \begin{array}{c} d \\ \text{[rectangle]} \\ z \end{array} \\ & & \xrightarrow{V} & \begin{array}{c} \text{[rectangle]} \\ \tilde{x} \end{array} \end{array}$$
$$\begin{array}{ll} z = Wx & \text{Encoder} \\ \tilde{x} = Vz & \text{Decoder} \end{array}$$

$(d \times m) \quad (m \times 1)$
 $(m \times d) \quad (d \times 1)$

Loss: Reconstruction error

$$\mathcal{L}(x, \tilde{x}) = \text{MSE}(x, \tilde{x})$$

De-noising Auto-encoders add noise first

$$z = Wx' \quad x' = x + \epsilon$$
$$\epsilon \sim \mathcal{N}(\vec{0}, \sigma^2)$$

Interpolating in latent space

$$h_1 \leftarrow \text{enc}(x_1)$$

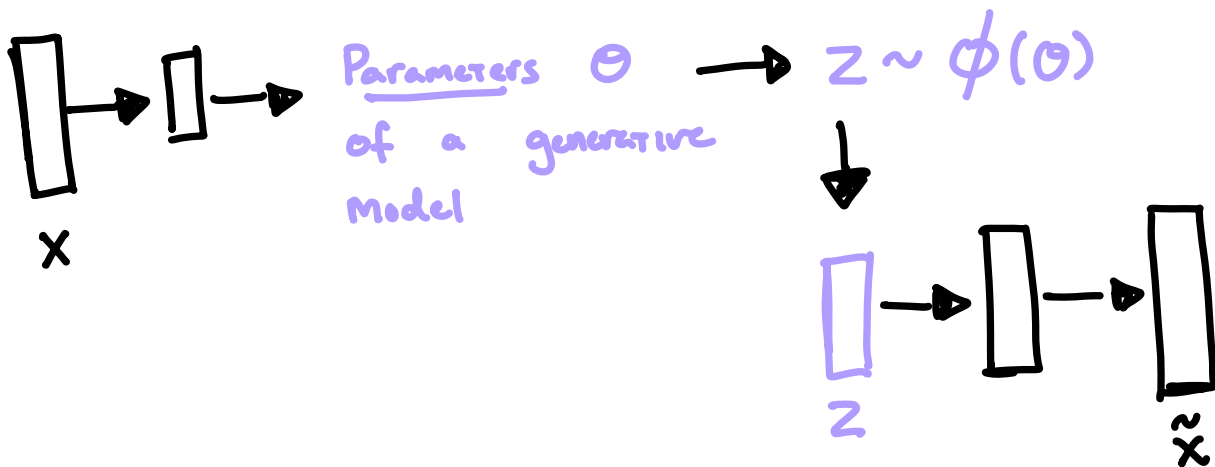
$$h_2 \leftarrow \text{enc}(x_2)$$

$$h_{1,2}^\lambda \leftarrow \lambda h_1 + (1-\lambda) h_2$$

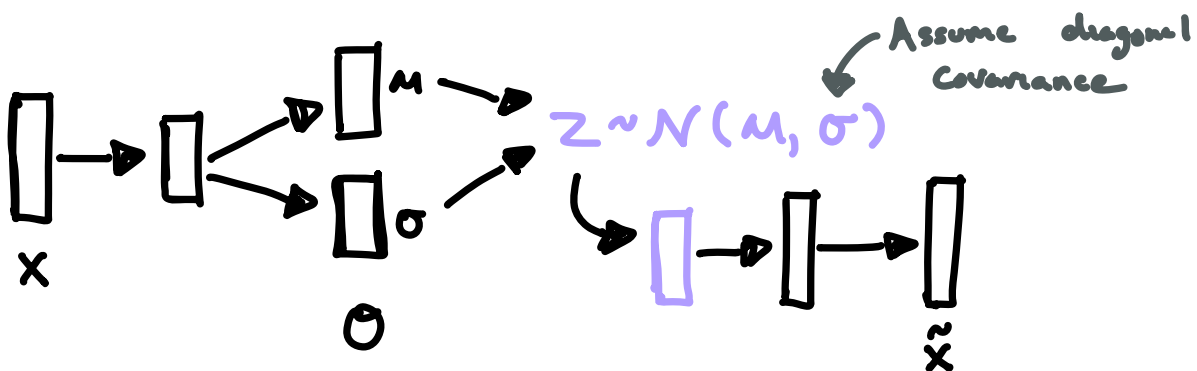
$$\tilde{x}_{1,2} \leftarrow \text{dec}(h_{1,2}^\lambda)$$

Variational Auto-encoders
 element into the process.

Inject a probabilistic



In practice: $\phi = \mathcal{N}(\mu, \sigma)$



For a Normal Distribution, This is equivalent to:

$$\begin{aligned}
 z &= \mu + \sigma \cdot \mathcal{N}(\phi, I) \\
 &= \underbrace{\mu + \sigma \cdot \epsilon}_{\text{ Torch Params. }}
 \end{aligned}$$

[See CoLab]

Diffusion Models approximate P
by iteratively sampling from
simple distributions.

Suppose we have a sample cat
image x_0 . In Gaussian Diffusion
we repeatedly noise this

$$x_{T+1} \leftarrow x_T + \eta_T$$

$$\eta_T \sim N(0, \sigma^2)$$



Now consider de-noising iteratively:
Given P_T , produce distribution P_{T-1} .
This is a Reverse Sampler.

IDEA Learn to reverse noising
one step at a time.

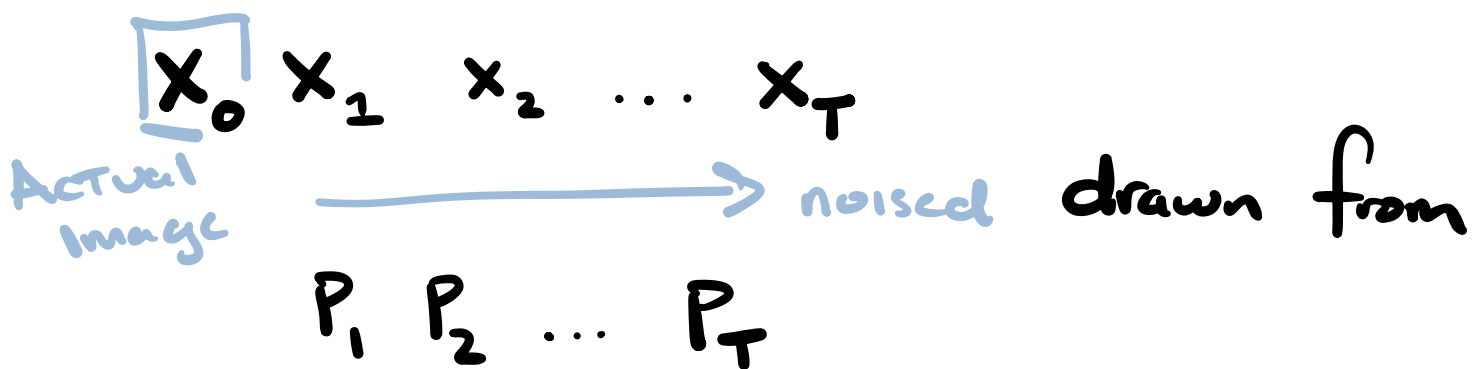
$$P(x_{T-1} | x_T = z)$$

Sample from P_T

But how? For Gaussian noising it
turns out (for small σ)

$$P(x_{T-1} | x_T = z) \approx \mathcal{N}(x_{T-1} | \mu_T, \sigma^2)$$

So we have



We want a Reverse Sampler to

$$\text{yield } \mu_T(z) \doteq \mathbb{E}[x_T | x_{T+\Delta T} = z]$$

Sample ΔT steps away

Estimate via learned $f_{\theta} : \mathbb{R}^d \rightarrow \mathbb{R}^d$

Learn via denoising objective

$$\min_{x_t, x_{t+\Delta t}} \mathbb{E} \left\| \underbrace{f_{\theta}(x_t + \Delta_t) - x_t}_{\text{Learn to denoise a sample}} \right\|_2^2$$

Because we
are sampling

Learn to denoise a sample

If we learn this, we can perform
one denoising step.

But then we can generate!

$$\hat{x}_T \sim \mathcal{N}(0, \sigma^2) \quad \text{sample pure noise}$$

For t in $\{T-1 \dots 0\}$

$$\hat{x}_t \sim \mathcal{N}(f_{\theta}(\hat{x}_{t+\Delta_t}), \sigma^2)$$

Return \hat{x}_T
Prior sample

How could we condition the
generation?