

# Generative Diffusion Models

Shervin Khalafi \*and Alejandro Ribeiro

## 1 Project Description

In this project you will first learn a bit about Diffusion models. You will then implement a basic Diffusion Model and train it on the MNIST dataset so that it learns to generate new images of handwritten digits that are similar (but not identical) to those found in the dataset.

## 2 What Is a Diffusion Model?

In this section we will go over everything you need to know for implementing a diffusion model. For a thorough explanation of a Diffusion model including all the math, see "Understanding Diffusion Models: A Unified Perspective" by Calvin Luo.

### 2.1 Generative Models

A generative model is a type of statistical model that is capable of generating new data instances that resemble the training data. These models learn the underlying distribution of the data in order to produce new instances that could have plausibly come from the same dataset.

---

\*If you have any questions regarding the project you can reach me at [shervink@seas.upenn.edu](mailto:shervink@seas.upenn.edu)

To be a bit more precise consider  $q(x)$  to be the probability density assigned to the point  $x \in \mathbb{R}^d$  by the underlying distribution of the data. But it often is the case that we don't know the underlying distribution  $q(\cdot)$ . We only have access to some data points sampled from it in the form of a dataset  $\mathcal{D} = \{x^{(1)}, x^{(1)}, \dots, x^{(M)}\}$ . The goal of a generative model is to learn a distribution  $p_\theta(x)$ , parameterized by  $\theta$  such that  $p_\theta(x)$  is "as close as possible" to  $q(x)$ . In other words, the goal is to learn to generate new samples that "could have been from the same dataset".

## 2.2 Latent Representation of the Data

For many of the types of data we typically encounter, we can think of the data we observe,  $x$  as represented or generated by an associated unseen *latent* variable, which we can denote by random variable  $z$ . The best intuition for expressing this idea is through Plato's Allegory of the Cave. In the allegory, a group of people are chained inside a cave their entire life and can only see the two-dimensional shadows projected onto a wall in front of them, which are generated by unseen three-dimensional objects passed before a fire. To such people, everything they observe is actually determined by higher-dimensional abstract concepts that they can never behold.

A diffusion model is one of many ways to model the relationship between the observed data  $x$  and the corresponding latents  $z$ .

## 2.3 Diffusion Models

In a Diffusion model, we call the true data observation  $x_0$  (for example a natural image from our dataset) and we add small amounts of noise to it in multiple time steps to get a series of noisy versions of the original sample  $x_1, x_2, \dots, x_T$ . Here, our latent  $z$  is essentially all of the noisy images  $x_t$  for  $t = 1, \dots, T$ .

Each noising step can be described as a conditional distribution:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)\mathbf{I}) \quad (1)$$

For those of you that are less familiar with probability, it suffices to know

that:

$$\mathbf{x}_t = \sqrt{\alpha_t} \cdot \mathbf{x}_{t-1} + (1 - \alpha_t) \cdot \epsilon \quad (2)$$

where  $\epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I})$  is just a noise sample from the standard Gaussian distribution. The parameters  $\alpha_t$  determine the size of the noise added in each time step and in this project are chosen as hyperparameters and are fixed during training. This choice of  $\alpha_t$  is called the "noise schedule".

In a nutshell we start from initial image  $\mathbf{x}_0$  and add a small amount of noise to it in each step until we reach an image  $\mathbf{x}_T$  in the final step that's essentially just pure Gaussian noise (see Figure (1)). This is called the forward process. It is relevant that the properties of the Gaussian distribution lets us sample an  $\mathbf{x}_t$  directly from  $\mathbf{x}_0$  without having to go through all the steps as follows:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \cdot \mathbf{x}_{t-1} + (1 - \bar{\alpha}_t) \cdot \epsilon_0 \quad (3)$$

where  $\bar{\alpha} = \prod_{i=1}^t \alpha_i$  and  $\epsilon_0 \sim \mathcal{N}(\epsilon_0; \mathbf{0}, \mathbf{I})$ . Now comes the hard part which is: How can we start from pure noise i.e.  $\mathbf{x}_T$  and iteratively denoise the image until we reach  $\mathbf{x}_0$ ? We call this the backward process and this is what we will train a model to do.

Assume we have  $\mathbf{x}_t$  and we want to find the most likely  $\mathbf{x}_{t-1}$  that could have generated this  $\mathbf{x}_t$ . Well, if we also knew the  $\mathbf{x}_0$  that was the starting point for generating both  $\mathbf{x}_t$  and  $\mathbf{x}_{t-1}$  we could compute this most likely  $\mathbf{x}_{t-1}$  since we know all the transition distributions  $q(\mathbf{x}_t | \mathbf{x}_{t-1})$ . We call this "most likely" prediction  $\mu_q(\mathbf{x}_t, \mathbf{x}_0, t)$  and we can derive it as:

$$\mu_q(\mathbf{x}_t, \mathbf{x}_0, t) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t} \quad (4)$$

But the problem is, during the backward process when we start from random noise  $\mathbf{x}_T$  and denoise it step by step, we don't know what  $\mathbf{x}_0$  is. This is where the neural network finally comes into play. We will call our best prediction for  $\mathbf{x}_{t-1}$  when we only have access to the current noisy image,  $\mu_\theta(\mathbf{x}_t, t)$ . In order to make this prediction as close as we can to  $\mu_q(\mathbf{x}_t, \mathbf{x}_0)$ , we choose it to be the following:

$$\mu_\theta(\mathbf{x}_t, t) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t) \cdot \hat{\mathbf{x}}_\theta(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t} \quad (5)$$

Note that the only difference between (4) and (5) is that in (5) we have replaced  $x_0$ , which we don't have access to in the backwards process, with  $\hat{x}_\theta(x_t, t)$ . Here,  $\hat{x}_\theta(x_t, t)$  represents a neural network parameterized by  $\theta$  that seeks to predict the original image  $x_0$  from noisy image  $x_t$  and time index  $t$ . So all we have to do is to train the neural network  $\hat{x}_\theta(x_t, t)$  to be good at predicting the original image  $x_0$  given noisy image  $x_t$  and time index  $t$ . With all this in mind, it makes sense to choose our loss function evaluated for a single image  $x_0$  to be the average error of our network in predicting  $x_0$ :

$$\hat{\ell}(\theta; x_0) = \mathbb{E}_{t \sim U\{2, T\}} \left[ \mathbb{E}_{q(x_t | x_0)} \left[ \|\hat{x}_\theta(x_t, t) - x_0\|_2^2 \right] \right] \quad (6)$$

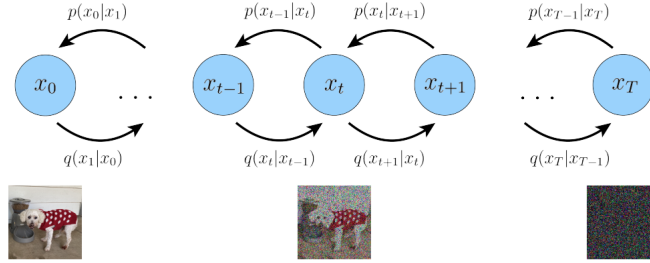
Note that there are two averages being computed in the loss above.

- The first one  $\mathbb{E}_{t \sim U\{2, T\}}$  is a uniform average over all time steps. This makes sense because we want our network to learn to predict the original image well at *all* time steps not just for any single time step. You can estimate this by sampling some time indices uniformly from 2 to  $T$ .<sup>1</sup>
- The second one  $\mathbb{E}_{q(x_t | x_0)}$  tells us that for any given  $t$ , we should average the loss over different possible noisy images  $x_t$  that could be generated by the forward process. We do this because since the forward process is stochastic, even when starting from the same  $x_0$  you could get a different  $x_t$  each time you run the forward process. Our model should learn to denoise all of these potential  $x_t$  images well.

In the implementation of our diffusion model, we will train a neural network  $\hat{x}_\theta(x_t, t)$  to minimize the loss  $\hat{\ell}(\theta; x_0)$ , averaged on samples  $x_0$  drawn from the MNIST dataset.

---

<sup>1</sup>The fact that we don't consider  $t = 1$  is a bit technical. Don't worry about it in this project



**Figure 1.** A Graphical representation of a Diffusion Model.

### 3 Notes on Implementation

To make your life a little easier, keep the following in mind when you are implementing the project:

- While you are free to use any neural network for  $\hat{x}_\theta(x_t, t)$  that takes as input a noisy image  $x_t$  and time index  $t$  and outputs an image with the same dimensions as  $x_t$ , we suggest that you use the "UNet2DModel" from the "diffusers" python library. This network is a time-conditioned UNet. For the purposes of this project you can consider a UNet to be similar to a CNN where the dimensions of the output and input image are the same.
- If training on the entire MNIST dataset is taking too long on your computer, you can train on a randomly selected subset of the dataset. Try to make this subset as large as you can given the limitations of your hardware.
- Regardless of the previous point, training might take a long time so account for that in your time management.
- For generating new samples using our models there are two different ways you could do it (Known as DDPM and DDIM sampling). You could either sample a standard Gaussian random noise  $x_T$  and then iteratively apply

$$x_{t-1} = \mu_\theta(x_t, t) \quad (7)$$

for  $t = T, T-1, \dots, 1$  until you reach  $x_0$ . Alternatively, you could do the stochastic version of (7) and add small mean-zero Gaussian

noise  $\epsilon_t$  at each denoising step:

$$x_{t-1} = \mu_\theta(x_t, t) + \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \epsilon_t \quad (8)$$

where  $\epsilon_t \sim \mathcal{N}(\epsilon_t; \mathbf{0}, \mathbf{I})$ . You can try both. What advantage/disadvantage do you think each method has?

- For the noise schedule, feel free to play around with different schedules. You can try the one used in the paper “Denoising Diffusion Probabilistic Models” by Ho et al. as a good baseline (note that  $\beta_t$  in the paper is equivalent to our  $1 - \alpha_t$ ).

## 4 Your Report

Your report should include the following for a full mark:

- Clear and concise explanations of each part of your implementation (comments in the code do not count).
- Attached code for each part of the implementation.
- Clear denoting of all the hyperparameters that you chose for your final model and the reasoning behind your choices.
- Comprehensive evaluation of the training process, including but not limited to multiple images generated by your model at different epochs/steps during the training process.