

DDPM 2020

Sunday, June 8, 2025 4:44 AM

DDPM Diffusion Model

High Level Overview

We seek to learn a distribution such that we can generate new images by sampling from this distribution.

High Level Visualization

Training

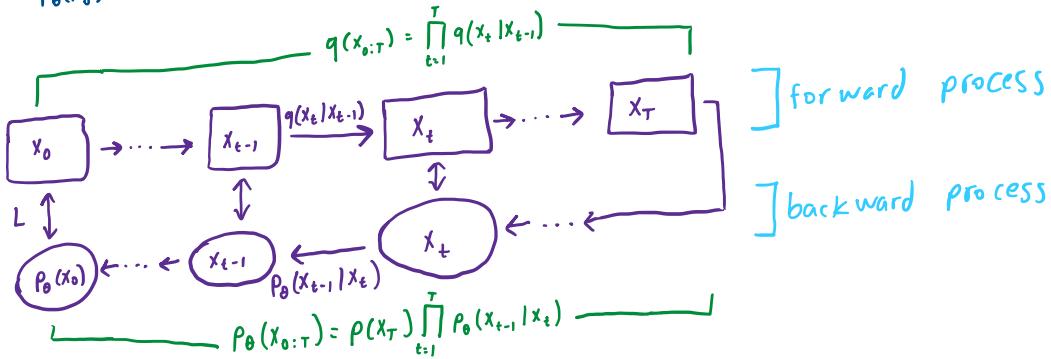


Inference



Terminology

- $q(x_0)$ ← Original data distribution
- x_0 ← Real image sampled from q
- x_T ← Pure noise image
- $p(x_T)$ ← Prior Distribution we start with for de-noising: $N(0, I)$
- $p_\theta(x_0)$ ← Generated denoised image



Forward Process (Adding Noise)

Formula 1: $q(x_t | x_{t-1}) = N(x_t; \sqrt{1-\beta_t} \cdot x_{t-1}, \beta_t I)$

conditional probability distribution of a noisy image given a less noisy image

mean variance

we sample x_t from a Normal distribution, where the mean is the previous (less noisy) image scaled by $\sqrt{1-\beta_t}$, and the variance is $\beta_t I$

Meaning: How likely is it to see a particular noisy image x_t , given that we start from x_{t-1} ?

Purpose: Defines how we add noise at each time step

Insight: We are not trying to learn images directly (there are infinitely many)—we are learning the distribution of images. We can then sample images from the learned distribution

$$\text{Formula 2 (using Formula 1): } q(x_{1:T} | x_0) = \prod_{t=1}^T q(x_t | x_{t-1})$$

joint probability of entire sequence of noisy images product of every conditional probability of a noisy image given its less noisy image
↓ ↓
Formula 1 given its less noisy image

Meaning: How likely is this sequence of noisy images (x_1, x_2, \dots, x_T) if we start from x_0 and apply the forward diffusion process?

Purpose: It tells us a probability distribution of the noisy images given a clean image—which lets us generate training data and define our training objective.

Insight: $q(x_{1:T} | x_0) =$ the true distribution. We are trying to estimate p_θ to be as close to q as possible, but we don't know p_θ exactly, because it's too complex. But since we defined the forward process, we know q . So we estimate p_θ by sampling from q , and use the samples from p_θ and q to compute loss.

Forward Process Trick

Instead of adding noise at each timestep, we can add all noise at once.

Given:

- d_1, d_2, \dots, d_T

we can compute:

$$\bar{d}_t = \prod_{s=1}^t d_s$$

\bar{d}_t = cumulative product of each d_s from 1 to t

we get:

$$\bar{d}_1, \bar{d}_2, \dots, \bar{d}_T$$

We need \bar{d}_t at each timestep for computing loss, but we

Reverse Process (Removing Noise)

$$\text{Formula 3: } p(x_T) = N(x_T; 0, I)$$

prior distribution of x_T Standard Gaussian centered at 0 with unit variance

Meaning: we define a prior distribution for the pure noise image, which is just the standard normal distribution

Purpose: we need a starting point for the reverse diffusion process. Since x_T is assumed to be pure noise, we can approximate it as $N(0, I)$. From here, we estimate p_θ

Insight: We need a prior because during inference, we don't have $q(x_T | x_0)$ because we don't have x_0 (this is what we are trying to generate), so we must start with $p(x_T)$, which is our best guess for $q(x_T)$.

$$\text{Formula 4: } p_\theta(x_{t-1} | x_t) = N(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

conditional probability of a less noisy image given a noisy image

mean → variance

We sample x_{t-1} from a normal distribution:

- the mean is a learnable parameter and depends on the noisy image and the timestep
- the variance is a fixed parameter that only depends on the timestep

Meaning: How likely is it that we see a less noisy image, x_{t-1} , given a more noisy image, x_t ?

Purpose: Tells us how to remove noise from an image

... we sample from it during training, where we then

Insight: We are learning a distribution, p_θ , and evaluating probabilities using it to compute loss, which then updates p_θ .

$$\text{Formula 5: } p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)$$

joint probability of generating the entire sequence of images

Formula 3
prior probability multiplied by the products of the conditional probabilities of less noisy images given their more noisy images

Formula 4

Meaning: Probability of generating this entire sequence starting from x_T

Purpose: We compare this probability value to that of the forward process, enabling us to estimate $p_\theta(x_0)$

Insight: The goal of the diffusion model is to estimate the distribution of the images, and p_θ represents our estimation of the true (unknown) distribution. We can't directly compute $p_\theta(x_0)$, so we start with our best guess for the noise, $p(x_T)$, and remove noise at each time step until we get $p_\theta(x_0)$. When we get $p_\theta(x_0)$, we compare $p_\theta(x_{0:T})$ with the true forward process, $q(x_{1:T}|x_0)$.

Training

We use the above formulas to optimize the negative log likelihood.

Negative Log Likelihood (NLL):

- Loss function
- If model is confident and correct, NLL is low (good)
- If model is wrong or unsure, NLL is high (bad)

Ex: We are trying to predict y given x , $p(y|x)$
Likelihood of observing data is $L = \prod_{i=1}^N p(y_i|x_i)$ (assuming independence)

Log-Likelihood

$$\log(L) = \sum_{i=1}^N \log(p(y_i|x_i)) \quad (\text{base-e})$$

$$\text{NLL} = -\log(L) = -\sum_{i=1}^N \log(p(y_i|x_i))$$

- we negate it because we want to minimise loss

computing NLL for our diffusion process:

$$E[-\log p_\theta(x_0)] \leq E_q \left[-\log \frac{p_\theta(x_0|x_T)}{q(x_0|x_T)} \right] = E_q \left[-\log p(x_T) - \sum_{t=1}^T \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_{t-1}|x_t)} \right] = L$$

our fully denoised image

KL Divergence

Given:

- $q(x)$: the true/target distribution
- $p(x)$: model/predicted distribution

$$D_{KL}(q(x)||p(x)) = \int q(x) \log \frac{q(x)}{p(x)} dx$$

- If $p(x)$ is close to $q(x)$, KL divergence is small
- If $p(x)$ is different from $q(x)$, KL divergence is large
- If $p(x)=q(x)$, KL divergence is 0
- Not symmetric

Computing Evidence Lower Bound (ELBO):

Problem:

We want to maximize $\log p_\theta(x_0)$, but we can't compute this directly.

$$\log p_\theta(x_0) = \log \int p_\theta(x_0, x_{1:T}) dx_{1:T}$$

what we are trying to learn

Solution: we compute a lower bound (ELBO):

$$\log p_\theta(x_0) \geq \text{ELBO}$$



$$L = E_q \left[D_{KL}(q(x_T | x_0) || p(x_T)) + \sum_{t>0} D_{KL}(q(x_{t-1} | x_t, x_0) || p_\theta(x_{t-1} | x_t)) - \log p_\theta(x_0 | x_1) \right]$$

KL divergence between the distribution that defines adding noise to x_0 and the prior distribution

KL divergence between the known reverse step distribution and the reverse step distribution modeled by the neural network

$$q(x_{t-1} | x_t, x_0) = N(x_{t-1}; M_\theta(x_t, x_0), \tilde{\beta}_t I)$$

$$\tilde{\beta}_t = \frac{1 - \bar{\beta}_{t-1}}{1 - \bar{\beta}_t} \beta_t$$

negative log likelihood of the true data point x_0 under the final reverse step

Purposes:

L_T : Ensures x_T (noisy image) matches pure gaussian noise

L_{t-1} : Trains model to denoise $x_t \rightarrow x_{t-1}$

$-\sum L_{t-1}$ represents all time steps we need

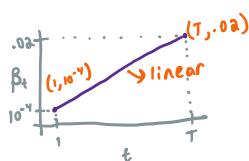
L_0 : Trains model to produce realistic final data $x_1 \rightarrow x_0$

Implementation Choices

forward choices and L_T

Authors choose to fix β_t instead of parameterizing them.

Instead, they fix them as constants:



Reverse Process and $L_{1:T-1}$

Given $p_\theta(x_{t-1} | x_t) = N(x_{t-1}; M_\theta(x_t, t), \Sigma_\theta(x_t, t))$, we have two choices for $\Sigma_\theta(x_t, t) = \sigma^2 I$, the learnable parameter for the variance given the image at the time step and the time step:

1) $\sigma^2 = \beta_t$ optimal if the original image $x_0 \sim N(0, I)$ (the images follow the gaussian distribution)

2) $\sigma^2 = \tilde{\beta}_t = \frac{1 - \bar{\beta}_{t-1}}{1 - \bar{\beta}_t} \beta_t$ optimal if the original image x_0 is not gaussian

Both choices have similar results per the paper

The mean, $M_\theta(x_t, t)$, needs to be parameterized, and the paper presents the following, given $p_\theta(x_{t-1} | x_t) = N(x_{t-1}; M_\theta(x_t, t), \sigma^2 I)$, we can rewrite L_{t-1} :

$$L_{t-1} = E_q \left[\frac{1}{\sigma^2} \| \tilde{M}(x_t, x_0) - M_\theta(x_t, t) \|^2 \right] + C$$

$M_\theta(x_t, t) \in \mathbb{R}^D$ ($D = \# \text{ of pixels/features}$)
to learn this vector with the Neural Network

norm of vector difference
between the known forward
posterior and what the neural
network predicts

C is a constant that doesn't depend on θ

$$L_{t-1} - C = E_{x_0, \epsilon} \left[\frac{1}{2\sigma_t^2} \left\| \tilde{M}_t(x_t(x_0, \epsilon), \frac{1}{\sqrt{\alpha_t}}(x_t(x_0, \epsilon) - \sqrt{1-\alpha_t}\epsilon) - \mu_\theta(x_t(x_0, \epsilon), t)) \right\|^2 \right]$$

$$= E_{x_0, \epsilon} \left[\frac{1}{2\sigma_t^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left(x_t(x_0, \epsilon) - \frac{\beta_t}{\sqrt{1-\alpha_t}} \epsilon \right) - \mu_\theta(x_t(x_0, \epsilon), t) \right\|^2 \right]$$

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1-\alpha_t}} \epsilon_\theta(x_t, t) \right)$$

$$E_{x_0, \epsilon} \left[\frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1-\alpha_t)} \| \epsilon - \epsilon_\theta(\sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \epsilon, t) \|^2 \right]$$

= mean square error between true noise ϵ and
neural network's predicted noise ϵ_θ

Data scaling, reverse process decoder, and L_0

We scale pixel values in $\{0, 1, \dots, 255\}$ linearly to $[-1, 1]$

Problem: $p(x_t | x_i)$ = density, not a probability, so the probability
of one point x_0 is zero.

Solution: we integrate over a small bin around x_0^i

$$p_\theta(x_0 | x_i) = \prod_{i=1}^D \frac{s_i(x_i)}{s_i(x_0^i)} \int_{x_0^i - \sigma_i}^{x_0^i + \sigma_i} N(x; M_\theta^i(x_i, 1), \sigma_i^2) dx$$

our choice for the variance

$s_+(x) = \begin{cases} 0 & \text{if } x=1 \\ x + \gamma_{\text{ass}} & \text{if } x<1 \end{cases}$

$s_-(x) = \begin{cases} -\infty & \text{if } x=-1 \\ x - \gamma_{\text{ass}} & \text{if } x>-1 \end{cases}$

probability of the clean image x_0 given the noisy image x_i

probability that pixel i matches x_0^i

what is the probability that every pixel of the model's Gaussian prediction would land in the correct bin

Meaning: The probability of the fully de-noised image, x_0 , given the last noisy image, x_i , is the product of probabilities that each generated pixel matches the true data pixel in x_0 for each pixel.

Purpose: Gives a way to compute L_0 in the loss equation

Paper Summary

- Our data distribution is $q(x_0)$
- We define $q(x_{0:T})$, which is a Markov chain that adds noise to the data
- Our goal is to train $p_\theta(x_{0:T})$, the reverse de-noising process such that we start with pure noise x_T and generate samples x_0 that match the true data distribution $q(x_0)$.
- Training is done by maximizing a lower bound on the data log-likelihood, which encourages each step of the reverse process to accurately reverse the noising process.

Questions/Points of Confusion

- How is the vanishing gradient problem avoided if we are taking the joint probability of thousands of small terms?
- Why do we need to structure the mean we are learning as $\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} (x_t - \frac{\beta_t}{\sqrt{1-\alpha_t}} \epsilon_\theta(x_t, t))$ instead of just $\mu_\theta(x_t, t) = \gamma$, where the neural network learns γ directly?
- I am still a little unclear on how predicting the probability of a sample or de-noising process, which are both numerical values, helps us ... distribution.