# **Applied ML** Diffusion Models











Troels C. Petersen (NBI)



"Statistics is merely a quantisation of common sense - Machine Learning is a sharpening of it!"

### Diffusion

Diffusion models is a class of models for computer vision that a capable of

- Image generation
- Image de-noising
- Image in-painting
- Image super-resolution

Diffusion models consist of three parts:

- 1. Forward diffusion
- 2. Reverse diffusion
- 3. Sampling procedure



They are build on the idea, that if you apply a known process (diffusion) on e.g. images, then by finding the reverse process, one can generate new images.

It is thus a **generative** model type. OpenAI's Dall-E 2 model - capable of generating images from text - is based on diffusion.

### Diffusion

Diffusion has a different ML architecture from the other "classic" methods of generative models. In diffusion one uses a repeated alteration (not NN layers).



From https://lilianweng.github.io/posts/2021-07-11-diffusion-models

### **Forward Process**



- No learning
  - Sample random numbers from a gaussian
  - Add them with some scaling to all of the pixels
  - Repeat T times
- End product pure noise  $\mathcal{N}(0,\mathbf{I})$

Noise

### **Forward Process**

- Noising is a Markov process and noise is Gaussian
- Can jump to any stage directly



 $\beta_t$  values schedule (i.e., the noise schedule) is designed such that  $\bar{\alpha}_T \to 0$  and  $q(\mathbf{x}_T | \mathbf{x}_0) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$ 

## **Reverse (de-noising) process**

We want to obtain  $q(x_{t-1} | x_t)$ . If  $\beta_t$  is small,  $q(x_{t-1} | x_t)$  will be Gaussian. This means that we can approximate  $q(x_{t-1} | x_t)$  with a neural network.



## Reverse (de-noising) process



- In practice, people predict amount of the noise added on the previous step directly
- $\mathscr{L} = \text{MSE}(\epsilon_{t-1} \epsilon_{\theta}(x_t, t))$ , where  $\epsilon_t$  is added noise

### **Diffusion Summary**

To some (simple) degree, the diffusion model setup can be summarised by the below two algorithms.

Algorithm 1 Training	Algorithm 2 Sampling
1: repeat 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 3: $t \sim \text{Uniform}(\{1, \dots, T\})$ 4: $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$ 5: Take gradient descent step on $\nabla_{\theta} \  \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}, t) \ ^2$ 6: until converged	1: $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$ 2: for $t = T,, 1$ do 3: $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$ if $t > 1$ , else $\mathbf{z} = 0$ 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\overline{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 5: end for 6: return $\mathbf{x}_0$

The training part is to train the algorithm to estimate the amount of noise added at each step.

The sampling part is to produce new images from sampling the noisy space and then using the trained part.

#### Diffusion example problem

Take images, and apply the diffusion forward and and then backwards in pixel space... i.e. NOT LATENT diffusion.

## Latent Diffusion "Overview"

Here is an attempt at labelling the whole process (& Latent) and all the parts...



This specific type of model can be considered a Stochastic Differential Equation (SDE).



Forward diffusion SDE:

$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_t dt + \sqrt{\beta(t)} d\boldsymbol{\omega}_t$$
  
drift term diffusion term

(pulls towards mode) (injects noise)

This specific type of model can be considered a Stochastic Differential Equation (SDE).



Putting the two parts together as ODE/SDE, one makes it possible to use advanced (existing) algorithms for solving these.





- Generative Reverse Diffusion SDE (stochastic):  $d\mathbf{x}_t = -\frac{1}{2}\beta(t) \left[\mathbf{x}_t + 2\mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x}_t, t)\right] dt + \sqrt{\beta(t)} \ d\bar{\boldsymbol{\omega}}_t$
- Accelerate generation
- Enables use of advanced ODE/SDE solvers
- Generative Probability Flow ODE (deterministic):  $d\mathbf{x}_t = -\frac{1}{2}\beta(t) \left[\mathbf{x}_t + \mathbf{s}_{\theta}(\mathbf{x}_t, t)\right] dt$

### **Example uses of DDIM**

Apart from images, it can be used in many other places for generation...



- <u>JetNet dataset</u> gluon jets, up to 30 constituents
- Number of constituents known in advance
- DDIM formulation solve SDE/ODE to generate data

### **Bonus Slides**

### Diffusion

Forward (noising) diffusion process (simple!):





From https://www.superannotate.com/blog/diffusion-models

### **Stochastic Differential Equation**

The process can be considered a Stochastic Differential Equation, which in small steps transforms the image to noise.



### **Score-based Generative Models**

In the reverse step, the model tries to reverse the process of noise addition. SGMs teaches the model to start from noisy data and progressively remove noise to reveal a more clear and detailed image.



Thus, the SDE maps data to a noise distribution (the prior), and reverse this SDE for generative modelling.

#### **Denoising diffusion probabilistic models**

DDPMs are a specific type of Diffusion Model, that focuses on removing noise from data in a probabilistic way.



During training, they learn how noise is added to data over time and how to reverse this process to recover the original data. This involves using probabilities to make educated guesses about what the data looked like before noise was added.

This approach is essential for the model's capability to accurately reconstruct data, ensuring the outputs aren't just noise-free but also closely resemble the original data.

The scoring is slightly complicated...

**Reverse diffusion SDE:** 
$$d\mathbf{x}_t = \left[-\frac{1}{2}\beta(t)\mathbf{x}_t - \beta(t)\nabla_{\mathbf{x}_t}\log q_t(\mathbf{x}_t)\right] dt + \sqrt{\beta(t)} d\bar{\boldsymbol{\omega}}_t$$

We can learn  $\nabla_{x_t} \log q_t(x_t)$  with NN, however direct regression is not possible. Instead, we can diffuse individual data points  $x_0$ . Diffused  $q_t(x_t | x_0)$  is tractable. **Denoising score matching:** 

$$\begin{split} \min_{\boldsymbol{\theta}} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_{0} \sim q_{0}(\mathbf{x}_{0})} \mathbb{E}_{\mathbf{x}_{t} \sim q_{t}(\mathbf{x}_{t} | \mathbf{x}_{0})} || \mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x}_{t}, t) - \nabla_{\mathbf{x}_{t}} \log q_{t}(\mathbf{x}_{t} | \mathbf{x}_{0}) ||_{2}^{2} \\ \\ \text{diffusion time } t \quad \text{sample } \mathbf{x}_{0} \quad \text{diffused data sample } \mathbf{x}_{t} \quad \text{neural network data sample} \\ \\ & \int_{\mathbf{x}_{t} = -\frac{1}{2}\beta(t)\mathbf{x}_{t} \, dt + \sqrt{\beta(t)} \, d\omega_{t}} \\ q_{t}(\mathbf{x}_{t} | \mathbf{x}_{0}) = \mathcal{N}(\mathbf{x}_{t}; \gamma_{t} \mathbf{x}_{0}, \sigma_{t}^{2} \mathbf{I}) \\ \\ & \min_{\boldsymbol{\theta}} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_{0} \sim q_{0}(\mathbf{x}_{0})} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \frac{1}{\sigma_{t}^{2}} || \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\mathbf{x}_{t}, t) ||_{2}^{2} \quad \mathbf{Sa} \end{split}$$

Same loss as before!