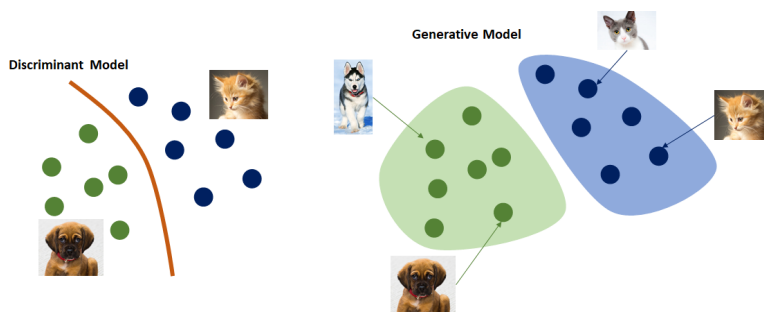# Generative Modelling | 11

## 11.1 Variational Autoencoders

### 11.1.1 Introduction

Variational Autoencoders (VAE) [9, 10] are one of the most commonly used and efficient approaches for the task of generative modelling. Unlike discriminative models, whose objective is to learn the conditional distribution $p(y|x)$, generative techniques aim to obtain the whole distribution of the data $\underline{x} = x_1, ..., x_n$, i.e. estimate $p(x)$ with the parametric distribution $p(x, \theta)$, typically differentiable with respect to $\theta$.

In this way it becomes possible to sample from such a probability distribution, thus generating new instances similar to the training dataset, or to assign to a given instance the probability of having been generated by the same process as the training dataset.

In general, generative modelling deals with a harder task with respect to discriminative models because in the former there is much more knowledge to learn. Intuitively, this is valid also for humans: by seeing many labelled images of dogs and cats it is much easier to distinguish dogs from cats than to draw a new image of a dog or a cat. In fact, in order to assign a label it is only required to find out some patterns which are different between the categories, while to generate a new instance it is necessary to capture correlations in the dataset, as for example the fact that dogs have two eyes, a tail, etc.. So, the discriminative model has to learn a decision boundary in the data space, while the generative one has to understand how data is placed throughout the space.

**Figure 11.1:** Image from *https://medium.com/@jordi299/about-generative-and-discriminative-models-d8958b67ad32*

It may be more convenient to avoid learning directly the distribution of the data (pixels in our example), but to introduce latent features $z$ which can describe the data and rewrite our distribution as:

$$p(x, z|\theta) = p(x|z, \theta)p(z)$$

We assume that the posterior distribution $p(z|x)$ is intractable and that we are working in the big data regime, so that we are forced to consider

mini-batches during the training phase. The goals that we would like to achieve are:

- ▶ Learning the parameters $\theta$
- ▶ Approximating the posterior $p(z|x)$ in order to understand which are the latent features extracted from a given instance $x$
- ▶ Performing approximate inference on $x$ to be able to fill holes in an instance

In this scenario, where we have a parametric distribution with latent variables, it might occur to us to use Expectation Maximization. But we must remember that the E-step requires being able to evaluate the conditional distribution $p(z|x, \theta_{OLD})$ which we have assumed here to be intractable. Other techniques that we have already seen, such as mean field variational approximation or sampling-based methods, turn out to be too computationally expensive as they do not scale well with large datasets. We will now see how to solve the problem by applying stochastic variational inference together with an encoder-decoder approach.

## 11.1.2 Autoencoding Variational Bayes (AEVB [11])

Following the idea of black-box variational inference, we start by considering a parametric variational distribution for the latent variables $q(z|\phi)$. We would like to optimize the ELBO jointly on $\theta$ and $\phi$ by stochastic gradient ascent and so we need to compute $\nabla_{\theta,\phi}\mathscr{L}(\phi, \theta)$. We fix the prior distribution of $z$ as a standard Gaussian:

$$p(z) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

and we start by rewriting the lower bound as:

$$
\begin{aligned}
\mathscr{L}(\phi, \theta) &= \mathbb{E}_{q(z|x,\phi)}[\log p(x, z|\theta) - \log q(z|x, \phi)] \\
&= \mathbb{E}_{q(z|x,\phi)}[\log p(x|z, \theta) + \log p(z) - \log q(z|x, \phi)] \\
&= \mathbb{E}_{q(z|x,\phi)}[\log p(x|z, \theta)] - \mathbb{E}_{q(z|x,\phi)}\left[\log \frac{q(z|x, \phi)}{p(z)}\right] \\
&= \mathbb{E}_{q(z|x,\phi)}[\log p(x|z, \theta)] - KL[q(z|x, \phi)||p(z)]
\end{aligned}
$$

As we said before, we are not able to work with the entire dataset and thus we have to consider a mini-batch $x_1, ..., x_m$ and approximate:

$$\mathscr{L}(\phi, \theta) \approx \frac{1}{m} \sum_{j=1}^{m} \mathbb{E}_{q(z|x_j,\phi)}[\log p(x_j|z, \theta)] - KL[q(z|x_j, \phi)||p(z)]$$
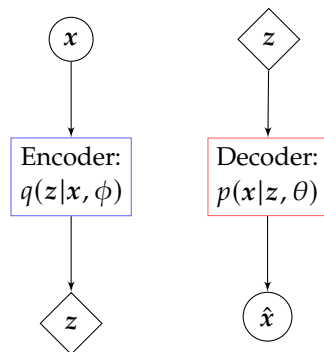
We can try to understand the meaning of this expression:

- ▶ The first term is the reconstruction error. It considers how well we are reconstructing $x$ given $z$ sampled from $q$. This is maximized when $p(x|z, \theta)$ assigns the highest probability to the original instance $x$. We call $p(x|z, \theta)$ **decoder**.

▶ The second term is a regularization term which encourages the variational distribution to look like a Gaussian distribution and not some kind of identity mapping. This avoids learning a mapping from inputs to latent features that just "stores" the inputs in different regions of the latent space, hence it favours generalization.

Both the terms involve the expectation with respect to the variational distribution of the latent variables that we interpret as features describing the data. So, we call $q(z|x, \phi)$ **encoder**.

We can represent the architecture in the following way:



where the input $x$ is mapped through the encoder into a useful latent space $z$ from which we can reconstruct $\hat{x}$ via the decoder. We would like $\hat{x}$ to be as similar as possible to $x$.
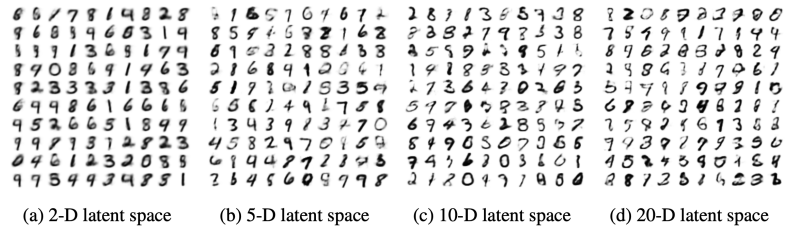
**Remark.** Notice that the variational parameters are shared across all the datapoints: using global parameters allows us to 'amortize' the cost of inference (**amortized inference**). Instead, in mean-field variational inference we had different parameters for each datapoint. This changes the behaviour of the model when we have a new datapoint: in mean-field VI we need to maximize the ELBO for each new point while here we can keep the global parameters fixed or run the variational inference again (maybe when many points are added).

So far we have obtained an expression for the ELBO, which is the objective function of our optimization. In order to maximize it, we need a good estimate of the gradient and we cannot obtain it directly from the previous expression because the gradient has to be computed on the same variables involved in the expectation. Therefore, it is not possible to swap the gradient and the expectation and this is where the **reparametrization trick** comes in.

We express the variational distribution $q(z|x, \phi)$ through a deterministic transformation $g(\varepsilon, x, \phi)$ which maps a noise variable $\varepsilon$ (distributed accordingly to a distribution we can easily sample from, as $\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$) to the distribution of $z$:

$$z = g(\varepsilon, x, \phi)$$

**Figure 11.2:** Random samples from learned generative models (AEVB) of MNIST for different dimensionalities of latent space [11].

(a) 2-D latent space     (b) 5-D latent space     (c) 10-D latent space     (d) 20-D latent space

Typically, the variational distribution is chosen to be a Gaussian:

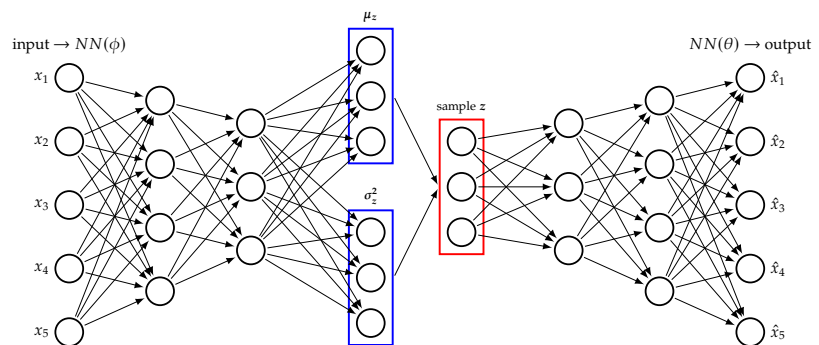$$q(z|x, \phi) = \mathcal{N}(\mu_z(x, \phi), \sigma_z{}^2(x, \phi) \cdot I)$$

so that we have an analytical expression for $KL[q(z|x, \phi)||p(z)]$ and we can rewrite $z$ as

$$z = \mu_z(x, \phi) + \sigma_z{}^2(x, \phi) \cdot \varepsilon$$

At this point, we can evaluate the gradient of the ELBO as:

$$\nabla_{\theta, \phi} \mathcal{L}(\phi, \theta) = \mathbb{E}_\varepsilon[\nabla_{\phi, \theta} \log p(x|g(\varepsilon, x, \phi), \theta)] - \nabla_\phi KL[q(z|x, \phi)||p(z)]$$

We still have to choose the functions $\mu_z(x, \phi)$ and $\sigma_z{}^2(x, \phi)$. To have high expressiveness and efficient optimization, we can opt for neural networks: we just built a **variational autoencoder**.



**Remark.** The variational autoencoder can be interpreted as a directed probabilistic graphical model with latent variables.

## 11.2 Diffusion Models

Diffusion models are a family of probabilistic generative models that progressively destruct data by injecting noise, then learn to reverse this process for sample generation. In particular, we will discuss denoising diffusion probabilistic models (DDPMs), but many concepts are common to other formulations.
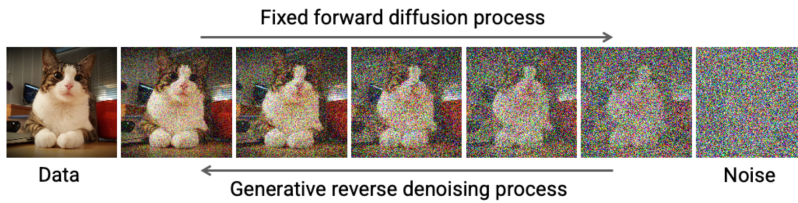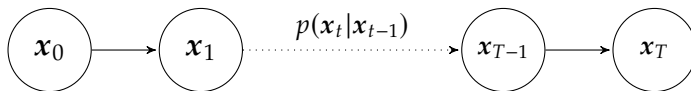


**Figure 11.3:** Image from [12], *Denoising Diffusion-based Generative Modeling: Foundations and Applications*

### 11.2.1 Forward Diffusion Process

Given a data point $x_0 \in \mathcal{X}$ sampled from the data distribution $p(x_0)$, consider a Markov chain $p(x_t|x_{t-1})$ with $t \in 1, \ldots, T$ satisfying the following properties:

▶ It is easy to sample from $p(x_t|x_{t-1})$
▶ For $T$ large enough, $p(x_T)$ is approximately a known distribution from which it is easy to sample and that does not depend on $x_0$, i.e, $p(x_T) \approx p(x_T|x_0)$. This distribution is referred as the **prior**.



Such Markov chain is referred as the **forward diffusion process**, and $p(x_t|x_{t-1})$ is often referred as the **forward transition kernel**. Using the chain rule of probability and the Markov property, we can factorize the joint distribution of $x_0, \ldots, x_T$ into:

$$p(x_0, \ldots, x_T) = p(x_0) \prod_{t=1}^{T} p(x_t|x_{t-1})$$

In practice, the most used forward process in continuous space is based on the injection of Gaussian noise:

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

where $\beta_t \in (0, 1)$ regulates the amount of noise that is injected at each step (the larger, the faster the convergence to the prior distribution). A useful property of the above process is that we can sample at any arbitrary time step in a closed form using the reparameterization trick. Given $\alpha_t := 1 - \beta_t$, $\bar{\alpha}_t := \prod_{i=1}^{t} \alpha_i$, and $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ we have

$$
\begin{aligned}
\mathbf{x}_t &= \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1 - \alpha_t}\epsilon_{t-1} \\
&= \sqrt{\alpha_t}\left(\sqrt{\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}}\epsilon_{t-2}\right) + \sqrt{1 - \alpha_t}\epsilon_{t-1} \\
&= \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{1 - \alpha_t\alpha_{t-1}}\hat{\epsilon}_{t-2} \\
&= \ldots \\
&= \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\hat{\epsilon}
\end{aligned}
$$

where $\epsilon_{t-2}$ and $\epsilon_{t-1}$ have been merged into $\hat{\epsilon}_{t-2}$.

Hence, we have

$$
p(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, 1 - \bar{\alpha}_t)
$$

This is also useful to show, as $\bar{\alpha}_t \to 0$, that

$$
\lim_{t \to \infty} p(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \mathbf{0}, \mathbf{I})
$$

So for a sufficiently large $T$, $\mathbf{x}_T \overset{\cdot}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$. In other words, this process ends up in corrupting the data into white noise.

### 11.2.2 Backward Process

Now that we are able to gradually corrupt data points $\mathbf{x}_0$ into noise, for generating new data samples we have to revert this process. We have seen that we can easily sample from $p(\mathbf{x}_T)$, but sampling from

$$
p(\mathbf{x}_0, \ldots, \mathbf{x}_T) = p(\mathbf{x}_T)\prod_{t=1}^{T} p(\mathbf{x}_{t-1}|\mathbf{x}_t)
$$

is non-trivial, since we do not have access to the reverse transition kernels $p(\mathbf{x}_{t-1}|\mathbf{x}_t)^*$ that would allow us to perform ancestral sampling as in the forward process. Therefore, the reverse transition kernels have to be learned. In order to do this, we fit parametric distributions $q_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$, from which it easy to sample and that can be easily computed, for example:

$$
q_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))
$$

where the mean and covariance functions are neural networks.

The objective is then to find parameters $\theta$ such that

$$
q_\theta(\mathbf{x}_0, \ldots, \mathbf{x}_T) := p(\mathbf{x}_T)\prod_{t=1}^{T} q_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \approx p(\mathbf{x}_0, \ldots, \mathbf{x}_T)
$$

---

* It is easy to prove that the reverse of a Markov chain is also a Markov chain.

The parametric kernels are learned by minimizing the Kullback-Leibler divergence between the forward and backward processes:

$$D_{KL}(p(\boldsymbol{x}_0,\ldots,\boldsymbol{x}_T)||q_\theta(\boldsymbol{x}_0,\ldots,\boldsymbol{x}_T)) = \mathbb{E}_{p(\boldsymbol{x}_0,\ldots,\boldsymbol{x}_T)}\left[\log\frac{p(\boldsymbol{x}_0,\ldots,\boldsymbol{x}_T)}{q_\theta(\boldsymbol{x}_0,\ldots,\boldsymbol{x}_T)}\right]$$

$$= -\mathbb{E}_{p(\boldsymbol{x}_0,\ldots,\boldsymbol{x}_T)}\left[\log q_\theta(\boldsymbol{x}_0,\ldots,\boldsymbol{x}_T)\right] + \text{const}$$

$$= -\mathbb{E}_{p(\boldsymbol{x}_0,\ldots,\boldsymbol{x}_T)}\left[\sum_{t=0}^{T}\log q_\theta(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t)\right] + \text{const}$$

**Remark.** The forward process $p(\boldsymbol{x}_0,\ldots,\boldsymbol{x}_T)$ was considered constant with respect to the parameters $\theta$. Although parameters of the forward process (as $\beta_t$) are usually considered fixed (hyperparameters), it is also possible to learn them. In that case we cannot consider the forward process as constant with respect to the parameters.

Minimizing this inverse KL divergence is equivalent to solving a maximum likelihood estimation problem: we can sample full trajectories starting from samples from the dataset and use them to fit parametric functions $q_\theta$ by stochastic optimization, as in a supervised learning problem.

**Remark.** It can be shown that, if forward transition kernels corrupts the data slowly enough ($T$ is large), then the reverse transition kernels will be approximately Gaussian. This is what makes the approximation of reverse transition kernels possible and approachable as a prediction problem. Conversely, if forward transition kernels add too much noise, the reverse transition probabilities can be multimodal, hence a parametric approximation would be intractable.

### 11.2.3 Denoising parametrization

There is empirical evidence that learning directly the backward transition kernels $q_\theta(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t)$ is not the most effective strategy. Alternatively, it is possible to write the backward transition kernel in the following way:

$$p(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t) = \int_{\boldsymbol{x}_0\in\mathcal{X}}p(\boldsymbol{x}_{t-1},\boldsymbol{x}_0|\boldsymbol{x}_t)d\boldsymbol{x}_0 = \int_{\boldsymbol{x}_0\in\mathcal{X}}p(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t,\boldsymbol{x}_0)p(\boldsymbol{x}_0|\boldsymbol{x}_t)d\boldsymbol{x}_0$$

This means that we can sample $\boldsymbol{x}_{t-1}$ if we can sample $\boldsymbol{x}_0 \sim p(\boldsymbol{x}_0|\boldsymbol{x}_t)$ and plug it into $p(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t,\boldsymbol{x}_0)$. It is possible to show that for the Gaussian diffusion process introduced above the distribution $p(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t,\boldsymbol{x}_0)$ is also a Gaussian and can be computed analytically:

$$p(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t,\boldsymbol{x}_0) = \mathcal{N}(\boldsymbol{x}_{t-1};\mu_t(\boldsymbol{x}_t,\boldsymbol{x}_0),\sigma_t\boldsymbol{I})$$

$$\mu_t(\boldsymbol{x}_t,\boldsymbol{x}_0) = \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}\boldsymbol{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t}\boldsymbol{x}_0$$

$$\sigma_t = \frac{(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \beta_t$$

Instead $p(x_0|x_t)$ (the denoising kernel) has to be learned, again by fitting a parametric distribution $q_\theta(x_0|x_t)$. Interestingly, this works despite $p(x_0|x_t)$ being highly multimodal and non-Gaussian, so only a "mean" approximation is possible.

However, it is common in DDPMs to avoid predicting directly $x_0$ given $x_t$. In fact, a model is trained to predict the noise that was added to $x_0$ in order to obtain $x_t$. Recall the property following from $p(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, 1 - \bar{\alpha}_t)$:

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t \text{ with } \epsilon_t \sim \mathcal{N}(0, I)$$

We can rewrite the mean of $p(x_{t-1}|x_t, x_0)$, or $p(x_{t-1}|x_t, \epsilon_t)$, as

$$\mu_t(x_t, \epsilon_t) = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_t\right)$$

Now we can fit a parametric model $\epsilon_\theta(x_t, t)$ to approximate $p(\epsilon_t|x_t)$. This is usually done in a simplified way by minimizing the squared error:

$$\mathcal{L}(\theta) = \mathbb{E}_{t,x_0,\epsilon_t}\left[\|\epsilon_t - \epsilon_\theta(x_t, t)\|^2\right]$$

$$= \mathbb{E}_{t,x_0,\epsilon_t}\left[\left\|\epsilon_t - \epsilon_\theta\left(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t, t\right)\right\|^2\right]$$

where the time step $t$ is sampled uniformly in $\{1, \ldots, T\}$, $x_0$ is randomly sampled from the dataset, and $\epsilon_t \sim \mathcal{N}(0, I)$. The training and sampling algorithm are then summarized as follows:

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

**Figure 11.4:** The training and sampling algorithms in DDPMs (image from [13]).

**Remark.** The empirical advantage of predicting the noise $\epsilon_t$ instead of $x_0$ could be due to the fact that it is easier to train a neural network when the target output is (marginally) distributed according to a normal distribution.

## Summary

A DDPM makes use of two Markov chains: a forward chain that perturbs data to noise, and a reverse chain that converts noise back to data. The forward Markov chain is designed with the goal to gradually transform any data distribution into a simple prior distribution (e.g., standard Gaussian). The backward Markov chain instead reverses the former starting from the known distribution and gradually converging to the

data distribution. Since the backward Markov chain is unknown, it is learned by a deep neural network using examples generated with the forward process. New data points are subsequently generated by first sampling a random vector from the prior distribution, followed by ancestral sampling through the reverse Markov chain.

### 11.2.4 Diffusion in discrete space

Despite being originally thought for data in continuous space, it is also possible to define a diffusion process in discrete space [14], even though its usefulness is questionable and subject of current research. Assuming without loss of generality data $x \in \{1, 2, \ldots, c\}^d$, there are two main kinds of discrete diffusion processes:

▶ **Uniform**: At each time step $t$, every component $x_t^i$ of $x_t$ switches to a random value with a given small probability $\epsilon_t$. The prior distribution of such diffusion process is the uniform over $\{1, 2, \ldots, c\}^d$.
▶ **Absorbing**: At each time step $t$, every component $x_t^i$ of $x_t$ switches to a particular value called the "mask", often referred as the [MASK] token. The prior distribution of such diffusion process is the single point $[MASK]^d$, that is the absorbing state of the Markov chain. A variant of the absorbing diffusion process is when at each time step $t$ the $t^{th}$ component $x_t^t$ is masked. Modelling the corresponding reverse process is then equivalent to fitting an autoregressive model.

### 11.2.5 Score-based diffusion models

Given the objective to sample from the distribution $p(x)$ that generated the data, score-based models models aim at estimating the so called *score* of $p(x)$, defined as $\nabla_x \log p(x)$ and then use sampling techniques that exploit the knowledge of the score of the distribution.

An example of score-based model is diffusion with stochastic differential equations (SDE), a generalization in continuous time of the Markov chain diffusion process discussed above, where the transition kernel $p(x_{t+\epsilon}|x_t)$ is implicitly defined through the following SDE:

$$\mathrm{d}x = f(x, t)\mathrm{d}t + g(t)\mathrm{d}w$$

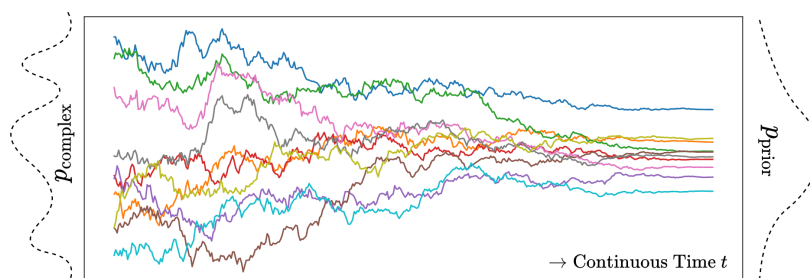where $\mathrm{d}w \sim \mathcal{N}(0, I\mathrm{d}t)$ is a Wiener process.



**Figure 11.5:** Image from [16], *An introduction to Diffusion Probabilistic Models.*

Generation of new data points reduces to solving the inverse-time SDE, that it can be proven to be:

$$\mathrm{d}x = \left[ f(x, t) - g^2(t) \nabla_x \log p_t(x) \right] \mathrm{d}t + g(t)\mathrm{d}w$$

For this we need to estimate the time dependent score $\nabla_x \log p_t(x)$, for example using a neural network $s_\theta(x, t)$. The most popular method is *denoising score matching*, where the following loss in minimized:

$$\mathbb{E}_{t \sim u(0,1), x_0 \sim p(x_0), x_t \sim p(x_t|x_0)} \left\| s_\theta(x_t, t) - \nabla_{x_t} \ln p(x_t \mid x_0) \right\|$$

that is deeply connected with the denoising loss of DDPMs.