

Advanced machine learning and data analysis for the physical sciences

Morten Hjorth-Jensen^{1,2}

Department of Physics and Center for Computing in Science Education,
University of Oslo, Norway¹

Department of Physics and Astronomy and Facility for Rare Isotope Beams,
Michigan State University, East Lansing, Michigan, USA²

May 7, 2024

© 1999-2024, Morten Hjorth-Jensen. Released under CC Attribution-NonCommercial 4.0 license

Plans for the week of May 6-10, 2024

Generative models

1. Finalizing discussion of Generative Adversarial Networks
2. Mathematics of diffusion models and selected examples
3. Video of lecture
4. Whiteboard notes

Reading on diffusion models

1. A central paper is the one by Sohl-Dickstein et al, Deep Unsupervised Learning using Nonequilibrium Thermodynamics, <https://arxiv.org/abs/1503.03585>
2. Calvin Luo at <https://arxiv.org/abs/2208.11970>
3. See also Diederik P. Kingma, Tim Salimans, Ben Poole, Jonathan Ho, Variational Diffusion Models, <https://arxiv.org/abs/2107.00630>

Reminder from last week, what is a GAN?

A GAN is a deep neural network which consists of two networks, a so-called generator network and a discriminating network, or just discriminator. Through several iterations of generation and discrimination, the idea is that these networks will train each other, while also trying to outsmart each other.

In its simplest version, the two networks could be two standard neural networks with a given number of hidden of hidden layers and parameters to train. The generator we have trained can then be used to produce new images.

Labeling the networks

For a GAN we have:

1. a discriminator D estimates the probability of a given sample coming from the real dataset. It attempts at discriminating the trained data by the generator and is optimized to tell the fake samples from the real ones (our data set). We say a discriminator tries to distinguish between real data and those generated by the abovementioned generator.
2. a generator G outputs synthetic samples given a noise variable input z (z brings in potential output diversity). It is trained to capture the real data distribution in order to generate samples that can be as real as possible, or in other words, can trick the discriminator to offer a high probability.

At the end of the training, the generator can be used to generate for example new images. In this sense we have trained a model which can produce new samples. We say that we have implicitly defined a probability.

Which data?

GANs are generally a form of unsupervised machine learning, although they also incorporate aspects of supervised learning. Internally the discriminator sets up a supervised learning problem. Its goal is to learn to distinguish between the two classes of generated data and original data. The generator then considers this classification problem and tries to find adversarial examples, that is samples which will be misclassified by the discriminator.

Semi-supervised learning

One can also design GAN architectures which work in a semi-supervised learning setting. A semi-supervised learning environment includes both labeled and unlabeled data. See https://proceedings.neurips.cc/paper_files/paper/2016/file/8a3363abe792db2d8761d6403605aeb7-Paper.pdf for a further discussion.

Thus, GANs can be used both on labeled and on unlabeled data and are used in three most commonly used contexts, that is

1. with labeled data (supervised training)
2. with unlabeled data (unsupervised learning)
3. a with a mix labeled and unlabeled data

Improving functionalities

These two models compete against each other during the training process: the generator G is trying hard to trick the discriminator, while the critic model D is trying hard not to be cheated. This interesting zero-sum game between two models motivates both to improve their functionalities.

Setup of the GAN

We define a probability $p_{\mathbf{h}}$ which is used by the generator. Usually it is given by a uniform distribution over the input \mathbf{h} . Thereafter we define the distribution of the generator which we want to train, p_g . This is the generator's distribution over the data \mathbf{x} . Finally, we have the distribution p_r over the real sample \mathbf{x} .

Optimization part

On one hand, we want to make sure the discriminator D 's decisions over real data are accurate by maximizing $\mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})}[\log D(\mathbf{x})]$.

Meanwhile, given a fake sample $G(\mathbf{h})$, $\mathbf{h} \sim p_h(\mathbf{h})$, the discriminator is expected to output a probability, $D(G(\mathbf{h}))$, close to zero by maximizing $\mathbb{E}_{\mathbf{h} \sim p_h(\mathbf{h})}[\log(1 - D(G(\mathbf{h})))]$.

On the other hand, the generator is trained to increase the chances of D producing a high probability for a fake example, thus to minimize $\mathbb{E}_{\mathbf{h} \sim p_h(\mathbf{h})}[\log(1 - D(G(\mathbf{h})))]$.

Minimax game

When combining both aspects together, D and G are playing a **minimax game** in which we should optimize the following loss function:

$$\begin{aligned}\min_G \max_D L(D, G) &= \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{h} \sim p_h(\mathbf{h})} [\log(1 - D(G(\mathbf{h})))] \\ &= \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g(\mathbf{x})} [\log(1 - D(\mathbf{x}))]\end{aligned}$$

where $\mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} [\log D(\mathbf{x})]$ has no impact on G during gradient descent updates.

Optimal value for D

Now we have a well-defined loss function. Let's first examine what is the best value for D .

$$L(G, D) = \int_{\mathbf{x}} \left(p_r(\mathbf{x}) \log(D(\mathbf{x})) + p_g(\mathbf{x}) \log(1 - D(\mathbf{x})) \right) d\mathbf{x}$$

Best value of D

Since we are interested in what is the best value of $D(\mathbf{x})$ to maximize $L(G, D)$, let us label

$$\tilde{\mathbf{x}} = D(\mathbf{x}), A = p_r(\mathbf{x}), B = p_g(\mathbf{x})$$

Integral evaluation

The integral (we can safely ignore the integral because \mathbf{x} is sampled over all the possible values) is:

$$\begin{aligned}f(\tilde{\mathbf{x}}) &= A \log \tilde{\mathbf{x}} + B \log (1 - \tilde{\mathbf{x}}) \\ \frac{df(\tilde{\mathbf{x}})}{d\tilde{\mathbf{x}}} &= A \frac{1}{\tilde{\mathbf{x}}} - B \frac{1}{1 - \tilde{\mathbf{x}}} \\ &= \frac{A - (A + B)\tilde{\mathbf{x}}}{\tilde{\mathbf{x}}(1 - \tilde{\mathbf{x}})}.\end{aligned}$$

Best values

If we set

$$\frac{df(\tilde{\mathbf{x}})}{d\tilde{\mathbf{x}}} = 0,$$

we get the best value of the discriminator:

$$D^*(\mathbf{x}) = \tilde{\mathbf{x}}^* = \frac{A}{A+B} = \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} \in [0, 1].$$

Once the generator is trained to its optimal, p_g gets very close to p_r . When $p_g = p_r$, $D^*(\mathbf{x})$ becomes 1/2. We will observe this when running the code from last week (see jupyter-notebook from week 15).

At their optimal values

When both G and D are at their optimal values, we have $p_g = p_r$ and $D^*(\mathbf{x}) = 1/2$, the loss function becomes

$$\begin{aligned} L(G, D^*) &= \int_{\mathbf{x}} \left(p_r(\mathbf{x}) \log(D^*(\mathbf{x})) + p_g(\mathbf{x}) \log(1 - D^*(\mathbf{x})) \right) d\mathbf{x} \\ &= \log \frac{1}{2} \int_{\mathbf{h}} p_r(\mathbf{x}) d\mathbf{x} + \log \frac{1}{2} \int_{\mathbf{x}} p_g(\mathbf{x}) d\mathbf{x} \\ &= -2 \log 2 \end{aligned}$$

What does the Loss Function Represent?

The JS divergence between p_r and p_g can be computed as:

$$\begin{aligned}D_{JS}(p_r \| p_g) &= \frac{1}{2} D_{KL}(p_r \| \frac{p_r + p_g}{2}) + \frac{1}{2} D_{KL}(p_g \| \frac{p_r + p_g}{2}) \\&= \frac{1}{2} \left(\log 2 + \int_{\mathbf{x}} p_r(\mathbf{x}) \log \frac{p_r(\mathbf{x})}{p_r + p_g(\mathbf{x})} d\mathbf{x} \right) + \\&\quad \frac{1}{2} \left(\log 2 + \int_{\mathbf{x}} p_g(\mathbf{x}) \log \frac{p_g(\mathbf{x})}{p_r + p_g(\mathbf{x})} d\mathbf{x} \right) \\&= \frac{1}{2} \left(\log 4 + L(G, D^*) \right)\end{aligned}$$

What does the loss function quantify?

We have

$$L(G, D^*) = 2D_{JS}(p_r \| p_g) - 2 \log 2.$$

The loss function of GANs quantifies the similarity between the generative data distribution p_g and the real sample distribution p_r by the so-called JS divergence when the discriminator is optimal. The best G^* that replicates the real data distribution leads to the minimum $L(G^*, D^*) = -2 \log 2$.

Problems with GANs

Although GANs have achieved great success in the generation of realistic images, the training is not easy; The process is known to be slow and unstable.

Hard to reach equilibrium.

Two models are trained simultaneously to an equilibrium to a two-player non-cooperative game. However, each model updates its cost independently with no respect to another player in the game. Updating the gradient of both models concurrently cannot guarantee a convergence.

Vanishing Gradient

When the discriminator is perfect, we are guaranteed with

$$D(\mathbf{x}) = 1, \forall \mathbf{x} \in p_r \text{ and } D(\mathbf{x}) = 0, \forall \mathbf{x} \in p_g.$$

Then, the loss function L falls to zero and we end up with no gradient to update the loss during learning iterations. One can encounter situations where the discriminator gets better and the gradient vanishes fast.

As a result, training GANs may face the following problems

1. If the discriminator behaves badly, the generator does not have accurate feedback and the loss function cannot represent the real data
2. If the discriminator does a great job, the gradient of the loss function drops down to close to zero and the learning can become slow

Improved GANs

One of the solutions to improved GANs training, is the introduction of what is called the Wasserstein distance, which is a way to compute the difference/distance between two probability distributions. For those interested in reading more, we recommend for example chapter 17 of Rashcka's et al textbook, Machine Learning with PyTorch and Scikit-Learn, chapter 17, see

<https://github.com/rasbt/python-machine-learning-book-3rd-edition/tree/master/ch17>

For a definition of the Wasserstein distance, see for example

<https://arxiv.org/pdf/2103.01678>

More codes

1. For codes and background, see Raschka et al, Machine Learning with PyTorch and Scikit-Learn, chapter 17, see <https://github.com/rasbt/python-machine-learning-book-3rd-edition/tree/master/ch17> for codes
2. Babcock and Bali, Generative AI with Python and TensorFlow2, chapter 6 and codes at https://github.com/raghavbali/generative_ai_with_tensorflow/blob/master/Chapter_6/conditional_gan.ipynb
3. See also Foster's text Generative Deep Learning and chapter 4 with codes at https://github.com/davidADSP/Generative_Deep_Learning_2nd_Edition/tree/main/notebooks/04_gan

Diffusion models, basics

Diffusion models are inspired by non-equilibrium thermodynamics. They define a Markov chain of diffusion steps to slowly add random noise to data and then learn to reverse the diffusion process to construct desired data samples from the noise. Unlike VAE or flow models, diffusion models are learned with a fixed procedure and the latent variable has high dimensionality (same as the original data).

Problems with probabilistic models

Historically, probabilistic models suffer from a tradeoff between two conflicting objectives: *tractability* and *flexibility*. Models that are *tractable* can be analytically evaluated and easily fit to data (e.g. a Gaussian or Laplace). However, these models are unable to aptly describe structure in rich datasets. On the other hand, models that are *flexible* can be molded to fit structure in arbitrary data. For example, we can define models in terms of any (non-negative) function $\phi(\mathbf{x})$ yielding the flexible distribution $p(\mathbf{x}) = \frac{\phi(\mathbf{x})}{Z}$, where Z is a normalization constant. However, computing this normalization constant is generally intractable. Evaluating, training, or drawing samples from such flexible models typically requires a very expensive Monte Carlo process.

Diffusion models

Diffusion models have several interesting features

- ▶ extreme flexibility in model structure,
- ▶ exact sampling,
- ▶ easy multiplication with other distributions, e.g. in order to compute a posterior, and
- ▶ the model log likelihood, and the probability of individual states, to be cheaply evaluated.

Original idea

In the original formulation, one uses a Markov chain to gradually convert one distribution into another, an idea used in non-equilibrium statistical physics and sequential Monte Carlo. Diffusion models build a generative Markov chain which converts a simple known distribution (e.g. a Gaussian) into a target (data) distribution using a diffusion process. Rather than use this Markov chain to approximately evaluate a model which has been otherwise defined, one can explicitly define the probabilistic model as the endpoint of the Markov chain. Since each step in the diffusion chain has an analytically evaluable probability, the full chain can also be analytically evaluated.

Diffusion learning

Learning in this framework involves estimating small perturbations to a diffusion process. Estimating small, analytically tractable, perturbations is more tractable than explicitly describing the full distribution with a single, non-analytically-normalizable, potential function. Furthermore, since a diffusion process exists for any smooth target distribution, this method can capture data distributions of arbitrary form.

Mathematics of diffusion models

Let us go back our discussions of the variational autoencoders from last week, see [https:](https://github.com/CompPhysics/AdvancedMachineLearning/blob/main/doc/pub/week15/ipynb/week15.ipynb)

[//github.com/CompPhysics/AdvancedMachineLearning/blob/main/doc/pub/week15/ipynb/week15.ipynb](https://github.com/CompPhysics/AdvancedMachineLearning/blob/main/doc/pub/week15/ipynb/week15.ipynb). As a first attempt at understanding diffusion models, we can think of these as stacked VAEs, or better, recursive VAEs.

Let us try to see why. As an intermediate step, we consider so-called hierarchical VAEs, which can be seen as a generalization of VAEs that include multiple hierarchies of latent spaces.

Note: Many of the derivations and figures here are inspired and borrowed from the excellent exposition of diffusion models by Calvin Luo at <https://arxiv.org/abs/2208.11970>.

Chains of VAEs

Markovian VAEs represent a generative process where we use Markov chain to build a hierarchy of VAEs.

Each transition down the hierarchy is Markovian, where we decode each latent set of variables \mathbf{h}_t in terms of the previous latent variable \mathbf{h}_{t-1} . Intuitively, and visually, this can be seen as simply stacking VAEs on top of each other (see figure next slide).

One can think of such a model as a recursive VAE.

Mathematical representation

Mathematically, we represent the joint distribution and the posterior of a Markovian VAE as

$$p(\mathbf{x}, \mathbf{h}_{1:T}) = p(\mathbf{h}_T) p_{\theta}(\mathbf{x} | \mathbf{h}_1) \prod_{t=2}^T p_{\theta}(\mathbf{h}_{t-1} | \mathbf{h}_t)$$

$$q_{\phi}(\mathbf{h}_{1:T} | \mathbf{x}) = q_{\phi}(\mathbf{h}_1 | \mathbf{x}) \prod_{t=2}^T q_{\phi}(\mathbf{h}_t | \mathbf{h}_{t-1})$$

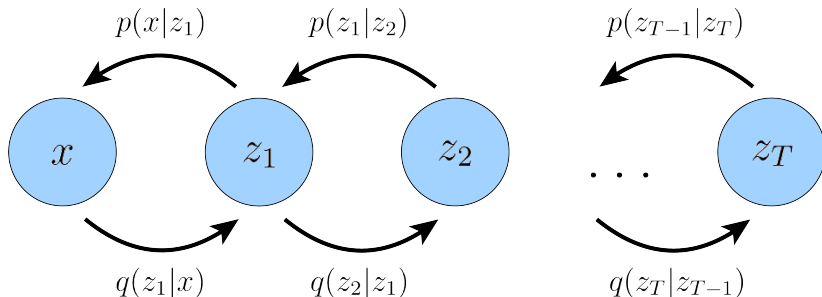
Back to the marginal probability

We can then define the marginal probability we want to optimize as

$$\begin{aligned}\log p(\mathbf{x}) &= \log \int p(\mathbf{x}, \mathbf{h}_{1:T}) d\mathbf{h}_{1:T} \\ &= \log \int \frac{p(\mathbf{x}, \mathbf{h}_{1:T}) q_\phi(\mathbf{h}_{1:T} | \mathbf{x})}{q_\phi(\mathbf{h}_{1:T} | \mathbf{x})} d\mathbf{h}_{1:T} && \text{(Multiply by } 1 = \frac{q_\phi(\mathbf{h}_{1:T} | \mathbf{x})}{q_\phi(\mathbf{h}_{1:T} | \mathbf{x})} \text{)} \\ &= \log \mathbb{E}_{q_\phi(\mathbf{h}_{1:T} | \mathbf{x})} \left[\frac{p(\mathbf{x}, \mathbf{h}_{1:T})}{q_\phi(\mathbf{h}_{1:T} | \mathbf{x})} \right] && \text{(Definition of Expectation)} \\ &\geq \mathbb{E}_{q_\phi(\mathbf{h}_{1:T} | \mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{h}_{1:T})}{q_\phi(\mathbf{h}_{1:T} | \mathbf{x})} \right] && \text{(Discussed last week)}\end{aligned}$$

Diffusion models for hierarchical VAE, from <https://arxiv.org/abs/2208.11970>

A Markovian hierarchical Variational Autoencoder with T hierarchical latents. The generative process is modeled as a Markov chain, where each latent \mathbf{h}_t is generated only from the previous latent \mathbf{h}_{t+1} . Here \mathbf{z} is our latent variable \mathbf{h} .



Equation for the Markovian hierarchical VAE

We obtain then

$$\mathbb{E}_{q_{\phi}(\mathbf{h}_{1:T}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{h}_{1:T})}{q_{\phi}(\mathbf{h}_{1:T}|\mathbf{x})} \right] = \mathbb{E}_{q_{\phi}(\mathbf{h}_{1:T}|\mathbf{x})} \left[\log \frac{p(\mathbf{h}_T) p_{\theta}(\mathbf{x}|\mathbf{h}_1) \prod_{t=2}^T p_{\theta}(\mathbf{h}_t|\mathbf{h}_{t-1})}{q_{\phi}(\mathbf{h}_1|\mathbf{x}) \prod_{t=2}^T q_{\phi}(\mathbf{h}_t|\mathbf{h}_{t-1})} \right]$$

We will modify this equation when we discuss what are normally called Variational Diffusion Models.

Variational Diffusion Models

The easiest way to think of a Variational Diffusion Model (VDM) is as a Markovian Hierarchical Variational Autoencoder with three key restrictions:

1. The latent dimension is exactly equal to the data dimension
2. The structure of the latent encoder at each timestep is not learned; it is pre-defined as a linear Gaussian model. In other words, it is a Gaussian distribution centered around the output of the previous timestep
3. The Gaussian parameters of the latent encoders vary over time in such a way that the distribution of the latent at final timestep T is a standard Gaussian

The VDM posterior is

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

Second assumption

The distribution of each latent variable in the encoder is a Gaussian centered around its previous hierarchical latent. Here then, the structure of the encoder at each timestep t is not learned; it is fixed as a linear Gaussian model, where the mean and standard deviation can be set beforehand as hyperparameters, or learned as parameters.

Parameterizing Gaussian encoder

We parameterize the Gaussian encoder with mean $\mu_t(\mathbf{x}_t) = \sqrt{\alpha_t}\mathbf{x}_{t-1}$, and variance $\Sigma_t(\mathbf{x}_t) = (1 - \alpha_t)\mathbf{I}$, where the form of the coefficients are chosen such that the variance of the latent variables stay at a similar scale; in other words, the encoding process is variance-preserving.

Note that alternate Gaussian parameterizations are allowed, and lead to similar derivations. The main takeaway is that α_t is a (potentially learnable) coefficient that can vary with the hierarchical depth t , for flexibility.

Encoder transitions

Mathematically, the encoder transitions are defined as

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I})$$

Third assumption

From the third assumption, we know that α_t evolves over time according to a fixed or learnable schedule structured such that the distribution of the final latent $p(\mathbf{x}_T)$ is a standard Gaussian. We can then update the joint distribution of a Markovian VAE to write the joint distribution for a VDM as

$$p(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$$

where,

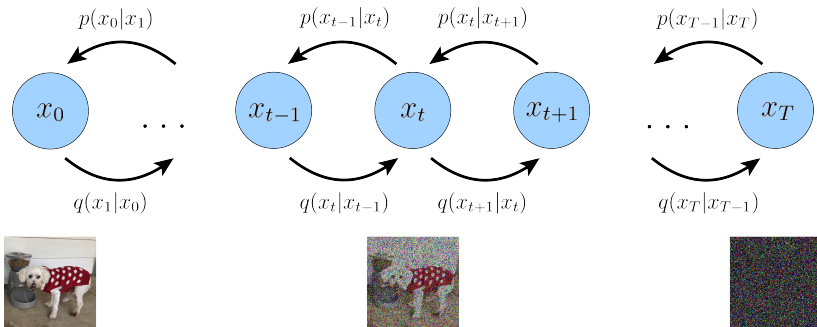
$$p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$$

Noisification

Collectively, what this set of assumptions describes is a steady noisification of an image input over time. We progressively corrupt an image by adding Gaussian noise until eventually it becomes completely identical to pure Gaussian noise. See figure on next slide.

Diffusion models, from

<https://arxiv.org/abs/2208.11970>



Gaussian modeling

Note that our encoder distributions $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ are no longer parameterized by ϕ , as they are completely modeled as Gaussians with defined mean and variance parameters at each timestep. Therefore, in a VDM, we are only interested in learning conditionals $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$, so that we can simulate new data. After optimizing the VDM, the sampling procedure is as simple as sampling Gaussian noise from $p(\mathbf{x}_T)$ and iteratively running the denoising transitions $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ for T steps to generate a novel \mathbf{x}_0 .

Optimizing the variational diffusion model

$$\begin{aligned}\log p(\mathbf{x}) &= \log \int p(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \\&= \log \int \frac{p(\mathbf{x}_{0:T}) q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T} \\&= \log \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\&\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_\theta(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=2}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_T|\mathbf{x}_{T-1}) \prod_{t=1}^{T-1} q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_\theta(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=1}^{T-1} p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_T|\mathbf{x}_{T-1}) \prod_{t=1}^{T-1} q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T) p_\theta(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{\prod_{t=1}^{T-1} p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})}{\prod_{t=1}^{T-1} q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right]\end{aligned}$$

Continues

$$\begin{aligned}\log p(\mathbf{x}) &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)p_\theta(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \prod_{t=1}^{T-1} \frac{p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t+1})} \right] \\&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \prod_{t=1}^{T-1} \frac{p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t+1})} \right] \\&= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \sum_{t=1}^{T-1} \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t+1})} \right] \\&= \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{T-1},\mathbf{x}_T|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \right] + \sum_{t=1}^{T-1} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_{t+1})} \left[\log \frac{p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})}{q(\mathbf{x}_t|\mathbf{x}_{t+1})} \right]\end{aligned}$$

Interpretations

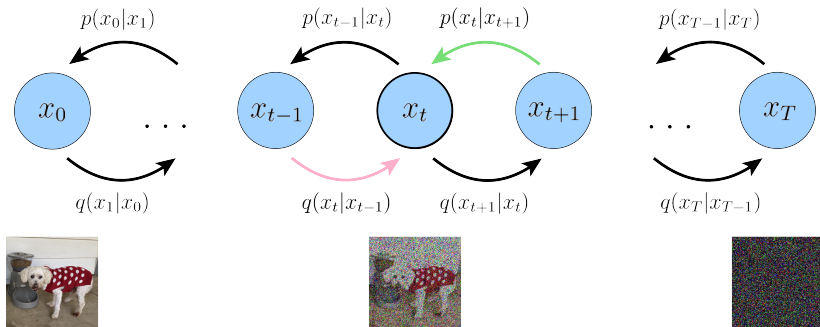
These equations can be interpreted as

- ▶ $\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]$ can be interpreted as a **reconstruction term**, predicting the log probability of the original data sample given the first-step latent. This term also appears in a vanilla VAE, and can be trained similarly.
- ▶ $\mathbb{E}_{q(\mathbf{x}_{T-1}|\mathbf{x}_0)} [D_{KL}(q(\mathbf{x}_T|\mathbf{x}_{T-1})||p(\mathbf{x}_T))]$ is a **prior matching term**; it is minimized when the final latent distribution matches the Gaussian prior. This term requires no optimization, as it has no trainable parameters; furthermore, as we have assumed a large enough T such that the final distribution is Gaussian, this term effectively becomes zero.

The last term

- ▶ $\mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_{t+1} | \mathbf{x}_0)} [D_{KL}(q(\mathbf{x}_t | \mathbf{x}_{t-1}) || p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1}))]$ is a *consistency term*; it endeavors to make the distribution at \mathbf{x}_t consistent, from both forward and backward processes. That is, a denoising step from a noisier image should match the corresponding noising step from a cleaner image, for every intermediate timestep; this is reflected mathematically by the KL Divergence. This term is minimized when we train $p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1})$ to match the Gaussian distribution $q(\mathbf{x}_t | \mathbf{x}_{t-1})$.

Diffusion models, part 2, from <https://arxiv.org/abs/2208.11970>



Optimization cost

The cost of optimizing a VDM is primarily dominated by the third term, since we must optimize over all timesteps t .

Under this derivation, all three terms are computed as expectations, and can therefore be approximated using Monte Carlo estimates.

However, actually optimizing the ELBO using the terms we just derived might be suboptimal; because the consistency term is computed as an expectation over two random variables

$\{\mathbf{x}_{t-1}, \mathbf{x}_{t+1}\}$ for every timestep, the variance of its Monte Carlo estimate could potentially be higher than a term that is estimated using only one random variable per timestep. As it is computed by summing up $T - 1$ consistency terms, the final estimated value may have high variance for large T values.

More details

For more details and implementaions, see Calvin Luo at
<https://arxiv.org/abs/2208.11970>