

Probabilistic Generative Models

Jonas Hübotter

May 7th, 2022

This is intended as an introduction to the use of uncertainty in generative models. We will start without assuming any prior knowledge of generative models and progress quickly to discuss the most fundamental ideas behind some of the most powerful generative models to date.

1 Generative Models

A generative model is a model that learns to imitate some pattern-generating process. That is, we have (high-dimensional) patterns $x \in \mathcal{X}$ that are generated according to a probability distribution p on \mathcal{X} . We want to learn a model p_θ that is indistinguishable from p . In particular, our model should

- (1) avoid missing patterns that are generated by p (this is called *mode collapse*) and
- (2) avoid generating patterns that are not supported by p .

A famous example is the example of generating human faces, where we want to generate facial features proportional to how frequently they occur in real human faces.

1.1 Density Estimation

The problem as stated is above is also called a *density estimation* problem, that is, we want to estimate a density p using a family of densities p_θ parametrized by θ . The canonical method for parameter estimation (also called inference) is *maximum likelihood estimation*, for short. In MLE, we pick θ such that it explains the patterns generated by p best,

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim p} [p_\theta(x)] \quad (1)$$

$$= \arg \max_{\theta} \mathbb{E}_{x \sim p} [\log p_\theta(x)]. \quad (2) \quad \text{using monotonicity of the logarithm}$$

$p_\theta(x)$ is also called the *evidence* of x . So maximum likelihood estimation corresponds to maximizing the evidence of x under the model p_θ when x is sampled according to the true distribution p .

Given that we can differentiate $p_\theta(x)$ with respect to θ , we can easily obtain unbiased gradient estimates, which we can use to find a sample-based approximation of θ^* .¹

¹ Because we can only use finitely many samples x , picking a single point estimate of $\hat{\theta}$ might not be the best choice. Often, it is helpful to encode the uncertainty in our choice of θ^* using a distribution.

Learning a density p directly — for example using standard techniques from deep learning — is hard. This is because our model has to fulfill structural constraints to be a valid density. Moreover, learning the densities directly does not allow us to efficiently sample from the resulting distribution.

Simple parametric models such as *autoregressive models*, which break down the problem of generating x into the problem of sequentially generating components $x(i)$ given previously generated components $x(1), \dots, x(i-1)$, turn out not to be expressive enough for challenging, high-dimensional generative tasks.² We have seen that throwing deep neural networks (DNNs) at the task directly is not very useful, but we do need their expressiveness for modeling p accurately.

1.2 Implicit Models

A useful approach to make DNNs work in a constrained setting, is to encode the constraints implicitly in our model. This idea leads to the family of *implicit models*, where we generate x by first generating a random vector z from a known distribution \mathcal{D} and then use a deterministic DNN f_θ with parameters θ to “lift” z into the set of patterns \mathcal{X} . We can think of this as a transformation of random vectors,

$$x = f_\theta(z), \quad z \sim \mathcal{D}. \quad (3)$$

Often, we sample the *encoding* z according to the standard normal distribution, $\mathcal{D} \doteq \mathcal{N}(\mathbf{0}, I)$.

For general DNNs f , parameter estimation becomes intractable. To perform parameter estimation, we need to recover the density p_θ of our model. This can be done using so-called *pushforward measures* when f is invertible and it is easy to compute $\det J_f$, where J_f is the Jacobian of f . Such f are also called *normalizing flows*.³

1.3 Outline

Normalizing flows restrict the architecture of the DNN f to ensure that recovering the density, and thus inference, remains tractable. An alternative approach is to perform inference approximately. This approach, we will discuss in the next section on *variational autoencoders*.

Finally, we will look at a different and very powerful class of models known as *generative adversarial models*.

2 Variational Autoencoders

One common approach to perform approximate inference is to maximize a lower bound on the evidence instead of maximizing

² Autoregressive models are very useful for time-series prediction.

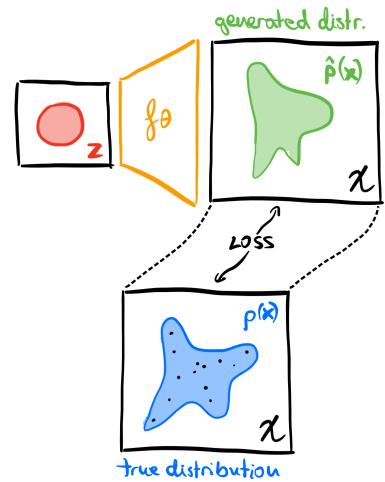


Figure 1: An implicit model “lifts” z from the *code space* (also called *latent space*) to the *pattern space* using a DNN f parametrized by θ .

³ We will not go into more depth on normalizing flows here.

this evidence directly. To obtain a lower bound on the log-evidence $\log p_\theta(x)$, we can relate this quantity back to the space of encodings:

$$\log p_\theta(x) = \log \int p_\theta(x | z) p(z) dz.$$

Intuitively, the probability of x is the probability of an encoding z being “lifted” to x times the probability of the encoding z to be sampled in the first place.⁴ We can rewrite this to,

$$\begin{aligned} &= \log \int p_\theta(x | z) \frac{p(z)}{q(z)} q(z) dz \\ &= \log \mathbb{E}_{z \sim q} \left[p_\theta(x | z) \frac{p(z)}{q(z)} \right], \end{aligned}$$

where q is any density on the code space. Using Jensen’s inequality and that \log is concave, we obtain,

$$\geq \mathbb{E}_{z \sim q} [\log p_\theta(x | z)] - \mathbb{E}_{z \sim q} \left[\log \frac{q(z)}{p(z)} \right]. \quad (4)$$

This lower bound on the log-evidence is also known as the *evidence lower bound* (or ELBO for short).

The regularization term $\mathbb{E}_{z \sim q} \left[\log \frac{q(z)}{p(z)} \right]$ is also called (*reverse*) *Kullback-Leibler divergence*, $\text{KL}(q \| p)$, and measures how well q approximates the true distribution of encodings p . Maximizing ELBO therefore corresponds to maximizing the likelihood of x given samples of encodings z , while z are sampled according to a distribution q , which is close to p . When q is from a class of distributions parametrized by ψ , then q_ψ is also called a *variational distribution*.

⁴ The true distribution over the latent space $p(z)$ is not to be confused with the true distribution over the pattern space $p(x)$.

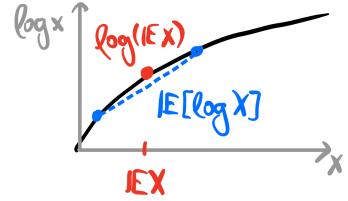


Figure 2: An illustration of Jensen’s inequality.

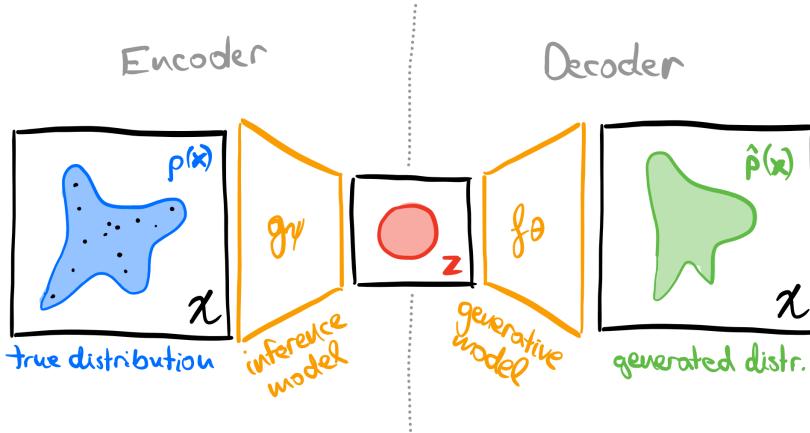


Figure 3: A variational autoencoder learns an inference model and generative model simultaneously. The learned inference model is used to produce encodings from which the generative model should produce patterns.

The key idea of variational autoencoders is to simultaneously learn a generative model p_θ and inference model q_ψ . So given a pattern x

sampled according to the true distribution of patterns p , the objective to be maximized is,

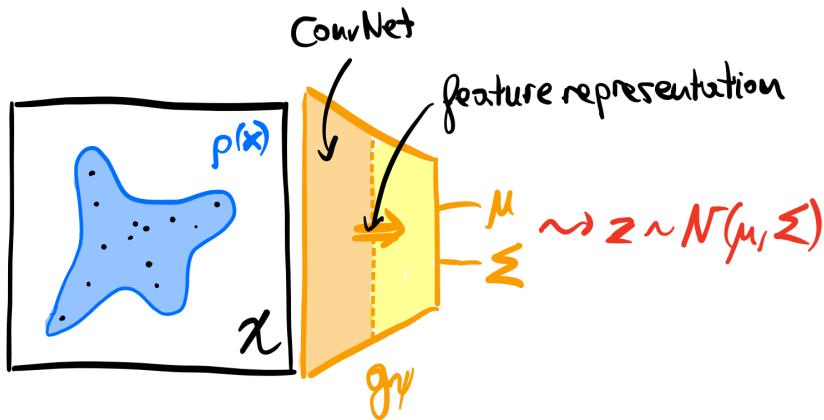
$$\log p_{\theta}(x) \gtrsim \mathbb{E}_{z \sim q_{\psi}(\cdot | x)} [\log p_{\theta}(x | z)] - \text{KL}(q_{\psi}(\cdot | x) \| p) \doteq \ell(\theta, \psi; x). \quad (5)$$

The KL-divergence can often be expressed analytically, and we can use a sample-based approximation for the expectation.

2.1 The Inference Model

You may (rightfully) object that we have simply extended our task of finding a generative model p_{θ} by finding yet another generative model g_{ψ} . However, the latter generative model is from a high- to a low-dimensional space, namely from the pattern space to the code space. It turns out that this task is much simpler than the task that we started with.

Typically, the inference model is chosen to be a non-isotropic statistical transformation of the outputs of a DNN. More concretely, we use a DNN g_{ψ} to first transform the high-dimensional pattern x to a low-dimensional feature representation and then to outputs μ_{ψ} and Σ_{ψ} .⁵ We then choose q_{ψ} to be a simple distribution with mean μ_{ψ} and covariance Σ_{ψ} (i.e., $\mathcal{N}(\mu_{\psi}, \Sigma_{\psi})$) that is easy to sample from.



We call such a model *non-isotropic* because the statistics μ_{ψ} and Σ_{ψ} depend on the pattern x .

2.2 Inference

Now, to perform inference, we need to be able to obtain gradients with respect to θ and ψ of $\ell(\theta, \psi; x)$. As we mentioned, the KL-divergence can usually be expressed analytically. This is also true

⁵ For example, by using a convolutional neural network to obtain the feature representation and a fully connected neural network to learn mappings to μ_{ψ} and Σ_{ψ} .

Figure 4: Illustration of a typical inference model. First, a pattern is transformed into a feature representation, which is then used to obtain statistics μ and Σ that parametrize a distribution over encodings.

for its gradient. The gradient of $p_\theta(x | z)$ can be computed using backpropagation. The difficulty lies in finding the gradient of the expectation $\mathbb{E}_{z \sim q_\psi(\cdot | x)}[\log p_\theta(x | z)]$ with respect to ψ . Because the expectation integrates over a measure that depends on the parameters ψ , we cannot simply draw the gradient into the expectation using linearity.

Here, we can use the so-called *reparameterization trick* by writing $z \sim q_\psi(\cdot, x) = \mathcal{N}(\mu_\psi, \Sigma_\psi)$ as,

$$z = \mu_\psi + \Sigma_\psi^{1/2} \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, I), \quad (6)$$

where μ_ψ and Σ_ψ are the outputs of the DNN g_ψ . We obtain,

$$\mathbb{E}_{z \sim q_\psi(\cdot | x)}[\log p_\theta(x | z)] = \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \left[\log p_\theta(x | \mu_\psi + \Sigma_\psi^{1/2} \epsilon) \right]. \quad (7)$$

Of this expectation, we can now easily obtain unbiased gradient estimates, which we can use to maximize the ELBO using stochastic gradient descent. For more details on the reparameterization trick for Gaussians or other distributions and the computation of gradients of the ELBO, refer to section 5.2.7 of my notes on probabilistic artificial intelligence.



Figure 5: An example of faces reconstructed by a variational autoencoder. Taken from [3].

3 Generative Adversarial Models

We will now introduce a different class of models that are based on a very simple and beautiful idea, and which are powerful enough to yield models that produce faces that are visually indistinguishable from real faces.

Recall that our original goal was to find a distribution p_θ that is indistinguishable from the distribution p . Our initial instinct was to use density estimation to find p_θ , but we have seen that we either have to limit the expressiveness of our model to retain tractability of the inference problem or we have to perform inference approximately. It is therefore useful to think whether we can characterize our problem in a different way.

The notion of indistinguishability leads to a very natural characterization: our goal is to make it virtually impossible for an agent to distinguish between the patterns produced by p and the patterns produced by p_θ . From the perspective of such an agent, this is a simple binary classification problem. Our original task can be interpreted as maximizing the loss of this agent, as a large loss corresponds to the scenario where the patterns produced by p_θ cannot be distinguished from the patterns produced by p .

We write $Y \doteq \mathbb{1}\{\text{pattern is "natural"}\}$. Assuming that the classification problem is balanced,⁶ it can be characterized using the joint density,

$$p(\mathbf{x}, y; \boldsymbol{\theta}) = \underbrace{\mathbb{P}[Y = y]}_{= \frac{1}{2}} \begin{cases} p(\mathbf{x}) & y = 1 \\ p_\theta(\mathbf{x}) & y = 0. \end{cases} \quad (8)$$

A classifier $\pi : \mathbb{R}^n \rightarrow [0, 1]$, $\pi(\mathbf{x}) \approx \mathbb{P}[Y = 1 | \mathbf{x}]$ is also called a *discriminator*. Intuitively, the discriminator provides an error signal to the generative model, hinting at promising regions in the pattern space.

The canonical approach is to model the loss of a discriminator π using *binary cross-entropy* (also called a *logistic loss*),

$$\ell(\boldsymbol{\theta}, \pi) \doteq \mathbb{E}_{\mathbf{x}, y \sim p(\cdot, \boldsymbol{\theta})} [y \log \pi(\mathbf{x}) + (1 - y) \log(1 - \pi(\mathbf{x}))] \quad (9)$$

$$= \frac{1}{2} \int p(\mathbf{x}) \log \pi(\mathbf{x}) d\mathbf{x} + \frac{1}{2} \int p_\theta(\mathbf{x}) \log(1 - \pi(\mathbf{x})) d\mathbf{x}. \quad (10)$$

Our optimization problem is then given as,

$$\min_{\boldsymbol{\theta}} \max_{\pi} \ell(\boldsymbol{\theta}, \pi). \quad (11)$$

In words, we want to make the classification task as difficult as possible for the best-possible classifier π by tuning $\boldsymbol{\theta}$.

The optimal discriminator is a Bayes-optimal classifier π^* ,

$$\pi^*(\mathbf{x}) = \mathbb{P}[Y = 1 | \mathbf{x}] = \frac{p(\mathbf{x})}{p(\mathbf{x}) + p_\theta(\mathbf{x})}. \quad (12)$$

In practice, we do not have access to this discriminator, as it depends on the true distribution p . Nevertheless, we can compute the expected loss of this classifier to confirm our intuition that when this quantity is large, p_θ must be “almost indistinguishable” from p . For the Bayes-optimal discriminator π^* , we obtain,

$$\begin{aligned} \ell(\boldsymbol{\theta}, \pi^*) &= \mathbb{E}_{\mathbf{x}, y \sim p(\cdot, \boldsymbol{\theta})} [y \log \pi^*(\mathbf{x}) + (1 - y) \log(1 - \pi^*(\mathbf{x}))] \\ &= \frac{1}{2} \int p(\mathbf{x}) [\log p(\mathbf{x}) - \log(p(\mathbf{x}) + p_\theta(\mathbf{x}))] d\mathbf{x} \\ &\quad + \frac{1}{2} \int p_\theta(\mathbf{x}) [\log p_\theta(\mathbf{x}) - \log(p(\mathbf{x}) + p_\theta(\mathbf{x}))] d\mathbf{x} \end{aligned}$$

⁶ That is, we present the agent with labels such that $\mathbb{P}[Y = 1] = \mathbb{P}[Y = 0] = 1/2$.

$$\begin{aligned}
&= \frac{1}{2} \int p(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x} + \frac{1}{2} \int p_{\theta}(\mathbf{x}) \log p_{\theta}(\mathbf{x}) d\mathbf{x} \\
&\quad - \int \frac{p(\mathbf{x}) + p_{\theta}(\mathbf{x})}{2} \log \left(\frac{p(\mathbf{x}) + p_{\theta}(\mathbf{x})}{2} \right) d\mathbf{x} - \log 2 \\
&= \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p} [\log p(\mathbf{x})] + \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\theta}} [\log p_{\theta}(\mathbf{x})] \\
&\quad - \mathbb{E}_{\mathbf{x} \sim \frac{1}{2}p + \frac{1}{2}p_{\theta}} \left[\log \left(\frac{1}{2}p(\mathbf{x}) + \frac{1}{2}p_{\theta}(\mathbf{x}) \right) \right] - \log 2 \\
&= -\frac{1}{2}H(p) - \frac{1}{2}H(p_{\theta}) + H\left(\frac{1}{2}p + \frac{1}{2}p_{\theta}\right) - \log 2, \tag{13}
\end{aligned}$$

where $H(p) \doteq \mathbb{E}_{\mathbf{x} \sim p} [-\log p(\mathbf{x})]$ is the entropy of distribution p and $\frac{1}{2}p + \frac{1}{2}p_{\theta}$ is a convex combination of the distributions p and p_{θ} . The sum of entropies is also known as the *Jensen-Shannon divergence*,

$$\text{JS}(p, q) = \frac{1}{2} \text{KL}\left(p \parallel \frac{1}{2}p + \frac{1}{2}q\right) + \frac{1}{2} \text{KL}\left(q \parallel \frac{1}{2}p + \frac{1}{2}q\right) \tag{14}$$

$$= -\frac{1}{2}H(p) - \frac{1}{2}H(q) + H\left(\frac{1}{2}p + \frac{1}{2}q\right). \tag{15}$$

Consider the case where p and q are the same distribution, then the JS-divergence is zero. The described optimization problem minimizes the JS-divergence, and hence, the “distance” between p and p_{θ} .

As mentioned, we cannot compute π^* directly, nor can we compute the Jensen-Shannon divergence. Nevertheless, the above derivation motivates that replacing $\arg \max_{\pi} \ell(\theta, \pi)$ by a learned discriminator that is close to the true optimum π^* yields a loss that can be used to estimate a good generator. Let us consider a class of discriminators π_{ϕ} that is parametrized by ϕ . Assuming this class is expressive enough to cover π^* , we have,

$$\ell(\theta, \pi^*) = \max_{\phi} \ell(\theta, \phi), \tag{16}$$

yielding the saddle point problem,

$$\theta^* = \arg \min_{\theta} \max_{\phi} \ell(\theta, \phi), \quad \phi^* = \arg \max_{\phi} \ell(\theta^*, \phi). \tag{17}$$

Under conditions of Nash’s existence theorem for Nash equilibria, the above is equivalent to,

$$\theta^* = \arg \min_{\theta} \ell(\theta, \phi^*), \quad \phi^* = \arg \max_{\phi} \ell(\theta^*, \phi). \tag{18}$$

This formulation of generative models is known as *Generative Adversarial Networks* (GANs). We can think of this problem as a two player game with the generator trying to choose θ so as to minimize the loss and the adversary (i.e., the discriminator) trying to choose ϕ so as to maximize the loss.



Figure 6: An example of faces constructed by a generative adversarial model, StyleGAN. Taken from [2].

References

- [1] Thomas Hofmann. Computational intelligence lab, 2021.
- [2] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019.
- [3] Rafael S Toledo and Eric A Antonelo. Face reconstruction with variational autoencoder and face masks. *arXiv preprint arXiv:2112.02139*, 2021.