

Tutorial on Autoencoders and VAEs

Matteo Lugli

January 3, 2024

1 Introduction

This topic will be quite dense in terms of statistical concepts. I will try to write a summary that is as simple as possible, and I will also link the main resources that I used at the end of the document (see 8). If you are reading this, you should be familiar with deep learning and Bayesian inference.

Sections 2 and 3 will be a brief introduction to autoencoders, a really important deep learning model that is used to compress data. In section 4 I will introduce generative modeling related to the concept of latent variables. In section 5 a framework to train variational autoencoders is provided, together with the explanation. In chapter 6 I explain how you can use a trained VAE to perform different tasks. In chapter 7 I give a brief introduction to the reparametrization trick, which is crucial to make these kind of models trainable with gradient descent.

2 Autoencoders

Autoencoders exploit the ability of neural networks to represent information in a contracted form. The basic idea is simple: given an input x , we forward it through a neural network that has a "bottleneck" layer called h with a lower dimension than the input, and which gives in output a vector with the same dimension as the original x . We call the output \bar{x} . If we train the network in a way that minimizes the distance between x and \bar{x} , we obtain a model that learns a "compressed" representation of the input data. Hopefully, the contracted representation (also called **code**) can be interpreted as a combination of the most important features that are observed in the training data (meaning the ones that vary the most). This approach is similar to PCA, but works also with samples that are distributed on a non-linear manifold. Think of the *two-moon dataset* (google it if you are not familiar with it): PCA will not work on that dataset because samples are not distributed on a linear manifold (a line in 2D) but on a more complicated one. In other words, PCA works on datasets where dependencies between points are quite simple. Our goal is to find another method that can summarize also more complex feature dependencies.

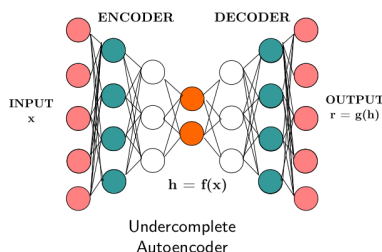


Figure 1: Autoencoder structure

In general, if we want a good lower dimensional representation of our data, we aim at:

- preserving as much information on the data as possible;
- the new representation has to be simpler and more accessible;

The key idea behind the training procedure is to minimize a loss function that gives

a penalty if $g(f(x))$ is dissimilar to x :

$$L(x, g(f(x))) \quad (1)$$

Even if the autoencoder learns successfully to minimize 1, some problems may arise:

- Learning the **identity function**: if the code layer is given too much capacity, the network will learn to copy and paste the original representation of x .

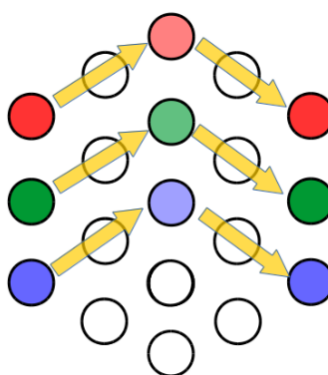


Figure 2: Copy-paste problem

- Learning a **mapping function**: if the encoder and decoder are given to deep and big layers, even with a hidden layer of 1 unit the network can learn to use a mapping procedure:

$$h = f(x^{(i)}) = i \quad (2)$$

\mathbf{x}	$\mathbf{x}^{(1)}$	$\mathbf{x}^{(2)}$...	$\mathbf{x}^{(i)}$...	$\mathbf{x}^{(N)}$
\mathbf{h}	1	2	...	i	...	N

Figure 3: Mapping problem

There are many ways to address these issues: they all involve a modification of the loss function needed to regularize the training phase, to avoid both overfitting and underfitting. In the next chapter I summarize the main regularization methods.

3 Regularized Autoencoders

The baseline loss function for autoencoders is the following:

$$L = D(x, g(f(x))) \quad (3)$$

with sparse autoencoders we also add some regularization terms according to what kind of regularization we want to apply.

3.1 Sparse Autoencoders

With sparse autoencoders we encourage the neurons of the layer \mathbf{h} to be mostly inactive. By using less neurons, we imply that the new features of the encoded sample will be more "orthogonal". In equation 4 $\bar{\rho}_j$ indicates the average activation of hidden unit j (of the hidden layer \mathbf{h}) calculated over the whole training set. In sparse autoencoders we want this value to be as close as possible to a sparsity parameter ρ , which is close to zero.

$$\bar{\rho}_j = \frac{1}{m} \sum_{i=1}^m a_j^{(h)}(x^{(i)}) \quad (4)$$

The loss function of a sparse autoencoder includes the sparsity penalty as well:

$$L = D(x, g(f(x))) + \lambda \sum_{j=1}^{|h|} KL(\rho || \bar{\rho}_j) \quad (5)$$

where $KL(\rho || \bar{\rho}_j)$ is the KL divergence between a bernoulli distribution with mean ρ and a bernoulli distribution with mean $\bar{\rho}_j$, which has a closed form solution. In the summation, j iterates over all of the hidden units in the layer h . λ controls the importance of the regularization term compared to the reconstruction term.

3.2 Denoising Autoencoders

The key idea behind denoising autoencoders is to train the network to reconstruct the original input x from a noisy version \hat{x} of the input itself.

$$L = D(x, g(f(\hat{x}))) \quad (6)$$

where $\hat{x} \sim \mathcal{N}(x, \Sigma I)$. This type of autoencoder doesn't run the risk of learning the identity function by construction, as otherwise, the loss function would be high. In other words, the autoencoder learns a vector field that brings back noisy samples on the underlying manifold on which the training data lies on. This means that denoising autoencoders implicitly learn the structure of that manifold!

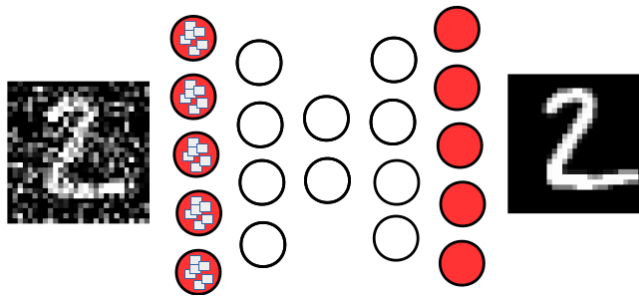


Figure 4: Denoising autoencoders

3.3 Contractive Autoencoders

With contractive autoencoders the main objective is to make the hidden representation as insensitive as possible to small changes of the input. Also in this case learning the copy function is not possible, because it would lead to big changes in the hidden representation with a small change in the input. In other words, we want the hidden representation to be sensitive only to the most important features of the training set. This kind of autoencoder somehow "simplifies" the structure of the manifold.

$$L = D(x, g(f(x))) + \lambda \left\| \frac{\partial f(x)}{\partial x} \right\|^2 \quad (7)$$

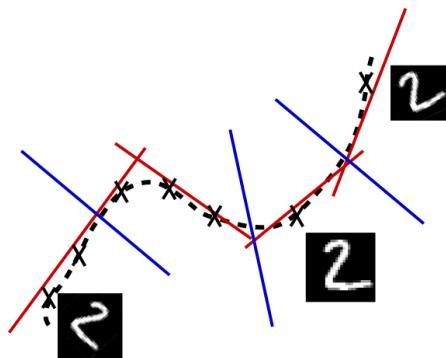


Figure 5: Contractive autoencoders

4 Latent variables models

The classic autoencoder structure used in (for example) denoising autoencoders is useful to implicitly learn the structure of the training data. The problem is that they do not provide a *explicit* representation of that distribution. If we want to generate new data, we need a way to work with a distribution that it's known and from which we can sample!

VAEs are also called *Latent Variable Models*, because they make the assumption that the data can be represented through a combination of **latent variables**, denotes as \mathbf{z} . If we have dataset made of images, latent variables might be the orientation of the subject, the color of the landscape, ecc... thanks to this assumption, instead of directly sampling from $p(x)$, we sample a set of latent variables from $p(z)$, then we use the produced z to sample from another distribution, $p(x|z)$.

The final objective is to maximize the likelihood of our data under the entire generative process:

$$p(x) = \int p(x|z)p(z)dz \quad (8)$$

this framework is usually called **maximum likelihood**. The problem is that computing that integral is too hard: theoretically we would need to take a infinite number of z , which is obviously impossible. What we might do is using some sort of Monte Carlo method, where we sample many z and only use those. With real world data this method is really uneffective and requires to many samples.

In practice, $p(x|z)$ will be zero for most of the combinations of the latent variables (among all the possible combinations), and will be useless when computing the integral. One of the core ideas behind variational autoencoders is trying to understand how the latent variables are distributed in our dataset, so that we can sample the z from that distribution instead. This way we will compute $p(x)$ with the z that are more likely to have produced x . Let's write the **posterior** that we want to compute, that describes how the latent variables are distributed given the dataset:

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)} \quad (9)$$

as we can see this is still intractable because we find $p(x)$ at the denominator. In **variational inference** we try to build a *surrogate* distribution $q(z|x)$ that approximates the *true posterior* $p(z|x)$ by using a much simpler structure (a gaussian). Ideally, we want the surrogate to be as similar as possible to the true posterior, so we want that the KL divergence to be small¹. Let's write the formula of the distance:

$$D_{KL}(q(z|x), p(z|x)) = E_q[\log q(z|x), \log p(z|x)] \quad (10)$$

$$= E_q[\log q(z|x), \log \frac{p(x|z)p(z)}{p(x)}] \quad (11)$$

$$= E_q[\log q(z|x)] - E_q[\log p(x, z)] + E_q[p(x)] \quad (12)$$

$$= E_q[\log q(z|x)] - E_q[\log p(x, z)] + p(x) \quad (13)$$

From 12 to 13 we removed the expectation on $\log p(x)$ because it is independent from z ; let's bring $\log p(x)$ to the left hand side and rewrite the rest:

$$\log p(x) = \underbrace{E_{z \sim q}[\log p(x, z)] - E_{z \sim q}[\log q(z|x)]}_{\text{ELBO}} + D_{KL}(q(z|x), p(z|x)) \quad (14)$$

In equation 14 we grouped the first 2 terms of the right hand side under the name of *ELBO*. This term is a **lower bound** of the left hand side, which is what we wanted to approximate at the beginning. By rewriting it we will see that it is a term that we can compute and work with instead of $p(x)$. For now, notice that if we find a way to maximize it, at the same time we will minimize the KL divergence between the proxy and the posterior: the left hand side is fixed and does not depend on any parameter, and the KL term is always positive because of Jensen inequalities.

¹Intuitively, it measures the distance between the two distributions

Now that we know that by maximizing the ELBO we also make $q(z|x)$ similar to $p(z|x)$, let's re-write the formula to analyze it better:

$$ELBO = E[\log(p(x|z)p(z))] - E[\log q(z|x)] \quad (15)$$

$$= E[\log p(x|z)] + E[\log p(z)] - E[\log q(z|x)] \quad (16)$$

$$= E[\log p(x|z)] - D_{KL}(\log q(z|x), p(z)) \quad (17)$$

$$ELBO = \underbrace{E_{z \sim q(z|x)}[\log p(x|z)]}_{\text{Reconstruction term}} - \underbrace{D_{KL}(\log q(z|x), p(z))}_{\text{Regularization term}} \quad (18)$$

- $p(z)$ is our prior distribution over the latent variables. It is set to be a gaussian: $\mathcal{N}(0, I)$;
- $q(z|x)$ is the proxy of the posterior, modeled as a normal distribution. The parameters of this distribution are the output of the **encoder network**, that taken a x value outputs the mean $\mu(x)$ and the s.t.d $\sigma^2(x)$ of the distribution instead of the plain encoding of the sample.
- $p(x|z)$ is the distribution that given a set of latent variables z describes the likelihood of each x , modeled as a normal distribution as well. The parameters of this distribution are the output of the **decoder network**.

5 Training procedure

In Figure 6 I try to explain visually the forward pass that is used during the training of the variational autoencoder.

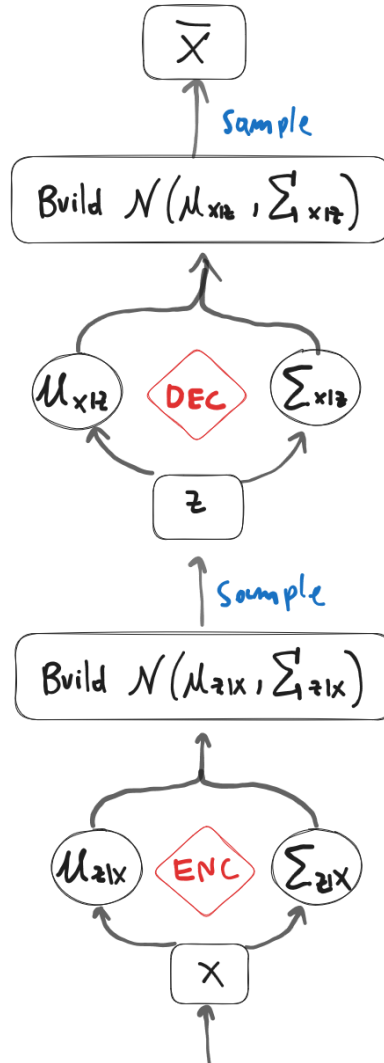


Figure 6: Training forward pass

Ideally we want \bar{x} to be similar to the original sample x , so we can minimize the MSE between the input and the reconstruction: $\|x - \bar{x}\|^2$ by using gradient descent. All of the operations except for the sampling are derivable, so we can perform standard backpropagation to train the network. To make also the sampling derivable, we use

a little trick called **reparametrization trick** (see section 7). We have now seen how we can maximize the reconstruction term by minimizing the MSE between the reconstruction and the real input. To maximize the ELBO, we also want to minimize the regularization term, that represents how much the proxy posterior diverges from the prior $p(z)$. Considering that we model both of them as normal distributions as explained in 4, the regularization term has a closed form solution, and the KL divergence between two gaussians can be easily computed.

We use the term "inference" because when we forward data through the encoder, we want to "infer" the underlying distribution of the latent variables. On the other hand, sometimes the term "generation network" is used when referring to the decoder, because we can generate new data by sampling from $p(z)$, do a forward pass and sample again from $q(x|z)$.

6 How to generate new data

If we want to generate variants of a sample, we can just do an entire forward pass and keep \bar{x} . If we want to generate data from scratch, we can sample $z \sim \mathcal{N}(0, I)$ and forward the latent variable into the decoder, that will output a new sample similar to the ones in the training dataset. One thing we can also do is calculate what is the probability of the decoder producing an object x . To do that we would need to compute $p(x)$, which we know is intractable. What we can do instead is just compute the ELBO with a forward pass, which is a **lower bound** of $p(x)$, so it is a valid approximation of what we want.

7 Reparametrization trick

As we said, we want each operation that we do for the forward pass to be derivable, so that we can backprop during training. Unfortunately, doing $z \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$ is not a derivable operation. With little effort we can make it derivable: we just sample ϵ from $\mathcal{N}(0, I)$ and do $z = \mu_{z|x} + \epsilon \Sigma_{z|x}$. This way we provide a way for the gradient to backpropagate also through the sampling operator.

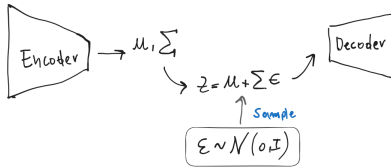


Figure 7: Reparametrization trick

8 References

- [1] Doersch, Carl. "Tutorial on variational autoencoders." arXiv preprint arXiv:1606.05908 (2016).
- [2] Slides from "Machine Learning and Deep Learning" course of AI Engineering, Unimore
- [3] Stanford lecture on generative models, <https://www.youtube.com/watch?v=5WoItGTWV54>
- [4] Introduction to variational inference, <https://www.youtube.com/watch?v=HxQ94L8n0vU>
- [5] Andrew Ng notes, https://web.stanford.edu/class/cs294a/sparseAutoencoder_2011new.pdf