

Advanced machine learning and data analysis for the physical sciences

Morten Hjorth-Jensen¹

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

April 9, 2024

© 1999-2025, Morten Hjorth-Jensen. Released under CC Attribution-NonCommercial 4.0 license

Plans for the week April 8-12, 2024

Generative methods, energy models and Boltzmann machines

1. Summary of discussions on Restricted Boltzmann machines, reminder from last week
2. Introduction to Variational Autoencoders (VAEs)

Reading recommendations

1. Boltzmann machines: Goodfellow et al chapters 18.1-18.2, 20.1-20.7; To create Boltzmann machine using Keras, see Babcock and Bali chapter 4, see https://github.com/PacktPublishing/Hands-On-Generative-AI-with-Python-and-TensorFlow-2/blob/master/Chapter_4/models/rbm.py
2. More on Boltzmann machines: see also Foster, chapter 7 on energy-based models at https://github.com/davidADSP/Generative_Deep_Learning_2nd_Edition/tree/main/notebooks/07_ebm/01_ebm
3. VAEs: Goodfellow et al, for VAEs see sections 20.10-20.11

Essential elements of generative models

The aim of generative methods is to train a probability distribution p . The methods we will focus on are:

1. Energy based models, with the family of Boltzmann distributions as a typical example
2. Variational autoencoders, based on our discussions on autoencoders
3. Generative adversarial networks (GANs) and
4. Diffusion models

Energy models, reminders from last two weeks

During the last two weeks we defined a domain \mathbf{X} of stochastic variables $\mathbf{X} = \{x_0, x_1, \dots, x_{n-1}\}$ with a pertinent probability distribution

$$p(\mathbf{X}) = \prod_{x_i \in \mathbf{X}} p(x_i),$$

where we have assumed that the random variables x_i are all independent and identically distributed (iid).

We will now assume that we can defined this function in terms of optimization parameters Θ , which could be the biases and weights of deep network, and a set of hidden variables we also assume to be random variables which also are iid. The domain of these variables is $\mathbf{H} = \{h_0, h_1, \dots, h_{m-1}\}$.

Probability model

We define a probability

$$p(x_i, h_j; \Theta) = \frac{f(x_i, h_j; \Theta)}{Z(\Theta)},$$

where $f(x_i, h_j; \Theta)$ is a function which we assume is larger or equal than zero and obeys all properties required for a probability distribution and $Z(\Theta)$ is a normalization constant. Inspired by statistical mechanics, we call it often for the partition function. It is defined as (assuming that we have discrete probability distributions)

$$Z(\Theta) = \sum_{x_i \in \mathbf{X}} \sum_{h_j \in \mathbf{H}} f(x_i, h_j; \Theta).$$

Marginal and conditional probabilities

We can in turn define the marginal probabilities

$$p(x_i; \Theta) = \frac{\sum_{h_j \in \mathbf{H}} f(x_i, h_j; \Theta)}{Z(\Theta)},$$

and

$$p(h_j; \Theta) = \frac{\sum_{x_i \in \mathbf{X}} f(x_i, h_j; \Theta)}{Z(\Theta)}.$$

Change of notation

Note the change to a vector notation. A variable like \mathbf{x} represents now a specific **configuration**. We can generate an infinity of such configurations. The final partition function is then the sum over all such possible configurations, that is

$$Z(\Theta) = \sum_{x_i \in \mathbf{X}} \sum_{h_j \in \mathbf{H}} f(x_i, h_j; \Theta),$$

changes to

$$Z(\Theta) = \sum_{\mathbf{x}} \sum_{\mathbf{h}} f(\mathbf{x}, \mathbf{h}; \Theta).$$

If we have a binary set of variable x_i and h_j and M values of x_i and N values of h_j we have in total 2^M and 2^N possible \mathbf{x} and \mathbf{h} configurations, respectively.

We see that even for the modest binary case, we can easily approach a number of configuration which is not possible to deal with.

Optimization problem

At the end, we are not interested in the probabilities of the hidden variables. The probability we thus want to optimize is

$$p(\mathbf{X}; \Theta) = \prod_{x_i \in \mathbf{X}} p(x_i; \Theta) = \prod_{x_i \in \mathbf{X}} \left(\frac{\sum_{h_j \in \mathbf{H}} f(x_i, h_j; \Theta)}{Z(\Theta)} \right),$$

which we rewrite as

$$p(\mathbf{X}; \Theta) = \frac{1}{Z(\Theta)} \prod_{x_i \in \mathbf{X}} \left(\sum_{h_j \in \mathbf{H}} f(x_i, h_j; \Theta) \right).$$

Further simplifications

We simplify further by rewriting it as

$$p(\mathbf{X}; \Theta) = \frac{1}{Z(\Theta)} \prod_{x_i \in \mathbf{X}} f(x_i; \Theta),$$

where we used $p(x_i; \Theta) = \sum_{h_j \in \mathbf{H}} f(x_i, h_j; \Theta)$. The optimization problem is then

$$\arg \max_{\Theta \in \mathbb{R}^p} p(\mathbf{X}; \Theta).$$

Optimizing the logarithm instead

Computing the derivatives with respect to the parameters Θ is easier (and equivalent) with taking the logarithm of the probability. We will thus optimize

$$\arg \max_{\Theta \in \mathbb{R}^p} \log p(\mathbf{X}; \Theta),$$

which leads to

$$\nabla_{\Theta} \log p(\mathbf{X}; \Theta) = 0.$$

Expression for the gradients

This leads to the following equation

$$\nabla_{\Theta} \log p(\mathbf{X}; \Theta) = \nabla_{\Theta} \left(\sum_{x_i \in \mathbf{X}} \log f(x_i; \Theta) \right) - \nabla_{\Theta} \log Z(\Theta) = 0.$$

The first term is called the positive phase and we assume that we have a model for the function f from which we can sample values. Below we will develop an explicit model for this. The second term is called the negative phase and is the one which leads to more difficulties.

The derivative of the partition function

The partition function, defined above as

$$Z(\Theta) = \sum_{x_i \in \mathbf{X}} \sum_{h_j \in \mathbf{H}} f(x_i, h_j; \Theta),$$

is in general the most problematic term. In principle both x and h can span large degrees of freedom, if not even infinitely many ones, and computing the partition function itself is often not desirable or even feasible. The above derivative of the partition function can however be written in terms of an expectation value which is in turn evaluated using Monte Carlo sampling and the theory of Markov chains, popularly shortened to MCMC (or just MC²).

Explicit expression for the derivative

We can rewrite

$$\nabla_{\Theta} \log Z(\Theta) = \frac{\nabla_{\Theta} Z(\Theta)}{Z(\Theta)},$$

which reads in more detail

$$\nabla_{\Theta} \log Z(\Theta) = \frac{\nabla_{\Theta} \sum_{x_i \in \mathbf{x}} f(x_i; \Theta)}{Z(\Theta)}.$$

We can rewrite the function f (we have assumed that is larger or equal than zero) as $f = \exp \log f$. We can then rewrite the last equation as

$$\nabla_{\Theta} \log Z(\Theta) = \frac{\sum_{x_i \in \mathbf{x}} \nabla_{\Theta} \exp \log f(x_i; \Theta)}{Z(\Theta)}.$$

Final expression

Taking the derivative gives us

$$\nabla_{\Theta} \log Z(\Theta) = \frac{\sum_{x_i \in \mathbf{X}} f(x_i; \Theta) \nabla_{\Theta} \log f(x_i; \Theta)}{Z(\Theta)},$$

which is the expectation value of $\log f$

$$\nabla_{\Theta} \log Z(\Theta) = \sum_{x_i \in \mathbf{X}} p(x_i; \Theta) \nabla_{\Theta} \log f(x_i; \Theta),$$

that is

$$\nabla_{\Theta} \log Z(\Theta) = \mathbb{E}(\log f(x_i; \Theta)).$$

This quantity is evaluated using Monte Carlo sampling, with Gibbs sampling as the standard sampling rule. Before we discuss the explicit algorithms, we need to remind ourselves about Markov chains and sampling rules like the Metropolis-Hastings algorithm and Gibbs sampling.

Positive and negative phases

As discussed earlier, the data-dependent term in the gradient is known as the positive phase of the gradient, while the model-dependent term is known as the negative phase of the gradient. The aim of the training is to lower the energy of configurations that are near observed data points (increasing their probability), and raising the energy of configurations that are far from observed data points (decreasing their probability).

Gradient examples

The gradient of the negative log-likelihood cost function of a Binary-Binary RBM is then

$$\frac{\partial \mathcal{C}(w_{ij}, a_i, b_j)}{\partial w_{ij}} = \langle x_i h_j \rangle_{data} - \langle x_i h_j \rangle_{model}$$

$$\frac{\partial \mathcal{C}(w_{ij}, a_i, b_j)}{\partial a_{ij}} = \langle x_i \rangle_{data} - \langle x_i \rangle_{model}$$

$$\frac{\partial \mathcal{C}(w_{ij}, a_i, b_j)}{\partial b_{ij}} = \langle h_i \rangle_{data} - \langle h_i \rangle_{model}.$$

To get the expectation values with respect to the *data*, we set the visible units to each of the observed samples in the training data, then update the hidden units according to the conditional probability found before. We then average over all samples in the training data to calculate expectation values with respect to the data.

Kullback-Leibler relative entropy

When the goal of the training is to approximate a probability distribution, as it is in generative modeling, another relevant measure is the **Kullback-Leibler divergence**, also known as the relative entropy or Shannon entropy. It is a non-symmetric measure of the dissimilarity between two probability density functions p and q . If p is the unknown probability which we approximate with q , we can measure the difference by

$$\text{KL}(p||q) = \int_{-\infty}^{\infty} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}.$$

Kullback-Leibler divergence

Thus, the Kullback-Leibler divergence between the distribution of the training data $f(\mathbf{x})$ and the model distribution $p(\mathbf{x}|\Theta)$ is

$$\begin{aligned}\text{KL}(f(\mathbf{x})||p(\mathbf{x}|\Theta)) &= \int_{-\infty}^{\infty} f(\mathbf{x}) \log \frac{f(\mathbf{x})}{p(\mathbf{x}|\Theta)} d\mathbf{x} \\ &= \int_{-\infty}^{\infty} f(\mathbf{x}) \log f(\mathbf{x}) d\mathbf{x} - \int_{-\infty}^{\infty} f(\mathbf{x}) \log p(\mathbf{x}|\Theta) d\mathbf{x} \\ &= \langle \log f(\mathbf{x}) \rangle_{f(\mathbf{x})} - \langle \log p(\mathbf{x}|\Theta) \rangle_{f(\mathbf{x})} \\ &= \langle \log f(\mathbf{x}) \rangle_{data} + \langle E(\mathbf{x}) \rangle_{data} + \log Z \\ &= \langle \log f(\mathbf{x}) \rangle_{data} + \mathcal{C}_{LL}.\end{aligned}$$

Maximizing log-likelihood

The first term is constant with respect to Θ since $f(\mathbf{x})$ is independent of Θ . Thus the Kullback-Leibler Divergence is minimal when the second term is minimal. The second term is the log-likelihood cost function, hence minimizing the Kullback-Leibler divergence is equivalent to maximizing the log-likelihood.

To further understand generative models it is useful to study the gradient of the cost function which is needed in order to minimize it using methods like stochastic gradient descent.

More on the partition function

The partition function is the generating function of expectation values, in particular there are mathematical relationships between expectation values and the log-partition function. In this case we have

$$\left\langle \frac{\partial E(\mathbf{x}; \Theta_i)}{\partial \Theta_i} \right\rangle_{model} = \int p(\mathbf{x}|\Theta) \frac{\partial E(\mathbf{x}; \Theta_i)}{\partial \Theta_i} d\mathbf{x} = -\frac{\partial \log Z(\Theta_i)}{\partial \Theta_i}.$$

Here $\langle \cdot \rangle_{model}$ is the expectation value over the model probability distribution $p(\mathbf{x}|\Theta)$.

Setting up for gradient descent calculations

Using the previous relationship we can express the gradient of the cost function as

$$\begin{aligned}\frac{\partial \mathcal{C}_{LL}}{\partial \Theta_i} &= \left\langle \frac{\partial E(\mathbf{x}; \Theta_i)}{\partial \Theta_i} \right\rangle_{data} + \frac{\partial \log Z(\Theta_i)}{\partial \Theta_i} \\ &= \left\langle \frac{\partial E(\mathbf{x}; \Theta_i)}{\partial \Theta_i} \right\rangle_{data} - \left\langle \frac{\partial E(\mathbf{x}; \Theta_i)}{\partial \Theta_i} \right\rangle_{model}\end{aligned}$$

Difference of moments

This expression shows that the gradient of the log-likelihood cost function is a **difference of moments**, with one calculated from the data and one calculated from the model. The data-dependent term is called the **positive phase** and the model-dependent term is called the **negative phase** of the gradient. We see now that minimizing the cost function results in lowering the energy of configurations \mathbf{x} near points in the training data and increasing the energy of configurations not observed in the training data. That means we increase the model's probability of configurations similar to those in the training data.

More observations

The gradient of the cost function also demonstrates why gradients of unsupervised, generative models must be computed differently from those of, for example, FNNs. While the data-dependent expectation value is easily calculated based on the samples \mathbf{x}_i in the training data, we must sample from the model in order to generate samples from which to calculate the model-dependent term. We sample from the model by using MCMC-based methods. We cannot sample from the model directly because the partition function Z is generally intractable.

Adding hyperparameters

As in supervised machine learning problems, the goal is also here to perform well on **unseen** data, that is to have good generalization from the training data. The distribution $f(x)$ we approximate is not the **true** distribution we wish to estimate, it is limited to the training data. Hence, in unsupervised training as well it is important to prevent overfitting to the training data. Thus it is common to add regularizers to the cost function in the same manner as we discussed for say linear regression.

Theory of Variational Autoencoders

Let us remind ourself about what an autoencoder is, see the jupyter-notebook at <https://github.com/CompPhysics/AdvancedMachineLearning/blob/main/doc/pub/week10/ipynb/week10.ipynb>.

The Autoencoder again

Autoencoders are neural networks where the outputs are its own inputs. They are split into an **encoder part** which maps the input \mathbf{x} via a function $f(\mathbf{x}, \mathbf{W})$ (this is the encoder part) to a **so-called code part** (or intermediate part) with the result \mathbf{h}

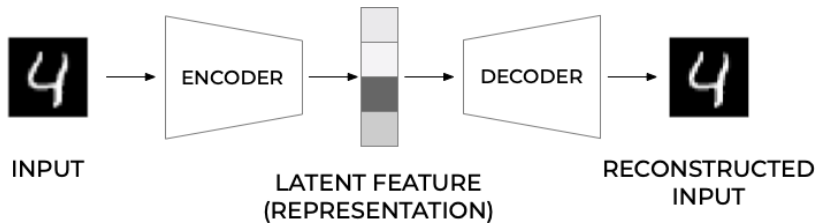
$$\mathbf{h} = f(\mathbf{x}, \mathbf{W}),$$

where \mathbf{W} are the weights to be determined. The **decoder** parts maps, via its own parameters (weights given by the matrix \mathbf{V} and its own biases) to the final output

$$\tilde{\mathbf{x}} = g(\mathbf{h}, \mathbf{V}).$$

The goal is to minimize the construction error, often done by optimizing the means squared error.

Schematic image of an Autoencoder



Mathematics of Variational Autoencoders

We have defined earlier a probability (marginal) distribution with hidden variables \mathbf{h} and parameters Θ as

$$p(\mathbf{x}; \Theta) = \int d\mathbf{h} p(\mathbf{x}, \mathbf{h}; \Theta),$$

for continuous variables \mathbf{h} and

$$p(\mathbf{x}; \Theta) = \sum_{\mathbf{h}} p(\mathbf{x}, \mathbf{h}; \Theta),$$

for discrete stochastic events \mathbf{h} . The variables \mathbf{h} are normally called the **latent variables** in the theory of autoencoders. We will also call them for that here.

Using the conditional probability

Using the the definition of the conditional probabilities $p(\mathbf{x}|\mathbf{h}; \Theta)$, $p(\mathbf{h}|\mathbf{x}; \Theta)$ and and the prior $p(\mathbf{h})$, we can rewrite the above equation as

$$p(\mathbf{x}; \Theta) = \sum_{\mathbf{h}} p(\mathbf{x}|\mathbf{h}; \Theta)p(\mathbf{h},$$

which allows us to make the dependence of \mathbf{x} on \mathbf{h} explicit by using the law of total probability. The intuition behind this approach for finding the marginal probability for \mathbf{x} is to optimize the above equations with respect to the parameters Θ . This is done normally by maximizing the probability, the so-called maximum-likelihood approach discussed earlier.

VAEs versus autoencoders

This trained probability is assumed to be able to produce similar samples as the input. In VAEs it is then common to compare via for example the mean-squared error or the cross-entropy the predicted values with the input values. Compared with autoencoders, we are now producing a probability instead of a functions which mimicks the input.

In VAEs, the choice of this output distribution is often Gaussian, meaning that the conditional probability is

$$p(\mathbf{x}|\mathbf{h}; \Theta) = N(\mathbf{x}|f(\mathbf{h}; \Theta), \sigma^2 \times \mathbf{I}),$$

with mean value given by the function $f(\mathbf{h}; \Theta)$ and a diagonal covariance matrix multiplied by a parameter σ^2 which is treated as a hyperparameter.

Gradient descent

By having a Gaussian distribution, we can use gradient descent (or any other optimization technique) to increase $p(\mathbf{x}; \Theta)$ by making $f(\mathbf{h}; \Theta)$ approach \mathbf{x} for some \mathbf{h} , gradually making the training data more likely under the generative model. The important property is simply that the marginal probability can be computed, and it is continuous in Θ .

Are VAEs just modified autoencoders?

The mathematical basis of VAEs actually has relatively little to do with classical autoencoders, for example the sparse autoencoders or denoising autoencoders discussed earlier.

VAEs approximately maximize the probability equation discussed above. They are called autoencoders only because the final training objective that derives from this setup does have an encoder and a decoder, and resembles a traditional autoencoder. Unlike sparse autoencoders, there are generally no tuning parameters analogous to the sparsity penalties. And unlike sparse and denoising autoencoders, we can sample directly from $p(\mathbf{x})$ without performing Markov Chain Monte Carlo.

Training VAEs

To solve the integral or sum for $p(\mathbf{x})$, there are two problems that VAEs must deal with: how to define the latent variables \mathbf{h} , that is decide what information they represent, and how to deal with the integral over \mathbf{h} . VAEs give a definite answer to both.

Kullback-Leibler relative entropy

When the goal of the training is to approximate a probability distribution, as it is in generative modeling, another relevant measure is the **Kullback-Leibler divergence**, also known as the relative entropy or Shannon entropy. It is a non-symmetric measure of the dissimilarity between two probability density functions p and q . If p is the unknown probability which we approximate with q , we can measure the difference by

$$\text{KL}(p||q) = \int_{-\infty}^{\infty} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}.$$

Kullback-Leibler divergence and RBMs

Thus, the Kullback-Leibler divergence between the distribution of the training data $f(\mathbf{x})$ and the model marginal distribution $p(\mathbf{x}; \Theta)$ from an RBM is

$$\begin{aligned}\text{KL}(f(\mathbf{x})||p(\mathbf{x}|\Theta)) &= \int_{-\infty}^{\infty} f(\mathbf{x}) \log \frac{f(\mathbf{x})}{p(\mathbf{x}; \Theta)} d\mathbf{x} \\ &= \int_{-\infty}^{\infty} f(\mathbf{x}) \log f(\mathbf{x}) d\mathbf{x} - \int_{-\infty}^{\infty} f(\mathbf{x}) \log p(\mathbf{x}; \Theta) d\mathbf{x} \\ &= \langle \log f(\mathbf{x}) \rangle_{f(\mathbf{x})} - \langle \log p(\mathbf{x}; \Theta) \rangle_{f(\mathbf{x})} \\ &= \langle \log f(\mathbf{x}) \rangle_{data} + \langle E(\mathbf{x}) \rangle_{data} + \log Z.\end{aligned}$$

Maximizing log-likelihood

The first term is constant with respect to Θ since $f(\mathbf{x})$ is independent of Θ . Thus the Kullback-Leibler divergence is minimal when the second term is minimal. The second term is the log-likelihood cost function, hence minimizing the Kullback-Leibler divergence is equivalent to maximizing the log-likelihood.

Back to VAEs

We want to train the marginal probability with some latent variables \mathbf{h}

$$p(\mathbf{x}; \Theta) = \int d\mathbf{h} p(\mathbf{x}, \mathbf{h}; \Theta),$$

for the continuous version (see previous slides for the discrete variant).

Using the KL divergence

In practice, for most \mathbf{h} , $p(\mathbf{x}|\mathbf{h}; \Theta)$ will be nearly zero, and hence contribute almost nothing to our estimate of $p(\mathbf{x})$.

The key idea behind the variational autoencoder is to attempt to sample values of \mathbf{h} that are likely to have produced \mathbf{x} , and compute $p(\mathbf{x})$ just from those.

This means that we need a new function $Q(\mathbf{h}|\mathbf{x})$ which can take a value of \mathbf{x} and give us a distribution over \mathbf{h} values that are likely to produce \mathbf{x} . Hopefully the space of \mathbf{h} values that are likely under Q will be much smaller than the space of all \mathbf{h} 's that are likely under the prior $p(\mathbf{h})$. This lets us, for example, compute $E_{\mathbf{h} \sim Q} p(\mathbf{x}|\mathbf{h})$ relatively easily. Note that we drop Θ from here and for notational simplicity.

Kullback-Leibler again

However, if \mathbf{h} is sampled from an arbitrary distribution with PDF $Q(\mathbf{h})$, which is not $\mathcal{N}(0, I)$, then how does that help us optimize $p(\mathbf{x})$?

The first thing we need to do is relate $E_{\mathbf{h} \sim Q} P(\mathbf{x}|\mathbf{h})$ and $p(\mathbf{x})$. We will see where Q comes from later.

The relationship between $E_{\mathbf{h} \sim Q} p(\mathbf{x}|\mathbf{h})$ and $p(\mathbf{x})$ is one of the cornerstones of variational Bayesian methods. We begin with the definition of Kullback-Leibler divergence (KL divergence or \mathcal{D}) between $p(\mathbf{h}|\mathbf{x})$ and $Q(\mathbf{h})$, for some arbitrary Q (which may or may not depend on \mathbf{x}):

$$\mathcal{D}[Q(\mathbf{h})||p(\mathbf{h}|\mathbf{x})] = E_{\mathbf{h} \sim Q} [\log Q(\mathbf{h}) - \log p(\mathbf{h}|\mathbf{x})] .$$

And applying Bayes rule

We can get both $p(\mathbf{x})$ and $p(\mathbf{x}|\mathbf{h})$ into this equation by applying Bayes rule to $p(\mathbf{h}|\mathbf{x})$

$$\mathcal{D}[Q(\mathbf{h})\|p(\mathbf{h}|\mathbf{x})] = E_{\mathbf{h}\sim Q} [\log Q(\mathbf{h}) - \log p(\mathbf{x}|\mathbf{h}) - \log p(\mathbf{h})] + \log p(\mathbf{x}).$$

Here, $\log p(\mathbf{x})$ comes out of the expectation because it does not depend on \mathbf{h} . Negating both sides, rearranging, and contracting part of $E_{\mathbf{h}\sim Q}$ into a KL-divergence terms yields:

$$\log p(\mathbf{x}) - \mathcal{D}[Q(\mathbf{h})\|p(\mathbf{h}|\mathbf{x})] = E_{\mathbf{h}\sim Q} [\log p(\mathbf{x}|\mathbf{h})] - \mathcal{D}[Q(\mathbf{h})\|P(\mathbf{h})].$$

Rearranging

Using Bayes rule we obtain

$$E_{\mathbf{h} \sim Q} [\log p(y_i | \mathbf{h}, x_i)] = E_{\mathbf{h} \sim Q} [\log p(\mathbf{h} | y_i, x_i) - \log p(\mathbf{h} | x_i) + \log p(y_i | x_i)]$$

Rearranging the terms and subtracting $E_{\mathbf{h} \sim Q} \log Q(\mathbf{h})$ from both sides gives

$$\begin{aligned} \log P(y_i | x_i) - E_{\mathbf{h} \sim Q} [\log Q(\mathbf{h}) - \log p(\mathbf{h} | x_i, y_i)] = \\ E_{\mathbf{h} \sim Q} [\log p(y_i | \mathbf{h}, x_i) + \log p(\mathbf{h} | x_i) - \log Q(\mathbf{h})] \end{aligned}$$

Note that \mathbf{x} is fixed, and Q can be *any* distribution, not just a distribution which does a good job mapping \mathbf{x} to the \mathbf{h} 's that can produce X .

Inferring the probability

Since we are interested in inferring $p(\mathbf{x})$, it makes sense to construct a Q which *does* depend on \mathbf{x} , and in particular, one which makes $\mathcal{D}[Q(\mathbf{h})||p(\mathbf{h}|\mathbf{x})]$ small

$$\log p(\mathbf{x}) - \mathcal{D}[Q(\mathbf{h}|\mathbf{x})||p(\mathbf{h}|\mathbf{x})] = E_{\mathbf{h} \sim Q} [\log p(\mathbf{x}|\mathbf{h})] - \mathcal{D}[Q(\mathbf{h}|\mathbf{x})||p(\mathbf{h})] .$$

Hence, during training, it makes sense to choose a Q which will make $E_{\mathbf{h} \sim Q} [\log Q(\mathbf{h}) - \log p(\mathbf{h}|\mathbf{x}_i, y_i)]$ (a \mathcal{D} -divergence) small, such that the right hand side is a close approximation to $\log p(y_i|\mathbf{x}_i)$.

Central equation of VAEs

This equation serves as the core of the variational autoencoder, and it is worth spending some time thinking about what it means.

1. The left hand side has the quantity we want to maximize, namely $\log p(\mathbf{x})$ plus an error term.
2. The right hand side is something we can optimize via stochastic gradient descent given the right choice of Q .

Setting up SGD

So how can we perform stochastic gradient descent?

First we need to be a bit more specific about the form that $Q(\mathbf{h}|\mathbf{x})$ will take. The usual choice is to say that

$Q(\mathbf{h}|\mathbf{x}) = \mathcal{N}(\mathbf{h}|\mu(\mathbf{x}; \vartheta), \Sigma(\mathbf{x}; \vartheta))$, where μ and Σ are arbitrary deterministic functions with parameters ϑ that can be learned from data (we will omit ϑ in later equations). In practice, μ and Σ are again implemented via neural networks, and Σ is constrained to be a diagonal matrix.

More on the SGD

The name variational “autoencoder” comes from the fact that μ and Σ are “encoding” \mathbf{x} into the latent space \mathbf{h} . The advantages of this choice are computational, as they make it clear how to compute the right hand side. The last term— $\mathcal{D}[Q(\mathbf{h}|\mathbf{x})||p(\mathbf{h})]$ —is now a KL-divergence between two multivariate Gaussian distributions, which can be computed in closed form as:

$$\mathcal{D}[\mathcal{N}(\mu_0, \Sigma_0)||\mathcal{N}(\mu_1, \Sigma_1)] = \frac{1}{2} \left(\text{tr}(\Sigma_1^{-1}\Sigma_0) + (\mu_1 - \mu_0)^\top \Sigma_1^{-1}(\mu_1 - \mu_0) - k + \log \left(\frac{\det \Sigma_1}{\det \Sigma_0} \right) \right)$$

where k is the dimensionality of the distribution.

Simplification

In our case, this simplifies to:

$$\mathcal{D}[\mathcal{N}(\mu(X), \Sigma(X)) \| \mathcal{N}(0, I)] = \frac{1}{2} \left(\text{tr}(\Sigma(X)) + (\mu(X))^{\top} (\mu(X)) - k - \log \det(\Sigma(X)) \right).$$

Terms to compute

The first term on the right hand side is a bit more tricky. We could use sampling to estimate $E_{z \sim Q} [\log P(X|z)]$, but getting a good estimate would require passing many samples of z through f , which would be expensive. Hence, as is standard in stochastic gradient descent, we take one sample of z and treat $\log P(X|z)$ for that z as an approximation of $E_{z \sim Q} [\log P(X|z)]$. After all, we are already doing stochastic gradient descent over different values of X sampled from a dataset D . The full equation we want to optimize is:

$$E_{X \sim D} [\log P(X) - \mathcal{D} [Q(z|X) \| P(z|X)]] = \\ E_{X \sim D} [E_{z \sim Q} [\log P(X|z)] - \mathcal{D} [Q(z|X) \| P(z)]] .$$

Computing the gradients

If we take the gradient of this equation, the gradient symbol can be moved into the expectations. Therefore, we can sample a single value of X and a single value of z from the distribution $Q(z|X)$, and compute the gradient of:

$$\log P(X|z) - \mathcal{D}[Q(z|X)||P(z)]. \quad (1)$$

We can then average the gradient of this function over arbitrarily many samples of X and z , and the result converges to the gradient. There is, however, a significant problem $E_{z \sim Q} [\log P(X|z)]$ depends not just on the parameters of P , but also on the parameters of Q . In order to make VAEs work, it is essential to drive Q to produce codes for X that P can reliably decode.

$$E_{X \sim D} \left[E_{\epsilon \sim \mathcal{N}(0, I)} [\log P(X|z = \mu(X) + \Sigma^{1/2}(X) * \epsilon)] - \mathcal{D}[Q(z|X)||P(z)] \right]$$

Motivation from Kingma and Welling, An Introduction to Variational Autoencoders,

<https://arxiv.org/abs/1906.02691>

There are many reasons why generative modeling is attractive. First, we can express physical laws and constraints into the generative process while details that we don't know or care about, i.e. nuisance variables, are treated as noise. The resulting models are usually highly intuitive and interpretable and by testing them against observations we can confirm or reject our theories about how the world works. Another reason for trying to understand the generative process of data is that it naturally expresses causal relations of the world. Causal relations have the great advantage that they generalize much better to new situations than mere correlations. For instance, once we understand the generative process of an earthquake, we can use that knowledge both in California and in Chile.

Mathematics of VAEs

We want to train the marginal probability with some latent variables \mathbf{h}

$$p(\mathbf{x}; \Theta) = \int d\mathbf{h} p(\mathbf{x}, \mathbf{h}; \Theta),$$

for the continuous version (see previous slides for the discrete variant).

Using the KL divergence

In practice, for most \mathbf{h} , $p(\mathbf{x}|\mathbf{h}; \Theta)$ will be nearly zero, and hence contributes almost nothing to our estimate of $p(\mathbf{x})$.

The key idea behind the variational autoencoder is to attempt to sample values of \mathbf{h} that are likely to have produced \mathbf{x} , and compute $p(\mathbf{x})$ just from those.

This means that we need a new function $Q(\mathbf{h}|\mathbf{x})$ which can take a value of \mathbf{x} and give us a distribution over \mathbf{h} values that are likely to produce \mathbf{x} . Hopefully the space of \mathbf{h} values that are likely under Q will be much smaller than the space of all \mathbf{h} 's that are likely under the prior $p(\mathbf{h})$. This lets us, for example, compute $E_{\mathbf{h} \sim Q} p(\mathbf{x}|\mathbf{h})$ relatively easily. Note that we drop Θ from here and for notational simplicity.

Kullback-Leibler again

However, if \mathbf{h} is sampled from an arbitrary distribution with PDF $Q(\mathbf{h})$, which is not $\mathcal{N}(0, I)$, then how does that help us optimize $p(\mathbf{x})$?

The first thing we need to do is relate $E_{\mathbf{h} \sim Q} P(\mathbf{x}|\mathbf{h})$ and $p(\mathbf{x})$. We will see where Q comes from later.

The relationship between $E_{\mathbf{h} \sim Q} p(\mathbf{x}|\mathbf{h})$ and $p(\mathbf{x})$ is one of the cornerstones of variational Bayesian methods. We begin with the definition of Kullback-Leibler divergence (KL divergence or \mathcal{D}) between $p(\mathbf{h}|\mathbf{x})$ and $Q(\mathbf{h})$, for some arbitrary Q (which may or may not depend on \mathbf{x}):

$$\mathcal{D}[Q(\mathbf{h})||p(\mathbf{h}|\mathbf{x})] = E_{\mathbf{h} \sim Q} [\log Q(\mathbf{h}) - \log p(\mathbf{h}|\mathbf{x})] .$$

And applying Bayes rule

We can get both $p(\mathbf{x})$ and $p(\mathbf{x}|\mathbf{h})$ into this equation by applying Bayes rule to $p(\mathbf{h}|\mathbf{x})$

$$\mathcal{D}[Q(\mathbf{h})\|p(\mathbf{h}|\mathbf{x})] = E_{\mathbf{h}\sim Q} [\log Q(\mathbf{h}) - \log p(\mathbf{x}|\mathbf{h}) - \log p(\mathbf{h})] + \log p(\mathbf{x}).$$

Here, $\log p(\mathbf{x})$ comes out of the expectation because it does not depend on \mathbf{h} . Negating both sides, rearranging, and contracting part of $E_{\mathbf{h}\sim Q}$ into a KL-divergence terms yields:

$$\log p(\mathbf{x}) - \mathcal{D}[Q(\mathbf{h})\|p(\mathbf{h}|\mathbf{x})] = E_{\mathbf{h}\sim Q} [\log p(\mathbf{x}|\mathbf{h})] - \mathcal{D}[Q(\mathbf{h})\|P(\mathbf{h})].$$

Rearranging

Using Bayes rule we obtain

$$E_{\mathbf{h} \sim Q} [\log p(y_i | \mathbf{h}, x_i)] = E_{\mathbf{h} \sim Q} [\log p(\mathbf{h} | y_i, x_i) - \log p(\mathbf{h} | x_i) + \log p(y_i | x_i)]$$

Rearranging the terms and subtracting $E_{\mathbf{h} \sim Q} \log Q(\mathbf{h})$ from both sides gives

$$\begin{aligned} \log P(y_i | x_i) - E_{\mathbf{h} \sim Q} [\log Q(\mathbf{h}) - \log p(\mathbf{h} | x_i, y_i)] = \\ E_{\mathbf{h} \sim Q} [\log p(y_i | \mathbf{h}, x_i) + \log p(\mathbf{h} | x_i) - \log Q(\mathbf{h})] \end{aligned}$$

Note that \mathbf{x} is fixed, and Q can be *any* distribution, not just a distribution which does a good job mapping \mathbf{x} to the \mathbf{h} 's that can produce X .

Inferring the probability

Since we are interested in inferring $p(\mathbf{x})$, it makes sense to construct a Q which *does* depend on \mathbf{x} , and in particular, one which makes $\mathcal{D}[Q(\mathbf{h})\|p(\mathbf{h}|\mathbf{x})]$ small

$$\log p(\mathbf{x}) - \mathcal{D}[Q(\mathbf{h}|\mathbf{x})\|p(\mathbf{h}|\mathbf{x})] = E_{\mathbf{h} \sim Q}[\log p(\mathbf{x}|\mathbf{h})] - \mathcal{D}[Q(\mathbf{h}|\mathbf{x})\|p(\mathbf{h})].$$

Hence, during training, it makes sense to choose a Q which will make $E_{\mathbf{h} \sim Q}[\log Q(\mathbf{h}) - \log p(\mathbf{h}|\mathbf{x}_i, y_i)]$ (a \mathcal{D} -divergence) small, such that the right hand side is a close approximation to $\log p(y_i|\mathbf{x}_i)$.

Central equation of VAEs

This equation serves as the core of the variational autoencoder, and it is worth spending some time thinking about what it means.

1. The left hand side has the quantity we want to maximize, namely $\log p(\mathbf{x})$ plus an error term.
2. The right hand side is something we can optimize via stochastic gradient descent given the right choice of Q .

Setting up SGD

So how can we perform stochastic gradient descent?

First we need to be a bit more specific about the form that $Q(\mathbf{h}|\mathbf{x})$ will take. The usual choice is to say that

$Q(\mathbf{h}|\mathbf{x}) = \mathcal{N}(\mathbf{h}|\mu(\mathbf{x}; \vartheta), \Sigma(\mathbf{x}; \vartheta))$, where μ and Σ are arbitrary deterministic functions with parameters ϑ that can be learned from data (we will omit ϑ in later equations). In practice, μ and Σ are again implemented via neural networks, and Σ is constrained to be a diagonal matrix.

More on the SGD

The name variational “autoencoder” comes from the fact that μ and Σ are “encoding” \mathbf{x} into the latent space \mathbf{h} . The advantages of this choice are computational, as they make it clear how to compute the right hand side. The last term— $\mathcal{D}[Q(\mathbf{h}|\mathbf{x})||p(\mathbf{h})]$ —is now a KL-divergence between two multivariate Gaussian distributions, which can be computed in closed form as:

$$\mathcal{D}[\mathcal{N}(\mu_0, \Sigma_0)||\mathcal{N}(\mu_1, \Sigma_1)] = \frac{1}{2} \left(\text{tr}(\Sigma_1^{-1}\Sigma_0) + (\mu_1 - \mu_0)^\top \Sigma_1^{-1}(\mu_1 - \mu_0) - k + \log \left(\frac{\det \Sigma_1}{\det \Sigma_0} \right) \right)$$

where k is the dimensionality of the distribution.

Simplification

In our case, this simplifies to:

$$\mathcal{D}[\mathcal{N}(\mu(X), \Sigma(X)) \| \mathcal{N}(0, I)] = \frac{1}{2} \left(\text{tr}(\Sigma(X)) + (\mu(X))^{\top} (\mu(X)) - k - \log \det(\Sigma(X)) \right).$$

Terms to compute

The first term on the right hand side is a bit more tricky. We could use sampling to estimate $E_{z \sim Q} [\log P(X|z)]$, but getting a good estimate would require passing many samples of z through f , which would be expensive. Hence, as is standard in stochastic gradient descent, we take one sample of z and treat $\log P(X|z)$ for that z as an approximation of $E_{z \sim Q} [\log P(X|z)]$. After all, we are already doing stochastic gradient descent over different values of X sampled from a dataset D . The full equation we want to optimize is:

$$E_{X \sim D} [\log P(X) - \mathcal{D} [Q(z|X) \| P(z|X)]] = \\ E_{X \sim D} [E_{z \sim Q} [\log P(X|z)] - \mathcal{D} [Q(z|X) \| P(z)]] .$$

Computing the gradients

If we take the gradient of this equation, the gradient symbol can be moved into the expectations. Therefore, we can sample a single value of X and a single value of z from the distribution $Q(z|X)$, and compute the gradient of:

$$\log P(X|z) - \mathcal{D}[Q(z|X)||P(z)]. \quad (2)$$

We can then average the gradient of this function over arbitrarily many samples of X and z , and the result converges to the gradient. There is, however, a significant problem $E_{z \sim Q} [\log P(X|z)]$ depends not just on the parameters of P , but also on the parameters of Q . In order to make VAEs work, it is essential to drive Q to produce codes for X that P can reliably decode.

$$E_{X \sim D} \left[E_{\epsilon \sim \mathcal{N}(0, I)} [\log P(X|z = \mu(X) + \Sigma^{1/2}(X) * \epsilon)] - \mathcal{D}[Q(z|X)||P(z)] \right]$$

Code examples using Keras

Code taken from

<https://keras.io/examples/generative/vae/>

"""

Title: Variational AutoEncoder

Author: [fchollet](https://twitter.com/fchollet)

Date created: 2020/05/03

Last modified: 2023/11/22

Description: Convolutional Variational AutoEncoder (VAE) trained on MNIST

Accelerator: GPU

"""

"""

Setup

"""

import os

os.environ["KERAS_BACKEND"] = "tensorflow"

import numpy as np

import tensorflow as tf

import keras

from keras import layers

"""

Create a sampling layer

"""

Code in PyTorch for VAEs

```
import torch
from torch.autograd import Variable
import numpy as np
import torch.nn.functional as F
import torchvision
from torchvision import transforms
import torch.optim as optim
from torch import nn
import matplotlib.pyplot as plt
from torch import distributions

class Encoder(torch.nn.Module):
    def __init__(self, D_in, H, latent_size):
        super(Encoder, self).__init__()
        self.linear1 = torch.nn.Linear(D_in, H)
        self.linear2 = torch.nn.Linear(H, H)
        self.enc_mu = torch.nn.Linear(H, latent_size)
        self.enc_log_sigma = torch.nn.Linear(H, latent_size)

    def forward(self, x):
        x = F.relu(self.linear1(x))
        x = F.relu(self.linear2(x))
        mu = self.enc_mu(x)
        log_sigma = self.enc_log_sigma(x)
        sigma = torch.exp(log_sigma)
        return torch.distributions.Normal(loc=mu, scale=sigma)
```