

AutoEncoder

ifding

Abstract

Probabilistic PCA, Variational AutoEncoder, ELBO surgery

1. Probabilistic PCA

Principal Component Analysis (known as PCA) is a technique used for dimensionality reduction, lossy data compression and data visualization. PCA is a technique that can be analyzed from two different points of view: linear algebra and probability. One of Probabilistic PCA is that it can be used as simple generative model.

PCA finds the principal components of data, or in other words, it finds the features, the directions where there is the most variance. Given a dataset $\mathbf{X} = \{\mathbf{x}^{(i)}\}_1^N$ with dimensionality D_x , the main idea behind PCA is to obtain a subspace (called the principal-component subspace) with dimension D_z , being $D_z \ll D_x$. So that each $\mathbf{x}^{(i)}$ is represented with a $\mathbf{z}^{(i)}$ (latent variable) in the best possible way. It's possible to recover $\mathbf{x}^{(i)}$ from $\mathbf{z}^{(i)}$.

$$\mathbf{x}^{(i)} \in \mathbb{R}^{D_x} \rightarrow \mathbf{z}^{(i)} \in \mathbb{R}^{D_z} \quad (1)$$

The probabilistic PCA is a maximum likelihood solution of a probabilistic latent variable model in which all marginal and conditional distributions are Gaussian. The likelihood function to be maximized is expressed as:

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \quad (2)$$

As all distributions are Gaussian, the probability distribution of \mathbf{x} can be expressed as $p(\mathbf{x}) = N(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C})$, where the covariance matrix is $\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}$. Therefore, the parameters of the model are the $D_x \times D_z$ matrix \mathbf{W} , the D_x dimensional vector $\boldsymbol{\mu}$ and the scalar σ^2 .

The optimal values of the parameters are obtained using Maximum Likelihood Estimation (MLE), maximizing

the log-likelihood $\log p(\mathbf{x})$.

$$\log p(\mathbf{x}|\mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \sum_{i=1}^N \log p(\mathbf{x}^{(i)}|\mathbf{W}, \boldsymbol{\mu}, \sigma^2) \quad (3)$$

The optimization can be computed in a closed form or using the EM algorithm. By closed form, the derivative of $\log p(\mathbf{x}|\mathbf{W}, \boldsymbol{\mu}, \sigma^2)$ with respect to each parameter is computed and equal to 0 to obtain the optimal values, namely $\boldsymbol{\mu}_{ML}$, \mathbf{W}_{ML} and σ_{ML}^2 .

The optimal mean $\boldsymbol{\mu}_{ML}$ is the sample mean of the data:

$$\boldsymbol{\mu}_{ML} = \sum_{i=1}^N \mathbf{x}^{(i)} = \bar{\mathbf{x}} \quad (4)$$

The σ_{ML}^2 represents the average variance associated with the discarded dimensions.

$$\sigma_{ML}^2 = \frac{1}{D_x - D_z} \sum_{i=D_z+1}^{D_x} \lambda_i \quad (5)$$

The generation phase comprises two steps. First, the variable \mathbf{z} is sampled from a priori known distribution $p(\mathbf{z})$. Then, new data samples are drawn from $p(\mathbf{x}|\mathbf{z}) = N(\mathbf{x}|\boldsymbol{\mu}(\mathbf{z}), \sigma^2\mathbf{I})$ where $\boldsymbol{\mu}(\mathbf{z}) = \mathbf{W}\mathbf{z} + \boldsymbol{\mu}$. This process is shown in figure 1.

This is a simple case of a generative model where the relation between the latent and the predictive variables is linear:

$$\mathbf{x} \sim \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon} \quad (6)$$

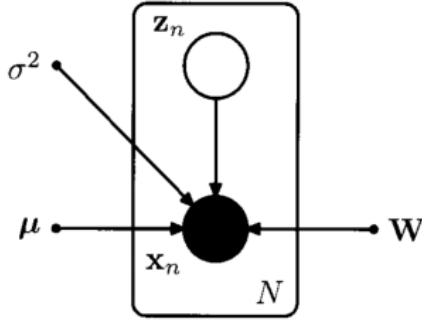


Figure 1: PCA generative model.

2. Variational AutoEncoder (VAE)

The VAE is a type of generative model, which is able to generate realistic data and obtain meaningful latent representations of the input. To achieve this, VAE is based on Bayesian inference. VAE models an approximation of the underlying probability distribution of input data $p(x)$ through a latent variable z . VAE can be understood as a non-linear Probabilistic PCA where the non-linearity is introduced using neural networks.

$$\begin{aligned} \text{Probabilistic PCA: } p(x|z) &= N(Wz + \mu, \sigma^2 I) \\ \text{VAE: } p_\theta(x|z) &= N(\mu_\theta(z), \sigma^2 I) \end{aligned} \quad (7)$$

What's a variational autoencoder? It's an autoencoder, with a random variable z defined in the latent space, that tries to maximize a lower bound of the data log-likelihood.

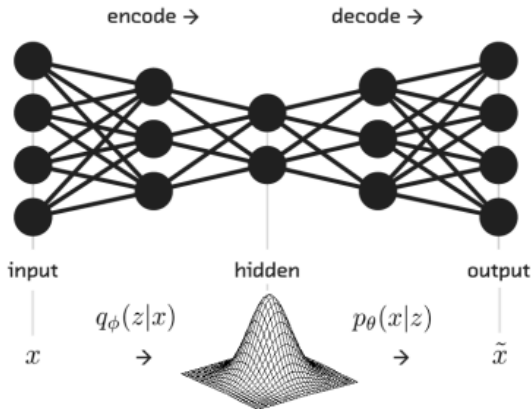


Figure 2: Probabilistic view of VAE.

The VAE is composed of two networks as shown in Figure 2 and Figure 3. The encoder or inference network (parameterized by ϕ) maps the input data into a latent representation that contains important attributes of the

data. This latent representation is then mapped, through the decoder of generative network (parameterized by θ), into the reconstructed data.



Figure 3: Graphical models for the VAE showing the generative model (left) and the inference model (right).

The main characteristic of VAE is that the latent representation of the input is motivated to be similar to a known distribution. The usual choice is to define the prior distribution $p(z) = N(z|0, I)$. In this way, if the latent variable is distributed similarly to $N(z|0, I)$ and it holds the correct information about the input data, it's possible to obtain samples from z and then generate realistic data.

Here are some important points about VAE:

- Similar inputs will have similar representations in the latent space.
- The dimensions of the latent variable z should be smaller than the dimension of the input so that it captures only the most important features.
- VAE builds an approximation of the true distribution of the input: $p_g(\mathbf{x}) \approx p(\mathbf{x})$.
- A lower bound of $\log p(\mathbf{x})$ called ELBO is the function that will be optimized.
- All density functions involved in the training process are prefixed. They are assumed to be Gaussian, it's possible to compute its derivatives and use optimization techniques.

2.1. Objective

Variational autoencoders approximately maximize the density function $p(\mathbf{x})$ of the training set according to the formula:

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

However, it's not possible to have infinite samples of z to compute the integral, the equation is intractable and

there will be losses. The objective function of the VAE, the Evidence Lower Bound (ELBO):

$$\log p(\mathbf{x}) \geq \mathbf{ELBO}(\theta, \phi, \mathbf{x}) = E_{q_\phi(z|\mathbf{x})} [\log p_\theta(\mathbf{x}|z)] - KL(q_\phi(z|\mathbf{x})||p(z)) \quad (9)$$

It's important to mention that only a subset of $\mathbf{X} = \{\mathbf{x}^{(i)}\}_1^N$ samples from the true distribution $p^*(\mathbf{x})$ is available. Then, the marginal log likelihood will be approximated as the following equation, which is an unbiased estimator:

$$E_{p^*(\mathbf{x})} [\log p_\theta(\mathbf{x})] \approx \frac{1}{N} \sum_{i=1}^N \log p_\theta(\mathbf{x}^{(i)}) \quad (10)$$

To reduce the memory needed to compute it, only a mini batch of size M ($M \ll N$) will be considered at a time, the estimator will remain unbiased but with greater variance.

$$E_{p^*(\mathbf{x})} [\log p_\theta(\mathbf{x})] \approx \frac{1}{M} \sum_{i=1}^M \log p_\theta(\mathbf{x}^{(i)}) \quad (11)$$

Finally, the ELBO will be estimated using M samples $\{\mathbf{x}^{(i)}\}_1^M$

$$\frac{1}{M} \sum_{i=1}^M \log p_\theta(\mathbf{x}^{(i)}) \geq \frac{1}{M} \sum_{i=1}^M \mathbf{ELBO}(\theta, \phi; \mathbf{x}^{(i)}) \quad (12)$$

1. Dealing with the Integral over z

To maximize $p(\mathbf{x})$, the key idea is to approximate the integral with n samples.

$$p(\mathbf{x}) = \int p(\mathbf{x}|z)p(z)dz = E_{z \sim p(z)} [p_\theta(\mathbf{x}|z)] \quad (13)$$

Then, if it is possible to obtain n samples from z , the equation can be written as:

$$p(\mathbf{x}) \approx \frac{1}{n} \sum_{i=1}^n p_\theta(\mathbf{x}|z^{(i)}) \quad (14)$$

By convenience, the $\log p_\theta(\mathbf{x}|z)$ will be used instead. It's assumed $p_\theta(\mathbf{x}|z)$ is distributed as a Gaussian $p_\theta(\mathbf{x}|z) = N(z|\mu_\theta(z), \Sigma)$. Therefore, its log likelihood is proportional to the euclidean distance between μ and \mathbf{x} which is simpler to manipulate analytically.

$$\log p_\theta(\mathbf{x}|z) = k - \frac{1}{2} (\mathbf{x} - \mu_\theta(z))^T \Sigma^{-1} (\mathbf{x} - \mu_\theta(z)) \quad (15)$$

where

$$k = -\log(|\Sigma|^{1/2}(2\pi)^{D_x/2}) \quad \Sigma = \sigma^2 \mathbf{I} \quad (16)$$

The value of σ^2 can be chosen, it's desirable to keep it small.

2. Inference Network: $q_\phi(z|\mathbf{x})$

According to equation 14, the approximation of $p(\mathbf{x})$ is better as the number of samples increases. However, the computation power needed also increases. In practice many values of z do not generate valid data. In conditional probability $p(\mathbf{x}|z)$ is close to 0 for the vast majority of z . This is why it is defined a family of conditional distributions $q_\phi(z|\mathbf{x})$ parameterized by ϕ that represents the distribution of z for each input datum. Thus, the new approximation is:

$$E_{z \sim q_\phi(z|\mathbf{x})} [\log(p_\theta(\mathbf{x}|z))] \quad (17)$$

In this way, only values of z that are likely to have produced \mathbf{x} are used to compute $p(\mathbf{x})$.

3. KL Divergence

The Kullback-Leibler divergence is a measure of how different are two probability functions p and q .

For discrete variables it is defines as follows:

$$KL(P(x)||Q(x)) = \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)} \quad (18)$$

For a continuous variable it is defined as follows:

$$KL(p(x)||q(x)) = \int p(x) \log \frac{p(x)}{q(x)} dx = E_{x \sim p(x)} [\log \frac{p(x)}{q(x)}] \quad (19)$$

During the VAE optimization, the $KL(q_\phi(z|\mathbf{x})||p(z|\mathbf{x}))$ should be close to 0.

$$KL(q_\phi(z|\mathbf{x})||p(z|\mathbf{x})) = E_{z \sim q_\phi} [\log q_\phi(z|\mathbf{x}) - \log p(z|\mathbf{x})] \quad (20)$$

4. Core Equation

Applying Bayes rule in $KL(q_\phi(z|\mathbf{x})||p(z|\mathbf{x}))$ and extracting $p(\mathbf{x})$ from the expectation because it does not depend on z :

$$p(z|\mathbf{x}) = \frac{p(\mathbf{x}|z)p(z)}{p(\mathbf{x})} \quad (21)$$

$$\begin{aligned} KL(q_\phi(z|\mathbf{x})||p(z|\mathbf{x})) &= E_{z \sim q_\phi} [\log q_\phi(z|\mathbf{x}) - \log p(z|\mathbf{x})] = \\ &= E_{z \sim q_\phi} [\log q_\phi(z|\mathbf{x}) - \log p_\theta(\mathbf{x}|z) - \log p(z)] + \log p(\mathbf{x}) = \\ &= -E_{z \sim q_\phi} [\log p_\theta(\mathbf{x}|z)] + E_{z \sim q_\phi} [\log q_\phi(z|\mathbf{x}) - \log p(z)] + \log p(\mathbf{x}) \end{aligned} \quad (22)$$

It's important to notice that the second term is a KL divergence. by rearranging terms, the fundamental equation of the VAE is obtained:

$$\log p(\mathbf{x}) - KL(q_\phi(z|\mathbf{x})||p(z|\mathbf{x})) = E_{z \sim q_\phi} [\log p_\theta(\mathbf{x}|z)] - KL(q_\phi(z|\mathbf{x})||p(z)) \quad (23)$$

Intuitively this equation tell us that the $\log p(\mathbf{x})$ minus the encoding error term, is equal to the mean of $\log p_\theta(\mathbf{x}|z)$ evaluated at z sampled from $q_\phi(z|\mathbf{x})$ minus an error term.

Because the $p(z|\mathbf{x})$ is unknown, it's not possible to calculate $KL(q_\phi(z|\mathbf{x})||p(z|\mathbf{x}))$. This is why this term will be removed from the equation, leading to the loss function of the VAE called ELBO:

$$ELBO(\theta, \phi) = E_{z \sim q_\phi} [\log p_\theta(\mathbf{x}|z)] - KL(q_\phi(z|\mathbf{x})||p(z)) \quad (24)$$

- The $E_{z \sim q_\phi} [\log p_\theta(\mathbf{x}|z)]$ term represents how good is the reconstruction of the input.
- The $KL(q_\phi(z|\mathbf{x})||p(z))$ term is a regularizer. In order to maximize the ELBO this term must be small. So this term forces $Q_\theta(z|\mathbf{x})$ to be similar to $p(z)$.

2.2. ELBO Optimization

In principle, all the distributions that appear in the ELBO are unknown, so they will be assumed to be Gaussian. This assumption is really convenience because it facilitates the computations since the ELBO will have a closed form. In addition, the Gaussian distribution has the minimal prior structure so it maximizes the entropy.

Then $q_\theta(z|\mathbf{x})$ is defined as:

$$q_\phi(z|\mathbf{x}) = N(\mu_\phi(\mathbf{x}), \Sigma_\phi(\mathbf{x})) \quad (25)$$

And $p(z)$ is defined as:

$$p(z) = N(0, \mathbf{I}) \quad (26)$$

Both μ and Σ are neural networks with parameters ϕ that are learnt during training. Moreover, Σ is defined as a diagonal matrix so that the $KL(q_\phi(z|\mathbf{x})||p(z|\mathbf{x}))$ has a closed form:

$$KL(N(\mu_\phi(\mathbf{x}), \Sigma_\phi(\mathbf{x}))||N(0, \mathbf{I})) = \frac{1}{2} \left[\text{tr}(\Sigma_\phi(\mathbf{x})) - D_z - \log \det(\Sigma_\phi(\mathbf{x})) + (\mu_\phi(\mathbf{x}))^T (\mu_\phi(\mathbf{x})) \right] \quad (27)$$

The other term involves solving the next integral:

$$E_{z \sim q_\phi} [\log p_\theta(\mathbf{x}|z)] = \int \log p_\theta(\mathbf{x}|z) q_\phi(z|\mathbf{x}) dz \quad (28)$$

It is approximated using the Monte Carlo estimator which is short it is an unbiased estimator, its variance is inversely proportional to the number of samples $\sigma^2 \propto \frac{1}{\sqrt{M}}$ and if M is large enough the estimated expectation is similar to the real one.

Then, the Monte Carlo estimator of $E_{z \sim q_\phi} [\log p_\theta(\mathbf{x}|z)]$ by taking M samples of z , $(z^{(1)}, \dots, z^{(M)})$, from $q_\phi(z|\mathbf{x})$ is:

$$E_{z \sim q_\phi} [\log p_\theta(\mathbf{x}|z)] \approx \frac{1}{M} \sum_{i=1}^M \log p_\theta(\mathbf{x}|z^{(i)}) \quad (29)$$

$$\log p_\theta(\mathbf{x}|z^{(i)}) = k - \frac{1}{2} (x - \mu_\theta(z^{(i)}))^T \Sigma^{-1} (x - \mu_\theta(z^{(i)})) \quad (30)$$

where

$$\Sigma = \sigma^2 \mathbf{I}_{D_x} \quad k = -\log \left(|\Sigma|^{1/2} (2\pi)^{D_x/2} \right) \quad |\Sigma| = (\sigma^2)^{D_x} \quad (31)$$

σ^2 acts as a weighting factor between the two terms. This parameter depends on the choice of $p_\theta(\mathbf{x}|z)$. Our choice of σ determines how accurately we expect the model to reconstruct \mathbf{x} . If σ^2 is small, the first term outweighs the second one forcing $\mathbf{x} \approx \mu_\theta(z)$.

This formula looks more tractable, but it can not be optimized via backpropagation. Because it's necessary to obtain z samples in an intermediate layer to calculate $\log p_\theta(\mathbf{x}|z^{(i)})$:

$$x \rightarrow \mu_\phi, \Sigma_\phi \rightarrow z \sim N(\mu_\phi, \Sigma_\phi) \rightarrow \mu_\theta \rightarrow x \quad (32)$$

The forward pass of this network works fine and, if the output is averaged over many samples of \mathbf{x} and z , it produces the correct expected value. However, it's not possible to back-propagate the error through stochastic layers. Sampling is a non-continuous operation and has no gradient.

2.3. Reparameterization Trick

In order to obtain the latent representation z of a given observation \mathbf{x} , it's necessary to sample it from the approximate posterior $q_\phi(z|\mathbf{x})$. The reparameterization trick consists in moving the sampling from an input layer. This is achieved defining a new input ϵ that introduces randomness into the model. Now, rather than sampling $z \sim N(\mu_\phi, \Sigma_\phi)$ directly, the sampling is done in $\epsilon \sim N(0, \mathbf{I})$

and then z is calculated in a deterministic way $z = \mu_\phi(\mathbf{x}) + \epsilon\sqrt{\Sigma_\phi(\mathbf{x})}$:

$$x, \epsilon \rightarrow \mu_\phi, \Sigma_\phi \rightarrow z = \mu_\phi + \sqrt{\Sigma_\phi} * \epsilon \rightarrow \mu_\theta \rightarrow x \quad (33)$$

In this way, backpropagation can be applied to optimize the network parameters.

The sampling of new data is quite straightforward: z is sampled from its prior distribution $p(z) = N(0, \mathbf{I})$ and injected into the decoder:

$$z \sim N(0, \mathbf{I}) \rightarrow \mu_\theta(z) \rightarrow x \quad (34)$$

3. ELBO surgery

The variational expectation-maximization (EM):

$$p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \quad (35)$$

where $p(\mathbf{z})$ is a prior on latent variables $\mathbf{z} = \{z_n\}_{n=1}^N$, $p_\theta(\mathbf{x}|\mathbf{z})$ is a likelihood on observations $\mathbf{x} = \{x_n\}_{n=1}^N$ parameterized by θ . The prior and likelihood follow from the generative model.

$$z_n \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \mathbf{I}) \quad (36)$$

$$x_n|z_n \sim \mathcal{N}(\mu(z_n; \theta), \Sigma(z_n; \theta)), \quad n = 1, 2, \dots, N$$

where the mean $\mu(z_n; \theta)$ and the covariance $\Sigma(z_n; \theta)$ depend on the latent variable z_n through a neural network with parameters θ .

We write joint densities as products of independent densities.

$$p(\mathbf{z}) = \prod_n p(z_n), \quad p_\theta(\mathbf{x}|\mathbf{z}) = \prod_n p_\theta(x_n|z_n) \quad (37)$$

The model is fit by maximizing the log evidence lower-bound (ELBO) \mathcal{L} ,

$$\begin{aligned} \log p_\theta(\mathbf{x}) &= \log \int p_\theta(\mathbf{z}, \mathbf{x})d\mathbf{z} = \log \int q_\phi(\mathbf{z}|\mathbf{x}) \frac{p_\theta(\mathbf{z}, \mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\ &\geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \log \frac{p_\theta(\mathbf{z}, \mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} \triangleq \mathcal{L}(\theta, \phi), \end{aligned} \quad (38)$$

where each term in the variational density $q_\phi(\mathbf{x}|\mathbf{z}) = \prod_n q_\phi(z_n|x_n)$ is a Gaussian in which the mean $\mu(x_n; \phi)$ and covariance $\Sigma(x_n; \phi)$ depend on the observation x_n through a neural network with free parameters ϕ .

The variational distribution $q_\phi(z_n|x_n)$ acts as a stochastic “encoder” from an observation x_n to a distribution on the latent variable z_n , and the likelihood $p_\theta(x_n|z_n)$ acts as a stochastic “decoder” from the latent variable z_n to a distribution on the observation x_n .

There are several ways to rewrite the objective $\mathcal{L}(\theta, \phi)$, and each provides its own perspective.

3.1. Evidence minus posterior KL

One form of $\mathcal{L}(\theta, \phi)$ emphasizes that the lower bound becomes tighter as the variational distribution better approximates the posterior:

$$\mathcal{L}(\theta, \phi) = \log p_\theta(\mathbf{x}) - \text{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})) \quad (39)$$

Thus we can improve the ELBO by improving the model log evidence $\log p_\theta(\mathbf{x})$, through the prior $p(\mathbf{z})$ or the likelihood $p_\theta(\mathbf{x}|\mathbf{z})$, or by improving the variational posterior approximation $q_\phi(\mathbf{z}|\mathbf{x})$.

3.2. Average negative energy plus entropy

Another way to rewrite the ELBO is as

$$\begin{aligned} \mathcal{L}(\theta, \phi) &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{z}, \mathbf{x})] + \mathbb{H}[q_\phi(\mathbf{z}|\mathbf{x})], \\ \mathbb{H}[q_\phi(\mathbf{z}|\mathbf{x})] &\triangleq -\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \log q_\phi(\mathbf{z}|\mathbf{x}), \end{aligned} \quad (40)$$

where the log joint $\log p_\theta(\mathbf{z}, \mathbf{x})$ is interpreted as the negative energy in a Boltzmann distribution. Since we choose (θ, ϕ) to maximize the ELBO, this version highlights that a good posterior approximation $q_\phi(\mathbf{z}|\mathbf{x})$ must assign most of its probability mass to regions of low energy (i.e. high joint probability density) while also maximizing the entropy of $q_\phi(\mathbf{z}|\mathbf{x})$.

This perspective is useful in contrasting variational EM with a maximum a-posteriori (MAP) approach; while MAP need only find a single value of \mathbf{z} that maximizes the joint density (even if it lies in a region with very low posterior mass), the entropy term in the ELBO prevents $q_\phi(\mathbf{z}|\mathbf{x})$ from collapsing to an atom.

3.3. Average term-by-term reconstruction minus KL to prior

Finally, we can write

$$\begin{aligned} \mathcal{L}(\theta, \phi) &= \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{q_\phi(z_n|x_n)} [\log p_\theta(x_n|z_n)] - \\ &\quad \text{KL}(q_\phi(z_n|x_n)||p(z_n)) \end{aligned} \quad (41)$$

For each observation index n , this version has a reconstruction term for the n th observation and a KL divergence from each encoding distribution to the prior. This KL-divergence term can be interpreted as a regularizer that “prune out” many of the latent dimensions in z .

3.4. Terms of the average encoding distribution

For the last decomposition, what’s a reasonable value for the KL-divergence term? Ideally it would be small, but we do not want it to approach 0 ($q_\phi(z_n|x_n) = p(z_n)$), since that would imply that x_n and z_n were almost independent, whereas virtually all of our modeling power comes from strongly coupling x_n to z_n . So if the KL term is large, is that a sign of underfitting, overfitting, or neither?

The *average encoding distribution* is defined as

$$q_\phi^{\text{avg}}(z) \triangleq \frac{1}{N} \sum_{n=1}^N q_\phi(z|x_n) \quad (42)$$

The *marginal* KL divergence $\text{KL}(q_\phi^{\text{avg}}(z)||p(z))$ is important because, unlike the individual terms $q_\phi(z_n|x_n)$, the average encoding distribution $q_\phi^{\text{avg}}(z)$ can be made arbitrarily close to the prior $p(z)$ without sacrificing model power. Indeed, if the data are drawn from the model, $x_n \sim p_\theta(x)$, and the posterior approximation is accurate, $q_\phi(z|x_n) \approx p_\theta(z|x_n)$, then for large N we would expect

$$\begin{aligned} p(z) &= \int p_\theta(z|x)p_\theta(x)dx = \mathbb{E}_{x \sim p_\theta(x)} p_\theta(z|x) \approx \\ &\frac{1}{N} \sum_n p_\theta(z|x_n) \approx \frac{1}{N} \sum_n q_\phi(z|x_n) = q_\phi^{\text{avg}}(z) \end{aligned} \quad (43)$$

To simplify the notation, we drop parameter subscripts, and it’s convenient to treat the index n as a random variable.

$$q(n, z) \triangleq q(n)q(z|n), \quad q(z|n) \triangleq q(z|x_n), \quad q(n) \triangleq \frac{1}{N} \quad (44)$$

$$p(n, z) \triangleq p(n)p(z|n), \quad p(z|n) \triangleq p(z), \quad p(n) \triangleq \frac{1}{N} \quad (45)$$

where $p(z)$ denotes a standard Gaussian prior density form $z \sim \mathcal{N}(0, I)$, the average encoder distribution $q^{\text{avg}}(z)$ is now simply the marginal $q(z) = \sum_{n=1}^N q(z, n)$.

$$\begin{aligned} \frac{1}{N} \sum_{n=1}^N \text{KL}(q(z_n|x_n)||p(z_n)) &= \sum_n q(n, z) \log \frac{q(n, z)}{p(n, z)} \\ &= q(z) \log \frac{q(z)}{p(z)} + \sum_n q(n|z)q(z) \log \frac{q(n|z)}{p(n|z)} \\ &= \text{KL}(q(z)||p(z)) + \mathbb{E}_{q(z)}[\text{KL}(q(n|z)||p(n))] \\ &= \text{KL}(q(z)||p(z)) + \mathbb{I}_{q(n, z)}[n, z] \\ &= \text{KL}(q(z)||p(z)) + (\log N - \mathbb{E}_{q(z)}[\mathbb{H}[q(n|z)]]) \end{aligned} \quad (46)$$

where $\mathbb{I}_{q(n, z)}[n, z] = \mathbb{E}_{q(n, z)} \left[\log \frac{q(n, z)}{q(n)q(z)} \right]$ denotes the mutual information of n and z in $q(n, z)$, $p(n) = \frac{1}{N}$. The mutual information expression:

$$\mathbb{I}_{q(n, z)}[n, z] = \mathbb{E}_{q(z)} \left[\mathbb{E}_{q(n|z)} \left[\log \frac{q(n|z)}{q(n)} \right] \right] = \log N - \mathbb{E}_{q(z)}[\mathbb{H}[q(n|z)]] \quad (47)$$

The ELBO can be written in three terms:

$$\begin{aligned} \mathcal{L}(\theta, \phi) &= \underbrace{\left[\frac{1}{N} \sum_{n=1}^N \mathbb{E}_{q(z_n|x_n)} [\log p(x_n|z_n)] \right]}_{(1) \text{ average reconstruction}} \\ &\quad - \underbrace{(\log N - \mathbb{E}_{q(z)}[\mathbb{H}[q(n|z)]])}_{(2) \text{ index-code mutual info.}} \\ &\quad - \underbrace{\text{KL}(q(z)||p(z))}_{(3) \text{ marginal KL to prior}} \end{aligned} \quad (48)$$

First, the two terms (1) and (2) are in tension with each other because to get a good average reconstruction score for (1), we typically need each encoding z_n to be specific to its corresponding observation x_n and hence $q(n|z)$ should have low entropy. Term (2) acts as a regularizer, in that it encourages the encodings $q(z|x_n)$ to overlap for distinct observation n , but this effect is likely to be weak relative to the reconstruction term (1). (2) is bounded above and below,

$$0 \leq \log N - \mathbb{E}_{q(z)} \mathbb{H}[q(n|z)] \leq \log N \quad (49)$$

Empirically, the reconstruction is precise and correspondingly, $q(z|n)$ is concentrated relative to $q(z)$, resulting in (2) is close to its maximum value of $\log N$.

Second, while $q(z)$ appears in all terms, $p(z)$ only appears in (3). Thus when considering choosing priors $p(z)$ to optimize the ELBO, only this term is affected. Observe that we could set (3) to zero without sacrificing model power by simply defining the prior to be $q(z)$. This choice would not be amenable to scalable computation because it is difficult to evaluate (2) in isolation: to normalize $q(n|z)$

at each evaluation requires accessing all N observations (and the normalization also precludes us from making unbiased Monte Carlo estimates). Setting (3) to zero may also be undesirable due to the potential for overfitting or the inability to use the prior to sculpt the latent representation. Because (3) can in principle be set to zero, whenever it is large it indicates a very strong and potentially unwanted regularization effect from the prior.

References

- [1] Sánchez Martín, P., 2018. Unsupervised deep learning research and implementation of variational autoencoders (Bachelor's thesis).
- [2] Bishop, Christopher M. Pattern recognition and machine learning. springer, 2006.
- [3] Kingma, D.P. and Welling, M., 2013. Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114.
- [4] Hoffman, M.D. and Johnson, M.J., 2016, December. Elbo surgery: yet another way to carve up the variational evidence lower bound. In Workshop in Advances in Approximate Bayesian Inference, NIPS (Vol. 1).