

# Assignment 3: Variational Autoencoders

STA414/STA2014 and CSC412/CSC2506 Winter 2020

David Duvenaud and Jesse Bettencourt

Student Number

March 27, 2020

## 0.1 Background

In this assignment, we will implement and investigate the Variational Autoencoder on binarized MNIST digits, as introduced by the paper [Auto-Encoding Variational Bayes](#) by Kingma and Welling (2013). Before starting, we recommend reading this paper.

**Data.** Each datapoint in the [MNIST](#) dataset is a 28x28 grayscale image (i.e. pixels are values between 0 and 1) of a handwritten digit in  $\{0 \dots 9\}$ , and a label indicating which number. MNIST is the ‘fruit fly’ of machine learning – a simple standard problem useful for comparing the properties of different algorithms.

Use the first 10000 samples for training, and the second 10000 for testing. Hint: Also build a dataset of only 100 training samples to use when debugging, to make loading and training faster.

**Tools.** As before, you can (and should) use automatic differentiation provided by your package of choice. Whereas in previous assignments you implemented neural network layers and stochastic gradient descent manually, in this assignment feel free to use those provided by a machine learning framework. In Julia, these will be provided by `Flux.jl`. You can also freely copy and adapt the Python autograd starter code provided. If you do, you should probably remove batch normalization.

However, you **may not use any probabilistic modelling elements** from these frameworks. In particular, sampling from and evaluating densities under distributions must be written by you or provided by the starter code.

## 0.2 Model Definition

**Prior.** The prior over each digit’s latent representation is a multivariate standard normal distribution. For all questions, we’ll set the dimension of the latent space  $D_z$  to 2. A larger latent dimension would provide a more powerful model, but for this assignment we’ll use a two-dimensional latent space to make visualization and debugging easier..

**Likelihood.** Given the latent representation  $z$  for an image, the distribution over all 784 pixels in the image is given by a product of independent Bernoullis, whose means are given by the output of a neural network  $f_\theta(z)$ :

$$p(x|z, \theta) = \prod_{d=1}^{784} \text{Ber}(x_d | f_\theta(z)_d)$$

The neural network  $f_\theta$  is called the decoder, and its parameters  $\theta$  will be optimized to fit the data.

# 1 Implementing the Model [5 points]

For your convenience we have provided the following functions:

- `factorized_gaussian_log_density` that accepts the mean and **log** standard deviations for a product of independent gaussian distributions and computes the likelihood under them. This function will produce the log-likelihood for each batch element  $1 \times B$
- `bernoulli_log_density` that accepts the logits of a bernoulli distribution over  $D$ -dimensional data and returns  $D \times B$  log-likelihoods.
- `sample_diag_gaussian` that accepts above parameters for a factorized Gaussian distribution and samples with the reparameterization trick.
- `sample_bernoulli` that accepts above parameters for a Bernoulli distribution and samples from it.
- `load_binarized_mnist` that loads and binarizes the MNIST dataset.
- `batch_x` and `batch_y` that splits the data, and just the images, into batches.

Further, in the file `example_flux_model.jl` we demonstrate how to specify neural network layers with Flux library. Note that Flux provides convenience functions that allow us to take gradients of functions with respect to parameters that **are not passed around explicitly**. Other AD frameworks, of if you prefer to implement your own network layers, recycling code from previous assignments, you may need to explicitly provide the network parameters to the functions below.

- (a) [1 point] Implement a function `log_prior` that computes the log of the prior over a digit's representation  $\log p(z)$ .
- (b) [2 points] Implement a function `decoder` that, given a latent representation  $z$  and a set of neural network parameters  $\theta$  (again, implicitly in Flux), produces a 784-dimensional mean vector of a product of Bernoulli distributions, one for each pixel in a  $28 \times 28$  image. Make the decoder architecture a multi-layer perceptron (i.e. a fully-connected neural network) with a single hidden layer with 500 hidden units, and a `tanh` nonlinearity. Its input will be a batch two-dimensional latent vectors ( $zs$  in  $D_z \times B$ ) and its output will be a 784-dimensional vector representing the logits of the Bernoulli means for each dimension  $D_{\text{data}} \times B$ . For numerical stability, instead of outputting the mean  $\mu \in [0, 1]$ , you should output  $\log\left(\frac{\mu}{1-\mu}\right) \in \mathbb{R}$  called "logit".
- (c) [1 point] Implement a function `log_likelihood` that, given a latent representation  $z$  and a binarized digit  $x$ , computes the log-likelihood  $\log p(x|z)$ .
- (d) [1 point] Implement a function `joint_log_density` which combines the log-prior and log-likelihood of the observations to give  $\log p(z, x)$  for a single image.

All of the functions in this section must be able to be evaluated in parallel, vectorized and non-mutating, on a batch of  $B$  latent vectors and images, using the same parameters  $\theta$  for each image. In particular, you can not use a for loop over the batch elements.

## 2 Amortized Approximate Inference and training [13 points]

- (a) [2 points] Write a function `encoder` that, given an image  $x$  (or batch of images) and recognition parameters  $\phi$ , evaluates an MLP to outputs the mean and log-standard deviation of a factorized Gaussian of dimension  $D_z = 2$ . Make the encoder architecture a multi-layer perceptron (i.e. a fully-connected neural network) with a single hidden layer with 500 hidden units, and a `tanh` nonlinearity. This function must be able to be evaluated in parallel on a batch of images, using the same parameters  $\phi$  for each image.
- (b) [1 points] Write a function `log_q` that given the parameters of the variational distribution, evaluates the likelihood of  $z$ .
- (c) [5 points] Implement a function `elbo` which computes an unbiased estimate of the mean variational evidence lower bound on a batch of images. Use the output of `encoder` to give the parameters for  $q_\phi(z|\text{data})$ . This estimator takes the following arguments:
- `x`, an batch of  $B$  images,  $D_x \times B$ .
  - `encoder_params`, the parameters  $\phi$  of the encoder (recognition network). Note: these are not required with Flux as parameters are implicit.
  - `decoder_params`, the parameters  $\theta$  of the decoder (likelihood). Note: these are not required with Flux as parameters are implicit.

This function should return a single scalar. Hint: You will need to use the reparameterization trick when sampling `zs`. You can use any form of the ELBO estimator you prefer. (i.e., if you want you can write the KL divergence between  $q$  and the prior in closed form since they are both Gaussians). You only need to sample a single  $z$  for each image in the batch.

- (d) [2 points] Write a loss function called `loss` that returns the negative elbo estimate over a batch of data.
- (e) [3 points] Write a function that initializes and optimizes the encoder and decoder parameters jointly on the training set. Note that this function should optimize with gradients on the elbo estimate over batches of data, not the entire dataset. Train the data for 100 epochs (each epoch involves a loop over every batch). Report the final ELBO on the test set. Tip: Save your trained weights here (e.g. with `JSON.jl`, see starter code, or by pickling in Python) so that you don't need to train them again.

### 3 Visualizing Posteriors and Exploring the Model [15 points]

In this section we will investigate our model by visualizing the distribution over data given by the generative model, sampling from it, and interpolating between digits.

- (a) [5 points] Plot samples from the trained generative model using ancestral sampling:
  - (a) First sample a  $z$  from the prior.
  - (b) Use the generative model to compute the bernoulli means over the pixels of  $x$  given  $z$ . Plot these means as a greyscale image.
  - (c) Sample a binary image  $x$  from this product of Bernoullis. Plot this sample as an image.

Do this for 10 samples  $z$  from the prior.

Concatenate all your plots into one 2x10 figure where each image in the first row shows the Bernoulli means of  $p(x|z)$  for a separate sample of  $z$ , and each image in the the second row is a binary image, sampled from the distribution above it. Make each column an independent sample.

- (b) [5 points] One way to understand the meaning of latent representations is to see which parts of the latent space correspond to which kinds of data. Here we'll produce a scatter plot in the latent space, where each point in the plot represents a different image in the training set.
  - (a) Encode each image in the training set.
  - (b) Take the 2D mean vector of each encoding  $q_\phi(z|x)$ .
  - (c) Plot these mean vectors in the 2D latent space with a scatterplot.
  - (d) Colour each point according to the class label (0 to 9).

Hopefully our latent space will group images of different classes, even though we never provided class labels to the model!

- (c) [5 points] Another way to examine a latent variable model with continuous latent variables is to interpolate between the latent representations of two points.

Here we will encode 3 pairs of data points with different classes. Then we will linearly interpolate between the mean vectors of their encodings. We will plot the generative distributions along the linear interpolation.

- (a) First, write a function which takes two points  $z_a$  and  $z_b$ , and a value  $\alpha \in [0, 1]$ , and outputs the linear interpolation  $z_\alpha = \alpha z_a + (1 - \alpha) z_b$ .
- (b) Sample 3 pairs of images, each having a different class.
- (c) Encode the data in each pair, and take the mean vectors
- (d) Linearly interpolate between these mean vectors
- (e) At 10 equally-space points along the interpolation, plot the Bernoulli means  $p(x|z_\alpha)$
- (f) Concatenate these plots into one figure.

## 4 Predicting the Bottom of Images given the Top [15 points]

Now we'll use the trained generative model to perform inference for  $p(z|\text{top half of image } x)$ . Unfortunately, we can't re-use our recognition network, since it can only input entire images. However, we can still do approximate inference without the encoder.

To illustrate this, we'll approximately infer the distribution over the pixels in the bottom half an image conditioned on the top half of the image:

$$p(\text{bottom half of image } x|\text{top half of image } x) = \int p(\text{bottom half of image } x|z)p(z|\text{top half of image } x)dz$$

To approximate the posterior  $p(z|\text{top half of image } x)$ , we'll use stochastic variational inference.

- (a) **[5 points]** Write a function that computes  $p(z, \text{top half of image } x)$
- (a) First, write a function which returns only the top half of a 28x28 array. This will be useful for plotting, as well as selecting the correct Bernoulli parameters.
  - (b) Write a function that computes  $\log p(\text{top half of image } x|z)$ . Hint: Given  $z$ , the likelihood factorizes, and all the unobserved dimensions of  $x$  are leaf nodes, so can be integrated out exactly.
  - (c) Combine this likelihood with the prior to get a function that takes an  $x$  and an array of  $z$ s, and computes the log joint density  $\log p(z, \text{top half of image } x)$  for each  $z$  in the array.
- (b) **[5 points]** Now, to approximate  $p(z|\text{top half of image } x)$  in a scalable way, we'll use stochastic variational inference. For a digit of your choosing from the training set (choose one that is modelled well, i.e. the resulting plot looks reasonable):
- (a) Initialize variational parameters  $\phi_\mu$  and  $\phi_{\log \sigma}$  for a variational distribution  $q(z|\text{top half of } x)$ .
  - (b) Write a function that computes estimates the ELBO over  $K$  samples  $z \sim q(z|\text{top half of } x)$ . Use  $\log p(z)$ ,  $\log p(\text{top half of } x|z)$ , and  $\log q(z|\text{top half of } x)$ .
  - (c) Optimize  $\phi_\mu$  and  $\phi_{\log \sigma}$  to maximize the ELBO.
  - (d) On a single plot, show the isocontours of the joint distribution  $p(z, \text{top half of image } x)$ , and the optimized approximate posterior  $q_\phi(z|\text{top half of image } x)$ .
  - (e) Finally, take a sample  $z$  from your approximate posterior, and feed it to the decoder to find the Bernoulli means of  $p(\text{bottom half of image } x|z)$ . Contatenate this greyscale image to the true top of the image. Plot the original whole image beside it for comparison.
- (c) **[5 points]** True or false: Questions about the model and variational inference.

There is no need to explain your work in this section.

- (a) Does the distribution over  $p(\text{bottom half of image } x|z)$  factorize over the pixels of the bottom half of image  $x$ ?
- (b) Does the distribution over  $p(\text{bottom half of image } x|\text{top half of image } x)$  factorize over the pixels of the bottom half of image  $x$ ?
- (c) When jointly optimizing the model parameters  $\theta$  and variational parameters  $\phi$ , if the ELBO increases, has the KL divergence between the approximate posterior  $q_\phi(z|x)$  and the true posterior  $p_\theta(z|x)$  necessarily gotten smaller?
- (d) If  $p(x) = \mathcal{N}(x|\mu, \sigma^2)$ , for some  $x \in \mathbb{R}, \mu \in \mathbb{R}, \sigma \in \mathbb{R}^+$ , can  $p(x) < 0$ ?
- (e) If  $p(x) = \mathcal{N}(x|\mu, \sigma^2)$ , for some  $x \in \mathbb{R}, \mu \in \mathbb{R}, \sigma \in \mathbb{R}^+$ , can  $p(x) > 1$ ?