

Problem Set 9

published: 21 June 2023
discussion: 28 June 2023

Data Compression With And Without Deep Probabilistic Models

Prof. Robert Bamler, University of Tübingen

Course materials available at <https://robamler.github.io/teaching/compress23/>

Note: The two problems on this set can be solved independently and in arbitrary order.

Problem 9.1: Variational Autoencoder

The jupyter notebook `vae-lossless-bitsback.ipynb` contains our toy variational autoencoder (VAE) from the lecture. This problem is designed to guide you through the existing implementation. You won't yet implement any new code in this problem.

- (a) **Data set:** always acquire some basic understanding of the data set that you deal with before you start implementing any fancy models. Here, we have 28×28 -pixel black and white images of handwritten digits (binarized MNIST). Let's get an upper bound on how much information one of these images contains. We could trivially store it as a string of 28×28 bits, but we can certainly do better if we take the probability distribution of the underlying data generative process into account.

We can't formally specify the true data generative process because it involved humans who wrote digits by hand at some point. But we do have a data set of samples from the data generative process, and we can fit a simple probabilistic model to these samples. The data set is split into a training set (`train_set`) and a `test_set`, where `train_set.data[n, i, j]` $\in \{0, 1\}$ denotes whether the pixel at horizontal position $i \in \{0, \dots, 27\}$ and vertical position $j \in \{0, \dots, 27\}$ in the n^{th} image of the training set is black or a white (analogously for the `test_set`).

Consider a model $P_\alpha(\mathbf{X})$ where the random variable $\mathbf{X} = (X_{i,j})_{i,j=0}^{27}$ denotes a single image, which is composed of 28×28 pixels $X_{i,j} \in \{0, 1\}$. To obtain a very simple baseline, we model all pixels $X_{i,j}$ as i.i.d. (independent and identically distributed). Thus, $P_\alpha(\mathbf{X}) = \prod_{i,j} P_\alpha(X_{i,j})$ factorizes, with the same marginal distribution $P_\alpha(X_{i,j})$ at every position (i, j) . Let's parameterize this distribution with a single parameter $\alpha \in [0, 1]$ such that $P_\alpha(X_{i,j}=1) = \alpha$ and thus $P_\alpha(X_{i,j}=0) = 1-\alpha$.

- (i) Which model parameter $\alpha^* \in [0, 1]$ minimizes the total information content under P_α of all images in the `train_set`? You should be able to formally derive an extremely simple and easily interpretable expression for α^* .
- (ii) Find the section "*Problem 9.1 (a): Trivial Baselines*" in the notebook. Read it, make sure you understand it, and execute the cells. This section fits the i.i.d. model and a slightly more general model to the training set and evaluates the information content of the training and test set under these two models. Why do you get a lower bit rate for the more general model?

- (b) **Entropy Bottleneck:** on Problem Set 8, you derived several equivalent formulations of the evidence lower bound (ELBO). One of them expressed the ELBO as a regularized maximum likelihood estimation,

$$\text{ELBO}(\theta, \phi, \mathbf{x}) = \mathbb{E}_{Q_\phi(\mathbf{Z})}[\log P_\theta(\mathbf{X}=\mathbf{x} | \mathbf{Z})] - D_{\text{KL}}(Q_\phi(\mathbf{Z} | \mathbf{X}=\mathbf{x}) \parallel P(\mathbf{Z})) \quad (1)$$

where we slightly changed the notation to follow conventions in the literature for amortized variational inference. The regularizer (second term on the right-hand side of Eq. 1) is the KL-divergence from the prior $P(\mathbf{Z})$ to the variational distribution $Q_\phi(\mathbf{Z} | \mathbf{X}=\mathbf{x})$. In our VAE, we can calculate this KL-divergence analytically because both prior and likelihood are normal distributions: $P(\mathbf{Z}) = \mathcal{N}(0, I)$ and $Q_\phi(\mathbf{Z} | \mathbf{X}=\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_\phi(\mathbf{x}), \text{diag}(\sigma_\phi(\mathbf{x})_1^2, \dots, \sigma_\phi(\mathbf{x})_n^2))$. Wikipedia states¹ that the KL-divergence between two k -dimensional normal distributions is:

$$D_{\text{KL}}(\mathcal{N}_0 \parallel \mathcal{N}_1) = \frac{1}{2} \left(\text{tr}(\Sigma_1^{-1} \Sigma_0) - k + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \Sigma_1^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) + \ln \left(\frac{\det \Sigma_1}{\det \Sigma_0} \right) \right). \quad (2)$$

- (i) Translate Eq. 2, which was copied verbatim from Wikipedia, to our specific setup and simplify as much as possible. You should be able to express the result as a sum over terms that don't contain any matrices or expensive matrix-operations like determinants.
 - (ii) Find the definition of the class `EntropyBottleneck` in the notebook and compare the implementation of its method `forward` to your result from (i).
- (c) **Encoder Model and Decoder Model:** let's understand how the notebook implements the expected log likelihood (first term on the right-hand side of Eq. 1).
- (i) The expectation $\mathbb{E}_{Q_\phi(\mathbf{Z})}[\dots]$ is estimated by sampling $\mathbf{z} \sim Q_\phi(\mathbf{Z})$, using the reparameterization trick discussed in Problem 8.2. Find the definition of the class `EncoderModel` in the notebook and make sure you understand its method `reparameterize`.
 - (ii) We use a fully factorized likelihood $P_\theta(\mathbf{X} | \mathbf{Z}=\mathbf{z}) = \prod_{i,j} P_\theta(X_{i,j} | \mathbf{Z}=\mathbf{z})$ where $P_\theta(X_{i,j}=1 | \mathbf{Z}=\mathbf{z}) = \sigma(\xi_{\theta,i,j}(\mathbf{z}))$ with the sigmoid function $\sigma(\alpha) := \frac{1}{1+e^{-\alpha}}$. Here, $\xi_{\theta,i,j}(\mathbf{z}) \in \mathbb{R}$ (called “logit”) is the (i,j) -th output of a neural network with input \mathbf{z} and weights θ . Show that $1 - \sigma(\alpha) = \sigma(-\alpha)$. Then look up the class `DecoderModel` in the notebook and make sure you understand the implementation of the method `log_likelihood`.
- (d) **Tying it all together:** find the definition of the function `bit_rates_and_logits` in the notebook and make sure you understand it. Identify the first two return values (`bit_rate_z` and `bit_rate_x_given_z`) with corresponding terms on the right-hand side of Eq. 1, then explain why we want to minimize their sum (as implemented in the function `train_one_epoch` in the next cell).

¹https://en.wikipedia.org/wiki/Kullback-Leibler_divergence#Multivariate_normal_distributions

Problem 9.2: Lossless Compression With a Variational Autoencoder and Bits-Back Coding

In this problem, you will use the toy Variational Autoencoder (VAE) that we implemented in the lecture to actually compress some data. Execute all cells in the section titled “Part 1” in the notebook to train the model. This takes about 10 minutes. While the model is training, familiarize yourself with Part 2 of the notebook.

- (a) **Problem Setup:** locate the definition of the function `test_compression` in the notebook. This function defines the task that we want to achieve. The argument `images` is a tensor that contains multiple images, which the function all encodes into a single bit string by calling `encode_single_image` for each image. After reporting the observed bit rate (in bits per pixel, BPP), the function decodes from the bit string by calling `decode_single_image` multiple times, and it verifies that all images were recovered without errors. Make sure you understand this setup.
- (b) **Bits-back Encoder:** The function `encode_single_image` is already implemented for you. It is annotated with three comments of the form “BITSBACK ENCODER STEP i : <description>” where $i \in \{1, 2, 3\}$. Remind yourself that these comments describe the three steps of bits-back coding (see also the table in Problem 7.1 (a)). Then read the function body and make sure you get the general idea. Don’t bother understanding all the transformations by `(un)scale_z(...)` and `(un)quantize_scaled_z(...)` yet. We’ll discuss this in Parts (d) and (e) below.
- (c) **Bits-back Decoder:** Now fill in the gaps marked “TODO” in the function `decode_single_image` in the cell below. Test your implementation by running `test_compression` and compare the empirical bit rates to the estimates that were reported during model training.

Hint: here’s how to encode/decode some symbols with the `constriction` library:

- if the model is already fully specified, such as `quantized_prior`:
`ans_coder.encode_reverse(symbols, quantized_prior)`
`symbols = ans_coder.decode(quantized_prior, count)`
- if the model has free parameters, such as `bernoulli` or `quantized_gaussian`:
`ans_coder.encode_reverse(symbols, model, params1, params2, ...)`
`symbols = ans_coder.decode(model, params1, params2, ...)`
This encodes/decodes `len(params1)` symbols, where, for the i^{th} symbol, the model parameters are `(params1[i], params2[i], ...)`.

- (d) **Quantizing latent space:** let’s now understand what the functions `(un)scale_z` and `(un)quantize_scaled_z` do. Two of the three bits-back steps encode or decode a latent representation `z`. However, in our VAE, `z` is a real-valued vector. We cannot losslessly compress arbitrary real values because \mathbb{R} is uncountable, i.e., there is no injective mapping from \mathbb{R} to the (countable) set of bit strings.

To work around this limitation, we approximate the random variable \mathbf{Z} by a discrete random variable $\hat{\mathbf{Z}}$ that we obtain by rounding each vector component of \mathbf{Z} to the nearest integer multiple of some fixed scalar `GRID_SPACING`. Thus, $\hat{\mathbf{Z}}$ takes only discrete values on an evenly spaced grid $\mathcal{G} := \{\hat{\mathbf{z}} : \hat{z}_i / \text{GRID_SPACING} \in \mathbb{Z} \forall i\}$. Convince yourself that, if $\hat{\mathbf{Z}}$ has PDF p under some model P , then, $\forall \hat{\mathbf{z}} \in \mathcal{G}$,

$$P(\hat{\mathbf{Z}} = \hat{\mathbf{z}}) = \int_{\mathcal{V}(\hat{\mathbf{z}})} p(\mathbf{z}) d\mathbf{z} \quad (3)$$

where $\mathcal{V}(\hat{\mathbf{z}}) := [\hat{\mathbf{z}} - \frac{1}{2} \times \text{GRID_SPACING}, \hat{\mathbf{z}} + \frac{1}{2} \times \text{GRID_SPACING}]$ is a cube of size `GRID_SPACING` centered at the grid point $\hat{\mathbf{z}}$.

Hint: Recall that random variables are defined as functions from some sample space Ω to some value space. Thus, if $\mathbf{Z} : \Omega \rightarrow \mathbb{R}^d$ for some dimension d , then $\hat{\mathbf{Z}} : \Omega \rightarrow \mathcal{G}$ with $\hat{\mathbf{Z}}(\omega) = \lceil \mathbf{Z}(\omega) \rceil_{\mathcal{G}}$ where $\lceil \cdot \rceil_{\mathcal{G}}$ denotes rounding to the nearest point in the grid \mathcal{G} . We defined the notation $P(\hat{\mathbf{Z}} = \hat{\mathbf{z}})$ as the probability of the event $E := \{\omega \in \Omega : \hat{\mathbf{Z}}(\omega) = \hat{\mathbf{z}}\}$. What can you say about $\mathbf{Z}(\omega)$ for all $\omega \in E$?

- (e) **Quantizing to integers:** the `constriction` library that we use for entropy coding here provides adapters that approximate arbitrary probability densities by quantizing according to Eq. 3. However, the library always quantizes to integers rather than to integer multiples of some given `GRID_SPACING`. This shouldn't be an issue though since we can simply define yet another random variable

$$\tilde{\mathbf{Z}} := \lceil (1/\text{GRID_SPACING}) \times \mathbf{Z} \rceil_{\mathbb{Z}^d} \quad (4)$$

which scales \mathbf{Z} by $1/\text{GRID_SPACING}$ before rounding each component to an integer.

The functions `scale_z` and `quantize_scaled_z` implement the scaling and rounding from Eq. 4, respectively. The functions `unscale_z` and `unquantize_scaled_z` implement the respective inverses (as far as inverting is possible). Read their definitions, then read the function `encode_single_image` again and make sure you understand why the (un)scaling and (un)quantizing is done at each point. Then explain why we also scale `q_mean` to `scaled_q_mean` and `q_std` to `scaled_q_std` at the beginning of `encode_single_image` (but we don't need to quantize these).

- (f) **Generalization Performance:** the last section of the notebook allows you to explore how both the VAE itself and your compression method generalizes to a different data set with images that have different dimensions than the training images. It is meant to demonstrate that fully convolutional model architectures naturally generalize to arbitrary image dimensions. In practice, VAEs for image or video compression are often trained on random patches of images for performance reasons, and then deployed on larger images.

Execute the cells in Part 3 of the notebook and enjoy the fruit of your labor.