5.4. Probabilistic autoencoders

An autoencoder is a pair of neural networks, an *encoder* which takes an input (for example an image) and compresses it down to a low-dimensional *latent representation*; and a *decoder* which takes the latent representation and expands back to something like the original.



What's the point in training a neural network just to reproduce its input? There are all sorts of ways that this could be useful. If we're lucky, the latent representation learns to pick up 'meaningful axes' for our dataset. There may be a lot of redundancy in each datapoint, for example most pixels in an image are much like their neighbours, and it's useful to get a meaningful summary with all the redundancy removed.

As a corollary, the autoencoder can fill in gaps. For example, if you see a picture of a stoat and part of it is missing, you can generally fill in the rest because you know what a stoat looks like—you're making use of the redundancy. Similarly it can remove noise:



We can generate novel synthetic datapoints by simply sampling a random point in the lowdimensional space (called the *latent space*) and feeding it into the decoder.



If we have a huge amount of unlabelled data and only a little bit of labelled data, we can first train the encoder and decoder on the full dataset, and then we can train a classifier which takes as its input the latent representation. The encoder has already done the hard work of learning the key features, and it's had a huge dataset to work on so it can do a good job, and this makes it easier to train the classifier.

All this is the unicorns and sparkles view of autoencoders. But neural networks aren't magic, they're just probability models that are trained using maximum likelihood estimation. Anything we want them to do, we have to fight for by setting up the right probability model. The probabilistic model behind autoencoding was laid out in two seminal papers³⁴. They are nothing more than generative neural networks, with a clever trick based on importance sampling for computing the log likelihood function.

generative neural networks: section 3.4 page 57

AUTOENCODERS AS PROBABILISTIC GENERATIVE MODELS

Suppose we have a dataset x_1, \ldots, x_n and we decide to model it as independent samples from a random variable X, where X is generated according to a latent variable model



where Z is some standard random variable, f_{θ} is a neural network, and X is a parametric random variable whose distribution depends on $f_{\theta}(Z)$. As a concrete illustration, perhaps the datapoints are in \mathbb{R}^{e} , and we decide to model them with a d-dimensional latent representation plus noise:

³⁴Diederik P. Kingma and Max Welling. "Auto-Encoding Variational Bayes". In: *ICLR*. 2014. URL: https://arxiv. org/abs/1312.6114; Danilo Jiminez Rezende, Shakir Mohamed, and Daan Wierstra. "Stochastic backpropagation and approximate inference in deep generative models". In: *ICML*. 2014. URL: https://arxiv.org/abs/1401.4082. The present description is a simplified account that owes more to Yuri Burda, Roger Grosse, and Ruslan Salakhutdinov. "Importance Weighted Autoencoders". In: *ICLR*. 2016. URL: https://arxiv.org/abs/1509.00519.

$$Z \sim N(0, I_d) \longrightarrow f_{\theta} \longrightarrow X \sim f_{\theta}(Z) + N(0, \rho^2 I_e)$$

(The noise parameter ρ might be known, or it might be another parameter to be trained along with θ .) This is exactly the generative neural we looked at in section 3.4, for generating points scattered around a path, where we used d = 1 and e = 2. If we want to generate MNIST images, we should let $e = 28 \times 28$ and let d be as small or as large as we like. It's modeller's choice: it's up to us to pick d, and to decide the number of layers and nodes etc. for f_{θ} , to come up with a useful model for the dataset in front of us.

Training. We'd like to fit this model in the usual way, by choosing θ to maximize the log likelihood of the dataset:

$$\log \Pr(x_1, \dots, x_n; \theta)$$

= $\sum_i \log \Pr_X(x_i; \theta)$ modelling the data as independent samples
= $\sum_i \log \int_z \Pr_X(x_i \mid Z = z; \theta) \Pr_Z(z) dz$ by law of total prob.

We could approximate this integral using Monte Carlo integration, and indeed that's exactly how we trained this model when we first looked at generative neural networks in section 3.4. But it can take lots of samples for the Monte Carlo approximation to be any good, especially when the latent variable Z has more than one dimension. There is some remarkably clever maths that gives a better way of training the model, and this maths is described in the next section.

For now, we'll simply state the conclusion. The probabilistic autoencoder consists of two networks, the encoder network g_{ϕ} and the decoder network f_{θ} . There are two separate configurations, one for generating values, the other for training.



- For generating, the latent variable Z is some standard random variable, $f_{\theta}(Z)$ computes the parameters for a distribution, and X is generated from that distribution.
- For training, $g_{\phi}(x)$ computes the parameters for a distribution, and \tilde{Z} is generated from that distribution—generated explicitly as a deterministic function of $g_{\phi}(x)$ and F where F is some standard random variable.

Training consists in finding θ and ϕ to maximize

$$\mathcal{L}_{\rm lb}(\theta,\phi) = \sum_{i} \Bigl\{ \Bigl[\mathbb{E} \log \Pr_{X} \bigl(x_{i} \mid \tilde{Z}^{(i,\phi)} \ ; \ \theta \bigr) \Bigr] - \mathrm{KL}(\Pr_{\tilde{Z}^{(i,\phi)}} \parallel \Pr_{Z}) \Bigr\}$$

where $\tilde{Z}^{(i,\phi)}$ is short for $\tilde{Z}(g_{\phi}(x_i), F)$, where the expectation is over F, where the expectation can be approximated using Monte Carlo, and where

$$\mathrm{KL}\big(\mathrm{Pr}_{\tilde{Z}}\,\|\,\mathrm{Pr}_{Z}\big) = \mathbb{E}_{z\sim \tilde{Z}^{(i,\phi)}}\log\Big(\frac{\mathrm{Pr}_{\tilde{Z}^{(i,\phi)}}(z)}{\mathrm{Pr}_{Z}(z)}\Big).$$

The KL term is called the Kullback-Leibler divergence from \Pr_Z to $\Pr_{\tilde{Z}^{(i,\phi)}}$, and it measures the difference between the two distributions.

To make all this concrete, let's pick some explicit distributions. Here is an autoencoder that we might use to model the MNIST images at the beginning of this section, a dataset in which each datapoint is a black-and-white image ${}^{35}x_i \in \{0,1\}^{28 \times 28}$.

³⁵Technically the images are greyscale, $x_i \in [0, 1]^{28 \times 28}$. The log likelihood function in this code technically corresponds to a continuous analogue of the Bernoulli random variable. But the distribution of pixels is heavily concentrated—most are 0 or 1, and so there's little harm in just pretending they're Bernoulli. A more refined model might take account of the remaining 8% or so of pixels which are roughly uniform in [0, 1]. But the real deficiency of this model is that it treats the pixels as independent; it's this deficiency that leads to the blurred outputs of the autoencoder.

Example 5.4.1 (Bernoulli outputs, Gaussian latent variables).

Consider a dataset x_1, \ldots, x_n with datapoints $x_i \in \{0, 1\}^e$. We wish to model it using a Bernoulli generative model:

def rx(): $z = np.random.normal(size=d) # vector in \mathbb{R}^d$ $p = f(z) # vector in [0, 1]^e$ return np.random.binom(1,p) # vector in $\{0, 1\}^e$

$$Z \sim N(0, I_d) \longrightarrow f_{\theta} \longrightarrow X \sim \operatorname{Binom}(1, f_{\theta}(Z))$$

Here the neural network is a function $f_{\theta} : \mathbb{R}^d \mapsto [0, 1]^e$, and each component of $X \in \{0, 1\}^e$ is generated independently from the corresponding component of $f_{\theta}(Z)$.

Explain how to fit this model, using the following Gaussian encoder:

$$x \xrightarrow{g_{\phi}} \mu, \sigma \xrightarrow{\tilde{Z}} \mu + \sigma F$$
$$F \sim N(0, I_d)$$

where the neural network $g_{\phi}(x)$ outputs two vectors, $\mu \in \mathbb{R}^d$ and $\sigma \in \mathbb{R}^d_{\geq 0}$, and where the expression for \tilde{Z} is componentwise addition and multiplication.

First let's implement the generator and its log likelihood function. The conditional log likelihood is

$$\log \Pr_X(x \mid Z = z) = \sum_{k=1}^{e} \left[x_k \log p_k + (1 - x_k) \log(1 - p_k) \right] \quad \text{where} \quad p = f_{\theta}(z).$$

Here's code. The neural network f and the log likelihood function have been written to work on batches of datapoints, B of them per batch, since that's the PyTorch convention.

```
class BernoulliImageGenerator(nn.Module):
1
2
           def ____init___(self, d=4):
                super().\_init\_()
3
                 self.d = d
4
                 self. f = \ldots \# f : \mathbb{R}^{B \times d} \to [0, 1]^{B \times e}
5
6
          def forward(self, z):
7
8
                return self.f(z)
9
          def loglik(self, x, z): \#x \in \{0,1\}^{B 	imes e} and z \in \mathbb{R}^{B 	imes d}
10
                xr = self(z)
11
                \texttt{return} (\texttt{x*torch.log}(\texttt{xr}) + (1-\texttt{x})\texttt{*torch.log}(1-\texttt{xr})).\texttt{sum}(1) \ \# \in \mathbb{R}^B
12
```

Next, the encoder. The autoencoder training objective \mathcal{L}_{lb} is implemented on line 25. There are two terms in the expression for \mathcal{L}_{lb} , an expected log likelihood term and a KL term. The expected log likelihood is computed on lines 27–28, using Monte Carlo with just a single sample (!).

```
13 class GaussianEncoder(nn.Module):
14
           def __init__(self, decoder):
15
                 super().__init__()
                 self.d = decoder.d
16
                 self.f = decoder
17
                 \mathsf{self.g} = \dots \# g : \{0,1\}^{B \times e} \to \mathbb{R}^{B \times 2d}
18
19
           def forward(self, x):
20
                 \mu \tau = self.g(x)
21
                 \mu,\tau = \mu\tau[:,:\mathsf{self.d}], \ \mu\tau[:,\mathsf{self.d}:]
22
                 return \mu, torch.exp(\tau/2) \mu \in \mathbb{R}^{B \times d}, \sigma \in \mathbb{R}^{B \times d}_{>0}
23
24
           def loglik_lb(self, x): \#x \in \{0,1\}^{B \times e}
25
                 \mu,\sigma = self(x)
26
                 \varepsilon = \text{torch.randn}_{\text{like}}(\sigma)
27
                 II = self.f.loglik(x, z = \mu + \sigma * \varepsilon)
28
```

kl = 0.5 * (
$$\mu$$
**2 + σ **2 - torch.log(σ **2) - 1).sum(1)

30 return $|| - \mathbf{k}| \quad \# \in \mathbb{R}^B$

As for the KL term, the distributions we're using in this model are so simple that it's easy to calculate an exact formula, in line 29. Here's the derivation. We've chosen to let $\tilde{Z}^{(i,\phi)}$ consist of d independent Normal (μ_k, σ_k^2) random variables, $k = 1, \ldots, d$, and to let Z consist of d independent N(0, 1) random variables. (The concise notation for this is $\tilde{Z} \sim N(\mu, \sigma^2 I_d)$ and $Z \sim N(0, I_d)$.) The KL divergence KL($\Pr_{\tilde{Z}} \parallel \Pr_Z$) is

$$\begin{split} KL(\Pr_{\tilde{Z}} \| \Pr_{Z}) &= \mathbb{E}_{z \sim \tilde{Z}} \log \left(\frac{\Pr_{\tilde{Z}}(z)}{\Pr_{Z}(z)} \right) \\ &= \mathbb{E}_{z \sim \tilde{Z}} \log \prod_{k=1}^{d} \frac{\frac{1}{\sqrt{2\pi\sigma_{k}^{2}}} e^{-(z_{k}-\mu_{k})^{2}/2\sigma_{k}^{2}}}{\frac{1}{\sqrt{2\pi}} e^{-z_{k}^{2}/2}} \\ &= \mathbb{E}_{z \sim \tilde{Z}} \sum_{k=1}^{d} \left(-\frac{1}{2} \log \sigma_{k}^{2} - \frac{(z_{k}-\mu_{k})^{2}}{2\sigma_{k}^{2}} + \frac{z_{k}^{2}}{2} \right) \\ &= \frac{1}{2} \sum_{k=1}^{d} \left(\sigma_{k}^{2} + \mu_{k}^{2} - \log \sigma_{k}^{2} - 1 \right). \end{split}$$

Finally, we can train it in the standard PyTorch way. The objective is to maximize \mathcal{L}_{lb}

PyTorch optimization: see section 3.3 page 55

31 model = GaussianEncoder(BernoulliImageGenerator(d=4))optimizer = optim.Adam(model.parameters()) 32 33 for epoch in range(5): 34 for batch num, (imgs, lbls) in enumerate (mnist batched): 35 36 optimizer.zero_grad() 37 loglik_lb = torch.mean(model.loglik_lb(imgs)) 38 (-loglik_lb).backward() 39 optimizer.step()

WHY DOES IT EVEN WORK?

Where on earth does the objective function for an autoencoder come from? Let's look at the terms that make it up.

$$\mathcal{L}_{\rm lb}(\theta,\phi) = \sum_{i} \left\{ \left[\mathbb{E} \log \Pr_{X} \left(x_{i} \mid \tilde{Z}^{(i,\phi)} \; ; \; \theta \right) \right] - \mathrm{KL}(\Pr_{\tilde{Z}^{(i,\phi)}} \parallel \Pr_{Z}) \right\}$$

The first term is a likelihood term. It will be high when the latent variable \hat{Z} makes x_i more likely. This will happen when the encoder and decoder are matched: when the encoder maps x_i onto \tilde{Z} in some small part of the latent space, and when the decoder maps that specific part of the latent space to x_i .

Why does there need to be randomness in the encoder? This is to encourage the autoencoder to learn 'meaningful representations', in the sense that we'd like to be able to tweak components of the latent representation z, and get slightly tweaked outputs $f_{\theta}(z)$, not radically different outputs. If the decoder were discontinuous, then depending on the exact value that the random variable \tilde{Z} happens to take we might end up with a high likelihood $\Pr_X(x_i|\tilde{Z};\theta)$ or we might end up with a low likelihood. The only way the autoencoder can get a reliably high likelihood is if it learns a reasonably continuous decoder.

The KL term is a penalty for badly-distributed \tilde{Z} . Suppose the autoencoder learnt to put all the encoded values into weird disconnected parts of the latent space. This wouldn't be any good at generating new values. We want to be able to synthesize new values by sampling a random Z and feeding it into the generator—and so we need to penalize the network if it doesn't use all parts of the latent space. That's what the KL term does.

It's useful to have intuition about what the terms do, especially when it comes to debugging; but this is all arguing by analogy, and it's weak argument. Jesus taught in parables because he taught ineffable wisdom; but in machine learning we should stick to describing what things *are* not what they are *like*. It'd be almost impossible to *derive* the autoencoder from hand-waving intuition like this.

The autoencoder is a beautiful and precise balance, and its only real justification is the probabilistic derivation in section 5.5. If we want to modify the setup to achieve different goals, we should go via the probability derivation and not via intuition, else we're liable to mess things up.