

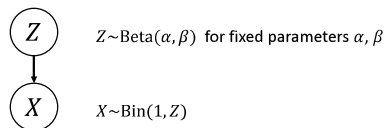
R250 Advanced topics in machine learning Topic 5: autoencoders

## Variational autoencoder: the foundations

Damon Wischik, Cambridge University

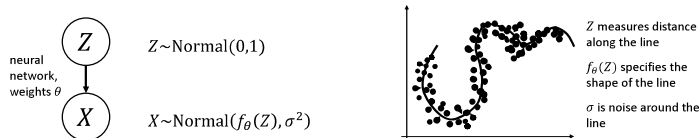
- Kingma and Welling, *Autoencoding variational Bayes*, ICLR 2014
- Burda, Grosse, Salakhutdinov, *Importance weighted autoencoders*, ICLR 2016

A Bayesian's favourite model, of a biased coin



The Bayesian likes to ask:  
given observations  $x_1, x_2, \dots, x_N$  drawn from distribution  $X$ ,  
what is the posterior distribution of  $Z$ ?

A latent-random-variable generative model



The machine-learning modeller might ask:  
given observations  $x_1, x_2, \dots, x_N$  drawn from distribution  $X$ ,  
how can I tune  $\theta$  so that the model fits the data?

What does "train a generative model" mean?

It means "Pick  $\theta$  so that the distribution  $\text{Pr}_X(x|\theta)$  is a good fit for the dataset  $x_1, \dots, x_N$ ".

maximize  $\log \text{lik}(\theta)$  over  $\theta$   
 where  $\log \text{lik}(\theta) = \frac{1}{N} \sum_{i=1}^N \log \text{Pr}_X(x_i|\theta)$

This also provides the evaluation metric,

$\frac{1}{M} \sum_{i=1}^M \log \text{Pr}_X(x_i|\hat{\theta})$  summed over the holdout dataset  $x_1, \dots, x_M$ , using the fitted  $\hat{\theta}$

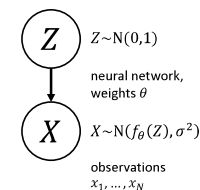
How to train a generative model using Monte Carlo

In toy examples, we can write down a formula for the log likelihood  $\log \text{Pr}_X(x|\theta)$ .

For interesting neural networks, this expression is intractable, so we approximate.

$$\begin{aligned} \log \text{lik}(\theta) &= \frac{1}{N} \sum_{i=1}^N \log \text{Pr}_X(x_i|\theta) \\ &\stackrel{\text{Law of total probability:}}{\text{Pr}_X(x) = \int_z \text{Pr}_X(x|Z=z) \text{Pr}_Z(z) dz = \mathbb{E}_Z \text{Pr}_X(x|Z)} \\ &= \frac{1}{N} \sum_{i=1}^N \log \mathbb{E}_Z \text{Pr}_X(x_i|Z, \theta) \\ &\stackrel{\text{Jensen's inequality, for concave functions:}}{\text{f}(\mathbb{E}X) \geq \mathbb{E}f(X)} \\ &\geq \frac{1}{N} \sum_{i=1}^N \mathbb{E}_Z \log \text{Pr}_X(x_i|Z, \theta) \\ &\stackrel{\text{Monte Carlo approximation, where } z_j \text{ are sampled from } Z}{\approx} \frac{1}{N} \sum_{i=1}^N \frac{1}{J} \sum_{j=1}^J \log \text{Pr}_X(x_i|Z = z_j, \theta) \end{aligned}$$

Our model gives us the explicit formula.  
Gradient descent does the rest.



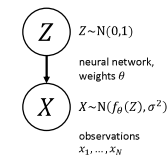
How to train a generative model using Monte Carlo

Our training objective: find  $\theta$  to maximize the lower bound

$\log \text{lik lb}(\theta) = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_Z \{\log \text{Pr}_X(x_i|Z, \theta)\}$

For each datapoint  $x_i$  (or each batch):

- Generate one or more random  $z$  samples from  $\text{Pr}_Z$
- Compute the loss function,  $L(\theta) = -\log \text{Pr}_X(x_i|Z = z, \theta)$  as well as its gradient  $dL/d\theta$
- Update  $\theta$  to reduce the loss function



PROBLEM: the Monte Carlo approximation is pretty terrible, for this problem, since most values of  $z$  give  $\text{Pr}_X(x_i|Z = z, \theta) \approx 0$ . This makes the lower bound very weak.

### Digression on Importance Sampling

Given a random distribution  $Z$  and a function  $h$ , how can we approximate  $\mathbb{E}_Z h(Z)$  ?

#### MONTE CARLO APPROXIMATION

Sample  $z_1, \dots, z_j$  from  $Z$ . Then  $\mathbb{E}_Z h(Z) \approx \frac{1}{j} \sum_1^j h(z_j)$

#### IMPORTANCE SAMPLING APPROXIMATION

Choose a distribution  $\tilde{Z}$ . Sample  $z_1, \dots, z_j$  from  $\tilde{Z}$ . Then  $\mathbb{E}_Z h(Z) \approx \frac{1}{j} \sum_1^j h(z_j) \frac{\Pr_Z(z_j)}{\Pr_{\tilde{Z}}(z_j)}$

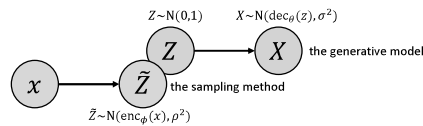
This approximation is valid for any distribution  $\tilde{Z}$ . It works best (i.e. is good for small  $j$ ) if  $\tilde{Z}$  is biased in favour of values where  $h(z)$  is large.

In our generative model, we picked  $Z \sim N(0,1)$  and want to compute

$$\mathbb{E}_Z \log \Pr_X(x|Z, \theta) = h(Z)$$

Let's sample instead from some other distribution  $\tilde{Z}$  (our choice). We just have to throw in a correction factor.

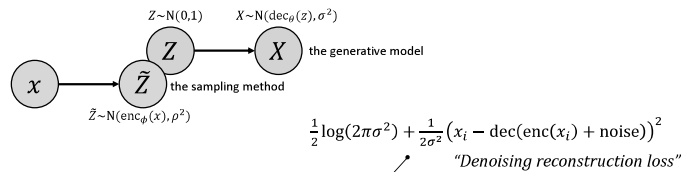
### How to train a generative model using importance sampling



For each datapoint  $x_i$  (or each batch):

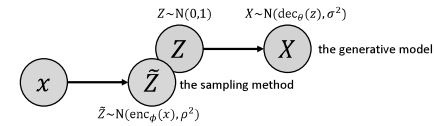
Generate one or more random  $z$  samples from  $\Pr_Z$  and compute the loss,  $L = -\log \left\{ \Pr_X(x_i | Z = z, \theta) \frac{\Pr_Z(z)}{\Pr_Z(z|\phi)} as well as the gradients  $\partial/\partial\theta$  and  $\partial/\partial\phi$ . Update  $(\theta, \phi)$  to reduce the loss$

### Our loss function can be thought of in terms of reconstruction error



Generate one or more random  $z$  samples from  $\Pr_Z$  and compute the loss,  $L = -\log \left\{ \Pr_X(x_i | Z = z, \theta) \frac{\Pr_Z(z)}{\Pr_Z(z|\phi)}$

### Backpropagation needs derivatives, and there's a trick needed ...



Generate one or more random  $z$  samples from  $\Pr_Z$

and compute the loss,  $L = \dots$

as well as the gradients  $\partial/\partial\theta$  and  $\partial/\partial\phi$

But  $\Pr_Z$  depends on  $\phi$ . How do we differentiate "sample from  $\Pr_Z$ "?

### Some neat maths ("variational inference")

$$\begin{aligned} \log \Pr_X(x|\theta) &= \log \mathbb{E}_Z \Pr_X(x|Z, \theta) \\ &= \log \mathbb{E}_{Z \sim \Pr_Z} \left\{ \Pr_X(x|Z, \theta) \frac{\Pr_Z(Z)}{\Pr_Z(Z)} \right\} \\ \text{With a little bit of algebra, the error in this approximation is } & \text{KL}(\Pr_Z \parallel \Pr_{Z|X=x}) \\ & \geq \mathbb{E}_{Z \sim \Pr_Z} \log \left\{ \Pr_X(x|Z, \theta) \frac{\Pr_Z(Z)}{\Pr_Z(Z)} \right\} \\ \text{So, if our encoder is a perfect match for the Bayesian posterior of } Z \text{ given } & X = x, \text{ the error is zero. (And also importance sampling is very efficient.)} \\ &= \mathbb{E}_{Z \sim \Pr_Z} \log \Pr_X(x|Z, \theta) - \underbrace{\mathbb{E}_{Z \sim \Pr_Z} \log \frac{\Pr_Z(Z)}{\Pr_Z(Z)}}_{\text{"denoising reconstruction loss"}} \\ & \qquad \underbrace{\log \frac{\Pr_Z(Z)}{\Pr_Z(Z)}}_{\text{The Kullback-Leibler divergence, KL}(\Pr_Z \parallel \Pr_Z) \geq 0, \text{ which is behaving here as a regularizer}} \end{aligned}$$