AutoEncoder (AE): High Level

TL;DR

- The Deep Learning analog of Principal Components (PCA)
 - Most of the lessons of AE apply equally to PCA
- Unsupervised: no labels (really semi-supervised)
- Create "synthetic features" from the original set of features
- May be able to use reduced set of synthetic features (dimensionality reduction)
- Generative (vs Discriminative)

Autoencoder

An Autoencoder network has two parts

- An Encoder, which takes input x and "encodes" it into z
- A Decoder, which takes the encoding z and tries to reproduce x

Each part has its own weights, which can be discovered through training, with examples

•
$$\langle X, y \rangle = \langle X, X \rangle$$

That is: we are asking the output to be identical to the input.

z is an alternative latent representation of x.

- Encoded by the Encoder
- Inverted by the Decoder

But when the dimension of z is less than the dimension of x.

- z is a bottle-neck
- the inversion by the Decoder will be imperfect

z becomes a reduced-dimensionality approximation of x.

This is reminiscent of the dimensionality reduction of Principal Components Analysis (PCA).

The main difference from PCA

- PCA uses a linear transformation
- NN can use non-linear transformations too
 - PCA as a special case of AE

Autoencoders: Uses

Dimension reduction

After training

- we can discard the Decoder
- use the Encoder output (synthetic features) as reduced dimension inputs to a new task
 - Encoder weights are frozen: non-learnable when training new task
 - \circ It may be easier to solve the new task given ${f z}$ rather than ${f x}$
 - have already discovered "structure" of x
 - Transfer Learning

Autoencoder: Encoder + New head

In PCA, we eliminated original features that were "less important"

- i.e., explained variation among only a small fraction of the training set
 - recall how we re-denominated explained variance in terms of "number of features"

There is no direct similar concept of feature importance in AE

• other than minimizing a Loss function, which *may* wind up focusing on "important" features

Layer-wise pre-training with Autoencoders

Autoencoders played a vital role in the development of Deep Learning:

- They made it possible to train otherwise untrainable NN's.
- Other innovations supplanted the need for AE's to assist training
 - better initialization
 - better activations functions
 - normalization

Although they are no longer needed for that purpose, it is interesting to see how (and why) they were used.

For a NN with $oldsymbol{L}$ layers that solves a Supervised Learning Problem

- Training attempts to learn the weights of all layers simultaneously
- Layer wise pre-training was an attempt
 - to initialize the weights of each layer
 - in succession
 - so that the task of simultaneously solving for optimal weights had a better chance of succeeding

The idea was to learn an initialization of $\mathbf{W}_{(l)}$, the weights of layer l.

• After having learned the weights $\mathbf{W}_{(l')}$ for all layers l' < l.

To initialize $\mathbf{W}_{(l)}$:

- Train an AE that takes $\mathbf{x^{(i)}}$ as input
- ullet Using initialized weights $\mathbf{W}_{(l')}$ for all layers l' < l
- Produces $\mathbf{\tilde{x}^{(i)}}$ at layer l's output $\mathbf{y}_{(l)}$

So weight initializations were learned layer by layer.

Note that the labels $y^{(i)}$ were not used!

wouldn't be useful for the shallow NN

It was thought

- ullet to be easier to learn the structure of the input ${f x}$ independent of the labels
- ullet to be easier to learn $\mathbf{W}_{(l)}$ incrementally

One the weights $\mathbf{W}_{(l)}$ were initialized via AE's

- training of the Supervised Learning task had a better chance of succeeding
- compared to any other initialization

Autoencoders and Transfer Learning

Today, autoencoders are useful for another purpose: Transfer Learning.

If we can train an AE network to create features that are useful for reconstruction

• it is possible that these features are useful for solving more complicated tasks.

This was in essence what

- Our dimension reduction example (replace the head) was doing
- Layerwise Pre-training was attempting.

| (and smaller) |
|---------------|
| |
| |
| |

Note that Autoencoders are unsupervised: they don't take labels.

So the encodings they produce stress syntactic similarity, rather than semantic similarity.

Their use in Transfer Learning depends on the hope that inputs that are syntactically similar also have the same labels.

Denoising

Very much like dimension reduction but with the assumption that

• "less important" features are just random noise that is added to the true example

Autoencoder: Denoising

Generative Artificial Intelligence

A less obvious use of AE (using the Decoder rather than the Encoder) is to generate examples.

Most of the Machine Learning we have studied thus far is discriminative

- $p(\hat{\mathbf{y}}^{(i)}|\mathbf{x}^{(i)})$
 - \blacksquare e.g., classifier: discriminate among the possible classes $y^{(i)}$, given example $\mathbf{x}^{(i)}$

We can use the Decoder on arbitrary z to generate a completely new x:

- $p(\mathbf{x}^{(i')}|\mathbf{z}^{(i')})$ for some i' not in training
- generate a new example $m{i'}$, in the domain of $m{x}$, that was not encountered during training

Generator

Autoencoder (AE): Details

The task that trains an Autoencoder

- Given input $\mathbf{x}^{(i)}$
- Output of Encoder: $\mathbf{z^{(i)}} = E(\mathbf{x^{(i)}})$
- Output of Decoder: $\mathbf{ ilde{x}^{(i)}} = D(\mathbf{z^{(i)}})$
- "Target": **x**⁽ⁱ⁾

Both the Encoder and Decoder are parameterized (learnable parameters)

• Goal: find the parameters such that

$$\tilde{\mathbf{x}}^{(i)} = D(E(\mathbf{x})) \approx \mathbf{x}$$

 $\mathbf{z^{(i)}} = E(\mathbf{x^{(i)}})$ is the latent representation of $\mathbf{x^{(i)}}$.

Loss function

The obvious loss functions compare the original $\mathbf{x^{(i)}}$ and reconstructed $\mathbf{\tilde{x}^{(i)}}$ feature by feature:

Mean Squared Error (MSE)

$$\mathcal{L}^{(\mathrm{i})} = \sum_{j=1}^{|\mathrm{x}|} (\mathrm{x}_j^{(\mathrm{i})} - ilde{\mathrm{x}}_j^{(\mathrm{i})})^2$$

Binary Cross Entropy

For the special case where each original feature is in the range [0,1] (e.g., an image)

$$\mathcal{L}^{(\mathrm{i})} = \sum_{j=1}^{|\mathrm{x}|} \left(\mathrm{x}_j^{(\mathrm{i})} \log(ilde{\mathrm{x}}_j^{(\mathrm{i})}) + (1-\mathrm{x}_j^{(\mathrm{i})}) \log(1- ilde{\mathrm{x}}_j^{(\mathrm{i})})
ight)$$

Variational Autoencoder (VAE): Generative ML

Observe that the Decoder part of the "vanilla" AE $D(\mathbf{z^{(i)}})$

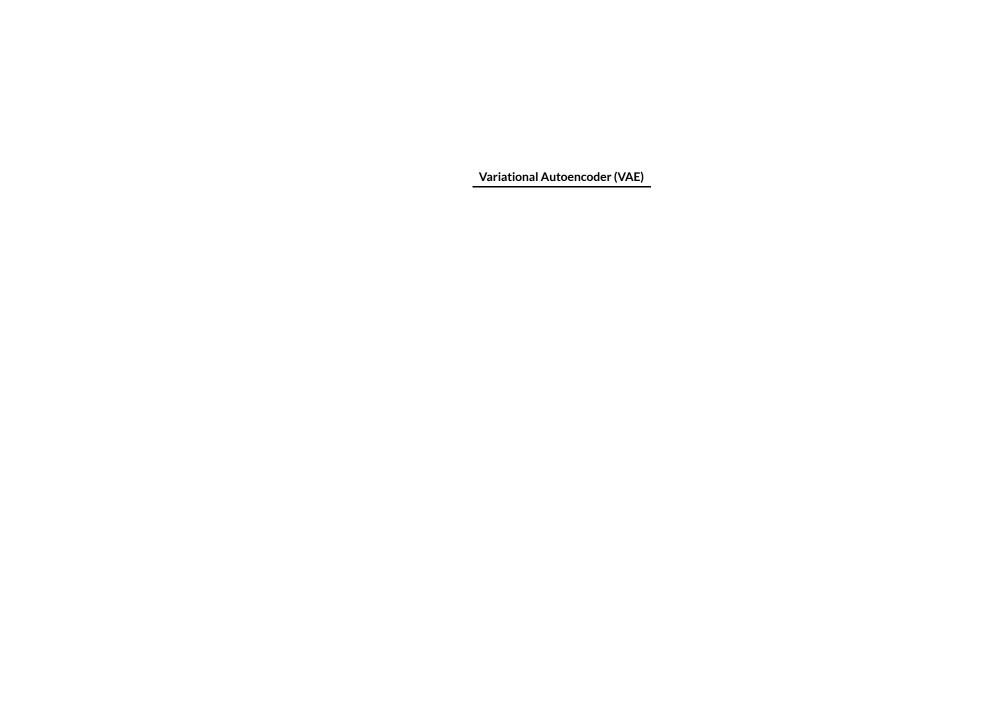
- has been trained to produce "realistic" $\mathbf{ ilde{x}^{(i)}}$ only for a $\mathbf{z^{(i)}} = E(\mathbf{x^{(i)}})$
 - lacktriangledown i.e., "realistic": appears to come from the distribution of training ${f X}$
- there is no guarantee that $D(\mathbf{z}^{(i')})$ for some i' not in training is realistic

That is: the AE has not been trained to extrapolate beyond the training inputs.

A VAE is able generate outputs

- that could have come from the training distribution from a latent representation $\mathbf{z}^{(i')}$
- but that did not come from X.

Our goal is constructing a Decoder that can extrapolate.



The Decoder will take a *latent vector* ${f z}$ and produce $D({f z})$, just as in a vanilla AE.

The difference is that z will be sampled from a distribution rather than being a unique mapping of a training example.

This will be done by modifying the Encoder

- It will indirectly create $z^{(i)}$
- It will compute variables $oldsymbol{\mu}^{(\mathrm{i})}$ and $oldsymbol{\sigma}^{(\mathrm{i})}$
 - $\mathbf{z^{(i)}}$ will be *sampled* from a distribution with mean $\boldsymbol{\mu^{(i)}}$ and standard deviations $\boldsymbol{\sigma^{(i)}}$

As long as ${\bf z}$ is sampled from this distribution, the decoder will produce a "realistic" output.

Note

 μ and σ are

- vectors
- computed values (and hence, functions of x) and not parameters so training learns a function from $x^{(i)}$ to $\mu^{(i)}$ and $\sigma^{(i)}$

To train a VAE:

- ullet pass input $\mathbf{x^{(i)}}$ through the Encoder, producing $oldsymbol{\mu^{(i)}, \sigma^{(i)}}$
 - use $\mu^{(i)}, \sigma^{(i)}$ to sample a latent representation $\mathbf{z}^{(i)}$ from the distribution
- ullet pass the sampled $\mathbf{z^{(i)}}$ through the decoder, producing $D(\mathbf{z^{(i)}})$
- ullet measure the reconstruction error $\mathbf{x^{(i)}} D(\mathbf{z^{(i)}})$, just as in a vanilla AE
- ullet back propagate the error, updating all weights and μ,σ

Essentially, each input $\mathbf{x^{(i)}}$ has many latent representations (with different probabilities): any sample from the distribution.

Training

Encoder produces

$$E(\mathbf{x}) = q_{\phi}(\mathbf{z}|\mathbf{x}) pprox p_{ heta}(\mathbf{z}|\mathbf{x})$$

We sample from

$$\hat{\mathbf{z}} \sim q_{\phi}(\mathbf{z}|\mathbf{x})$$

Decoder produces

$$D(\hat{\mathbf{z}}) = p_{\theta}(\mathbf{x}|\mathbf{z})$$

Each time (epoch) that we encounter the same training example, we select another random element from the distribution.

So the VAE learns to represent the same example from multiple latents.

Generative

- sample $\hat{\mathbf{z}} \sim \hat{p}(\mathbf{z})$
- use decoder to produce output $p_{ heta}(\mathbf{x}|\mathbf{z})$

This means we can feed in a z

- that doesn't correspond to any training example
- and perhaps get an output that *resembles* something from the training set, rather than noise.

Which came first: the VAE architecture or the loss function?

So far, we haven't told you the Loss function that is optimized during training.

It's a bit complicated so we'll save it for the end (for those who are interested).

The really interesting thing:

- The Loss function drove the architecture of the VAE, not vice-versa!
- As we've said many times in this course:
 - Deep Learning is about creating Loss functions that reflect our objectives

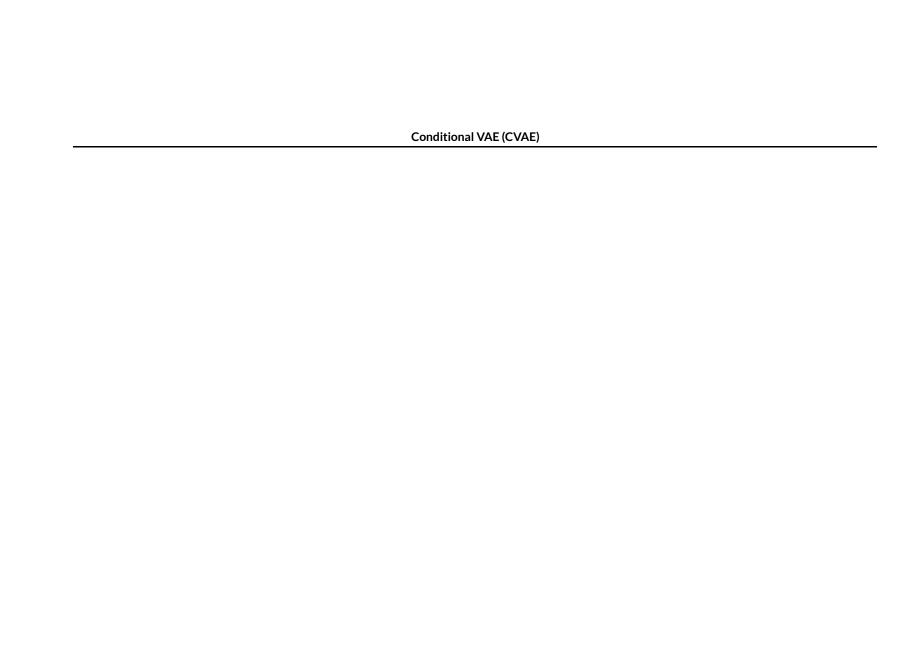
Conditional VAE

Once a VAE is trained, D(z) should produce a realistic output, for any z from the distribution.

However, if the distribution of X includes examples from many classes

- Assuming we have labels as auxilliary information (not used in training)
 - e.g., the 10 digits
- The VAE can't control which class the output will come from

A Conditional VAE allow our generator (Decoder) to control the class c of the output $\tilde{\mathbf{x}}$.



The class label c

- is given as part of training
 - So the Encoder produces a distribution that is conditioned on both ${\bf x}$ and ${\bf c}$.
- is an additional parameter of the Decoder
 - So the output class can be controlled

$$ilde{\mathbf{x}}^{(\mathrm{i})} = D(\mathbf{z}^{(\mathrm{i})}, c)$$

So now we

- create a latent z
- append a class label c
- and presumably have the decoder produce an output from the desired class.
- The encoding distribution is now conditional on class label \emph{c} : $q_{\phi}(z|x,c)$
- So is the decoding distribution $p_{ heta}(x|z,c)$

Again, by restricting the functional form of the prior distribution \hat{p} we can simplify the math.

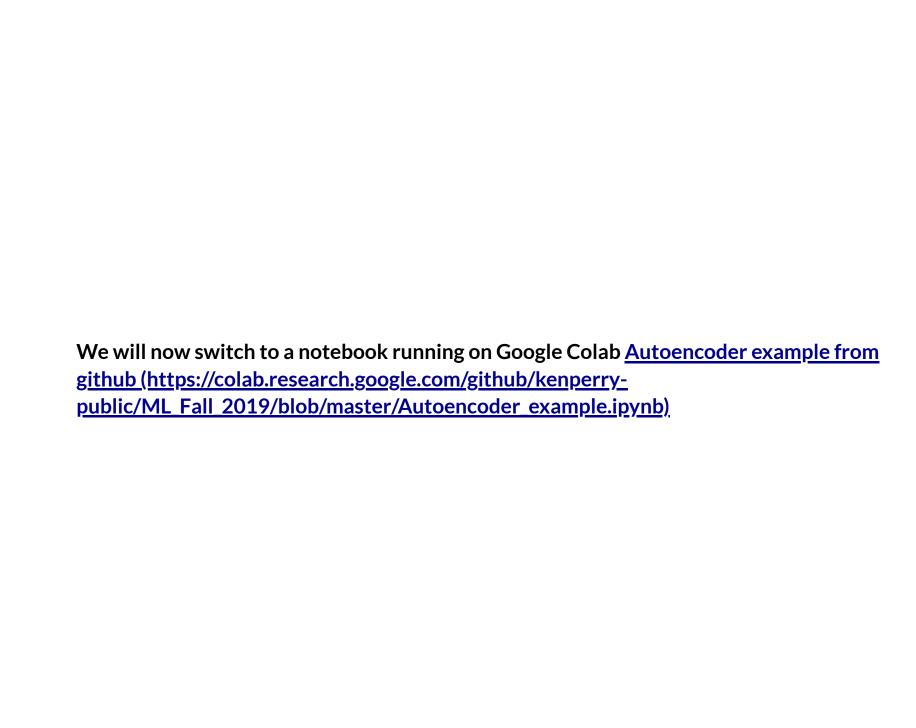
Detour: Autoencoder notebook on Colab

Let's examine some Keras code that implements several types of Autoencoders

- Vanilla
- Denoising
- VAE

We will write our AE's using the Keras Functional API rather than the Sequential model

- We could write the complete AE using the Sequential API
- But
- we want to extract the Encoder and Decoder parts as *separate models*
- we can do this with the Functional API



VAE: Probabilistic formulation

Note: Advanced material

Let's pretend: we don't already know that we will represent z as a function of $\mu_{ heta}(x)$ and $\sigma_{ heta}(x)$

this derivation will show why we made that choice

The mathematical derivation of a VAE is quite detailed

- it is interesting but not absolutely necessary to understand
- this is where we define the Loss function

The interested reader is referred to a highly recommended <u>VAE tutorial</u> (<u>https://arxiv.org/pdf/1606.05908.pdf</u>).

We will try to give the essence in the following slides.

TL;DR

- The VAE has a very interesting two part Loss Function
 - Reconstruction Loss, as in the Vanilla AE
 - Divergence Loss
- The Reconstruction Loss is not sufficient
 - Issues of intractability arise
 - The Divergence Loss skirts intractability
 - By constraining the Encoder to produce a tractable distribution

From the description of the VAE, observe that we are now dealing with distributions rather than deterministic values for

- the encoding (latent representation) **z**
- the output

So we will need to describe these distributions

We posit a joint probability distribution $p(\mathbf{x}, \mathbf{z})$ of examples and latents.

These are *empirical* distributions:

- they are defined by the data
- no closed form

We will approximate p, as usual, with a NN that we will parameterize with θ .

That is: we will train a model to learn θ .

So henceforth we write $p_{ heta}$ to denote the dependence of p on parameters heta.

We motivate these distributions as they relate to the VAE:

- ullet the encoder produces $p_{ heta}(\mathbf{z}|\mathbf{x})$
- the decoder produces $p_{ heta}(\mathbf{x}|\mathbf{z})$

Note that there are many reconstructions $\mathbf{\tilde{x}^{(i)}}$ of $\mathbf{x^{(i)}}$

ullet depending on the sampled $ilde{\mathbf{z}}^{(\mathrm{i})}$ drawn from $p_{ heta}(\mathbf{z}|\mathbf{z})$

The loss function: first attempt

Let's try to create an optimization objective function against which we choose our model weights.

We will use Maximum Likelihood as the objective

• Given the weights: how likely is the model to produce the training distribution?

Since our practice is to minimize Loss (rather than maximize an objective function) we write our loss as (negative of log) likelihood

$$\mathcal{L} = -\log(p_{\theta}(\mathbf{x}))$$

Minimizing ${\cal L}$ is equivalent to maximizing likelihood. T

Intractability

It turns out that things are not so simple:

- Some of the distributions we need to deal with may not be tractable
- They have no closed form, just empirical distributions
- Higher dimensional distributions may pose computational issues

The problem is that $p_{ heta}(\mathbf{z}|\mathbf{x})$ is assumed to be intractable

That is: for a given $\mathbf{x^{(i)}}$, we don't know which samples from \mathbf{z} can reconstruct $\mathbf{x^{(i)}}$.

$$p_{ heta}(\mathbf{z}|\mathbf{x}) = rac{p_{ heta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p_{ heta}(\mathbf{x})}$$
 (by Bayes rule)

- this represents the distribution for the encoder
- z is from an unknown distribution

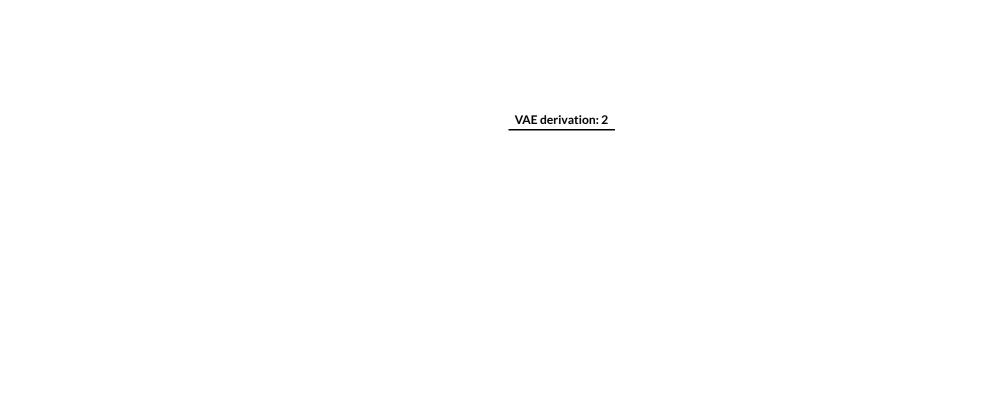
The solution is *approximate* the intractable $p_{ heta}(\mathbf{z}|\mathbf{x})$ with a tractable $q_{\phi}(\mathbf{z}|\mathbf{x})$

• where $q_{\phi}(\mathbf{z}|\mathbf{x})$ is a function of learnable parameters ϕ .

We use KL divergence as a measure of the difference between two distributions.

So we want to minimize

$$ext{KL}(q_{\phi}(ext{z}| ext{x}) \mid\mid p_{ heta}(ext{z}| ext{x}))$$



We add the KL divergence to our loss function

$$egin{array}{lll} \mathcal{L} &=& -\log(p_{ heta}(\mathbf{x})) + \mathrm{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \mid\mid p_{ heta}(\mathbf{z}|\mathbf{x})) \ &=& \mathcal{L}_R + \mathcal{L}_D \end{array}$$

which now has two objectives

- Reconstruction loss \mathcal{L}_R : maximize the likelihood (by minimizing the negative likelihood)
- Divergence constraint \mathcal{L}_D : $q_\phi(\mathbf{z}|\mathbf{x})$ must be close to $p_ heta(\mathbf{z}|\mathbf{x})$

$$egin{array}{lll} \mathcal{L}_R &=& -\log(p_{ heta}(\mathbf{x})) \ \mathcal{L}_D &=& \mathrm{KL}\left(q_{\phi}(\mathbf{z}|\mathbf{x}) \mid\mid p_{ heta}(\mathbf{z}|\mathbf{x})
ight) \end{array}$$

We will show (in the next section: lots of algebra!) that

$$\mathcal{L} = -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{x})} \left(\log(p_{ heta}(\mathbf{x}|\mathbf{z}))
ight) + \mathrm{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \mid\mid p_{ heta}(\mathbf{z}))$$

This is *almost* identical to \mathcal{L} except

- Re-write $\log(p_{ heta}(ext{x})) = \mathbb{E}_{z \sim q_{\phi}(ext{z}| ext{x})}\left(\log(p_{ heta}(ext{x}| ext{z}))
 ight)$
- the KL term becomes

$$ext{KL}\left(q_{\phi}(ext{z}| ext{x}) \mid\mid p_{ heta}(ext{z})
ight)$$

rather than the original

$$ext{KL}\left(q_{\phi}(ext{z}| ext{x}) \mid\mid p_{ heta}(ext{z}| ext{x})
ight)$$

That is, our new loss \mathcal{L} function has two components

- \mathcal{L}_R
- lacktriangledown the reconstruction loss, as before, but using the encoder $q_\phi(\mathbf{z}|\mathbf{x})$ instead of $p_ heta(\mathbf{z}|\mathbf{x})$
- \mathcal{L}_D
- lacktriangledown the "KL divergence" loss which constrains the approximate $q_\phi(\mathbf{z}|\mathbf{x})$

Advanced: Obtain ${\cal L}$ by rewriting ${f KL}(q_\phi({f z}|{f x}) \mid\mid p_ heta({f z}|{f x})$

Let's derive a simpler expression for ${\cal L}$ by manipulating ${
m KL}(q_\phi({f z}|{f x}) \mid\mid p_ heta({f z}|{f x}))$:

$$\begin{aligned} \operatorname{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \mid\mid p_{\theta}(\mathbf{z}|\mathbf{x})) &= \sum_{z} q_{\phi}(\mathbf{z}|\mathbf{x}) (\log(q_{\phi}(\mathbf{z}|\mathbf{x}) - \log(p_{\theta}(\mathbf{z}|\mathbf{x}))) \\ &= \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{x})} \left(\log(q_{\phi}(\mathbf{z}|\mathbf{x}) - \log(p_{\theta}(\mathbf{z}|\mathbf{x}))) \right) \\ &= \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{x})} \left(\log(q_{\phi}(\mathbf{z}|\mathbf{x})) - \log(p_{\theta}(\mathbf{z})) - \log(p_{\theta}(\mathbf{x})) \right) \\ &- \left(\log(p_{\theta}(\mathbf{x}|\mathbf{z})) + \log(p_{\theta}(\mathbf{z})) - \log(p_{\theta}(\mathbf{x})) \right) \end{aligned}$$

$$egin{array}{lll} \mathbf{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \mid\mid p_{ heta}(\mathbf{z}|\mathbf{x})) &=& \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{x})} \left(\ \log(q_{\phi}(\mathbf{z}|\mathbf{x})) - \left(\log(p_{ heta}(\mathbf{x}|\mathbf{z})) + \log(p_{ heta}(\mathbf{z}|\mathbf{z})) + \log(p_{ heta}(\mathbf{z}|\mathbf{z})) + \log(p_{ heta}(\mathbf{z}|\mathbf{z})) + \log(p_{ heta}(\mathbf{z}|\mathbf{z})) - \log(p_{ heta}(\mathbf{z}|\mathbf{z})) + \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{ heta}(\mathbf{z}|\mathbf{z})) + \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \right) \right) \\ \mathcal{L} &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{ heta}(\mathbf{z}|\mathbf{z})) + \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{ heta}(\mathbf{z}|\mathbf{z})) + \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \right) \right) \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{ heta}(\mathbf{z}|\mathbf{z})) + \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \right) \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{ heta}(\mathbf{z}|\mathbf{z})) + \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \right) \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{ heta}(\mathbf{z}|\mathbf{z})) + \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \right) \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{\phi}(\mathbf{z}|\mathbf{z})) + \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \right) \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{\phi}(\mathbf{z}|\mathbf{z})) + \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \right) \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{\phi}(\mathbf{z}|\mathbf{z})) + \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \right) \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{\phi}(\mathbf{z}|\mathbf{z})) + \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \right) \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{\phi}(\mathbf{z}|\mathbf{z})) + \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \right) \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{\phi}(\mathbf{z}|\mathbf{z})) + \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \right) \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{\phi}(\mathbf{z}|\mathbf{z})) + \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \right) \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{\phi}(\mathbf{z}|\mathbf{z}) \right) \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{\phi}(\mathbf{z}|\mathbf{z}) \right) \\ \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{\phi}(\mathbf{z}|\mathbf{z})) \right) \\ \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{\phi}(\mathbf{z}|\mathbf{z}) \right) \\ \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{\phi}(\mathbf{z}|\mathbf{z}) \right) \\ \\ \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p_{\phi}(\mathbf{z}|\mathbf{z}) \right) \\ \\ \\ &=& -\mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{z})} \left(\log(p$$

The LHS cannot be optimized via SGD (recall the tractability issue with $p_{ heta}(\mathbf{z}|\mathbf{x})$).

But the RHS can be made tractable giving a tractable choice of $p_{ heta}(\mathbf{z})$.

Choosing $p_{ heta}(\mathbf{z})$

So what distribution should we use for the prior $p_{ heta}(\mathbf{z})$?

- It should be differentiable, since we use Gradient Descent for optimization
- It would be advantageous if it had a tractable closed form (such as a normal)
- If we choose $p_{ heta}(\mathbf{z})$ as normal, we can require $q_{\phi}(\mathbf{z}|\mathbf{x})$ to be normal too
 - The KL divergence between two normals is an easy to compute function of their means and standard deviations.
 - See <u>VAE tutorial (https://arxiv.org/pdf/1606.05908.pdf)</u> Section 2.2

| So it may be fair to say that the idea for architecture of the VAE was obtained from the Loss function, rather than vice-versa. |
|---------------------------------------------------------------------------------------------------------------------------------|
| |
| |
| |

Re-parameterization trick

There is still one more problem for training:

• sampling $\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})$ which appears in the reconstruction loss \mathcal{L}_R term of the loss \mathcal{L} :

$$\mathcal{L}_R = \mathbb{E}_{z \sim q_\phi(\mathbf{z}|\mathbf{x})} \left(\log(p_ heta(\mathbf{x}|\mathbf{z}))
ight)$$

This is not a problem in the forward pass.

But optimization using Gradient Descent requires the ability to differentiate the loss with respect to trainable paramters.

And we don't know how to back propagate through a stochastic node

In particluar

$$rac{\mathcal{L}_R}{\partial \phi}$$

involves a random sample $z \sim q_\phi(\mathbf{z}|\mathbf{x})$ from a function of ϕ .

How do we take the derivative of a node involving a random choice?

The "reparameterization trick"

- redefines z
- ullet as the transformation of a distribution $p(\mathbf{z})$
 - (Looking ahead: which will wind up as unit normal distribution N(0,1))
- by scaling it by a standard deviation and adding a mean

$$\mathbf{z} = \mu_{\theta}(\mathbf{x}) + \sigma_{\theta}(\mathbf{x})\epsilon$$
 $\epsilon \sim p(\mathbf{z})$

Reparameterization trick

The key is that the mean $\mu_{ heta}(\mathbf{x})$ and standard deviation $\sigma_{ heta}(\mathbf{x})$

- are functions of input $\mathbf{x}^{(i)}$
 - so the encoding depends on the input
- ullet and a function of trainable parameters $oldsymbol{ heta}$ and not $oldsymbol{\phi}$
 - so neither affects $\frac{\mathcal{L}_R}{\partial \phi}$

We still can't take derivatives of L_R with respect to ϵ , but we don't need to !

• it's not a trainable parameter

We only need to take derivates of L_R with respect to $\phi, \theta, \mu, \sigma$, which we can do.

In evaluating derivatives, the ϵ that appears in the result (e.g., derivative wrt σ) can be treated as a constant.

• For a particular example, we can remember the drawn ϵ in the forward pass and use it in the backward pass

This gets us to the (near) final picture of the VAE:

Variational Autoencoder (VAE)

Reparameterization trick

ELBo lower bound

To summarize

- ullet we still have an intractable term (appears as another KL divergence after rewriting loss
 - this term appears as an additive term
 - by definition of **KL**, it is positive
- so we can't evaluate the full loss function
 - but, ignoring the intractable positive part, the remainder is a lower bound on the loss
 - so we optimize the lower bound
 - called the *ELBO* term (LB is lower bound)

Loss function: discussion

Let's examine the role of \mathcal{L}_R and \mathcal{L}_D in the loss function \mathcal{L} .

- What would happen if we dropped \mathcal{L}_D ?
 - We would wind up with a deterministic z and collapse to a vanilla VAE
- What would happen if we dropped \mathcal{L}_R ?
 - ullet the encoding approximation $q_\phi(\mathbf{z}|\mathbf{x})$ would be close to the empirical $p_ heta(\mathbf{z}|\mathbf{x})$ in distribution
 - but two variables with the same distribution are not necessarily the same?
 - e.g., get a distribution **p** by flipping a coin
 - \circ let distribution q be a relabelling of p by changing Heads to Tails and vice-versa
 - \circ **p** and **q** are equal in distribution but clearly different !

```
In [3]: print("Done")
```

Done