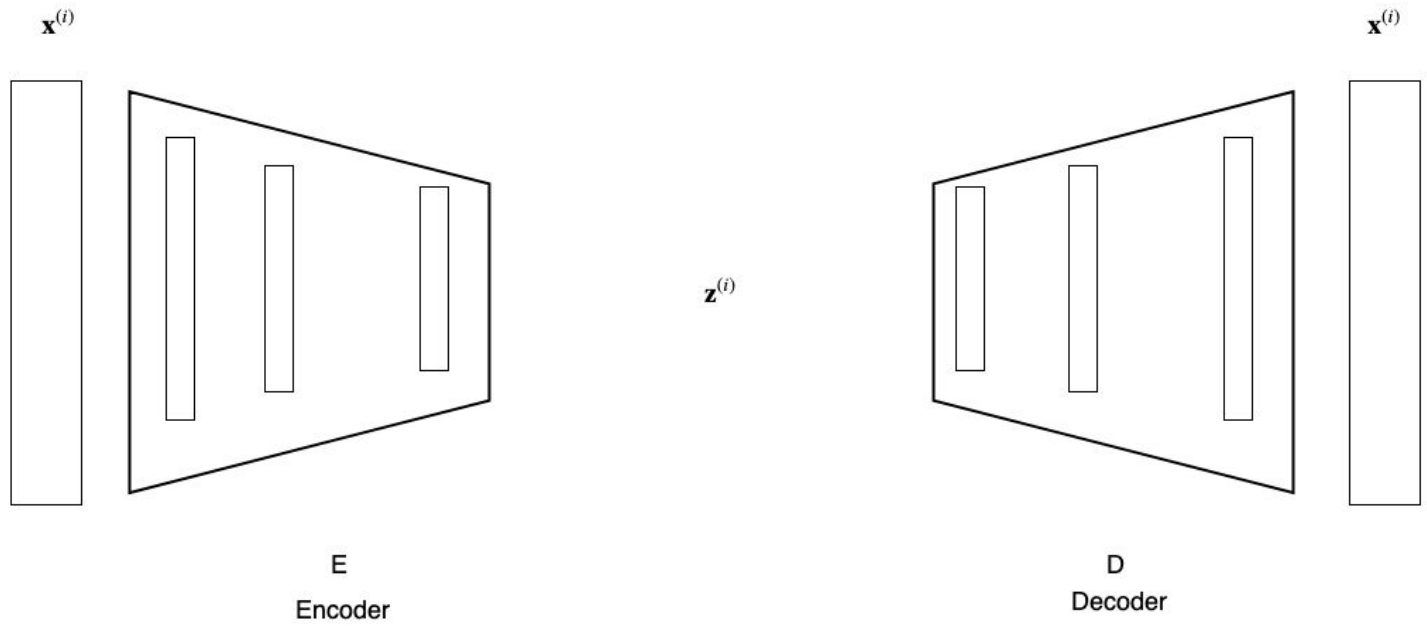


AutoEncoder (AE): High Level

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 is

- A NN consisting of two halves
- An *Encoder* sequence of layers
 - transforms inputs $\mathbf{x}^{(i)}$ to synthetic features $\mathbf{z}^{(i)}$
 - *latent representation*
- A *Decoder* sequence of layers
 - "inverts" latent representation $\mathbf{z}^{(i)}$ to recover $\mathbf{x}^{(i)}$

The Encoder and Decoder are *jointly trained*, not trained separately !

- No obvious target for training the Encoder
- Semi-supervised
 - No labels
 - But we use $\mathbf{x}^{(i)}$ as the "label" associated with example $\mathbf{x}^{(i)}$ in training

This should all be reminiscent of the definition of Principal Components.

Just as for PCA, we can perform dimension-reduction

- size of latent representation $|z| \leq n$

When $|z| < n$ we say

- that the input has been passed through a *bottleneck*
- that the AE is *under complete*

The *main difference* from PCA

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

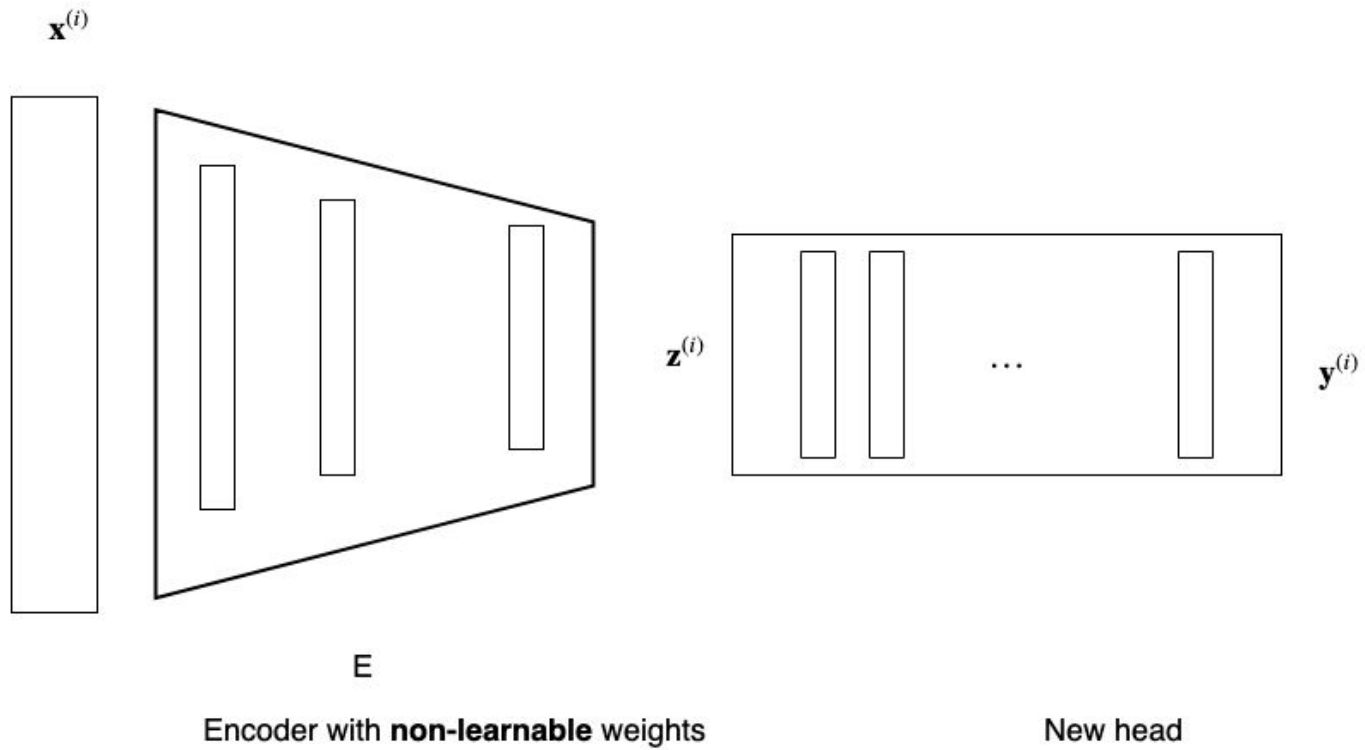
Autoencoders: Uses

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
 - It may be easier to solve the new task given z rather than 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 redominated explained variance in terms of "number of features"

There is no direct similar concept of feature importance in AE

- other than minimizing a Cost function, which *may* wind up focussing 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 \mathcal{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
- Using initialized weights $\mathbf{W}_{(l')}$ for all layers $l' < l$
- Produces $\tilde{\mathbf{x}}^{(i)}$ 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

- to be easier to learn the structure of the input x independent of the labels
- to be easier to learn $W_{(l)}$ incrementally

One the weights $W_{(l)}$ 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.

So it is not uncommon to approach a complicated problem

- by first constructing an autoencoder to come up with an alternate (and smaller) representation of the input space.

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

Generative Artificial Intelligence

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

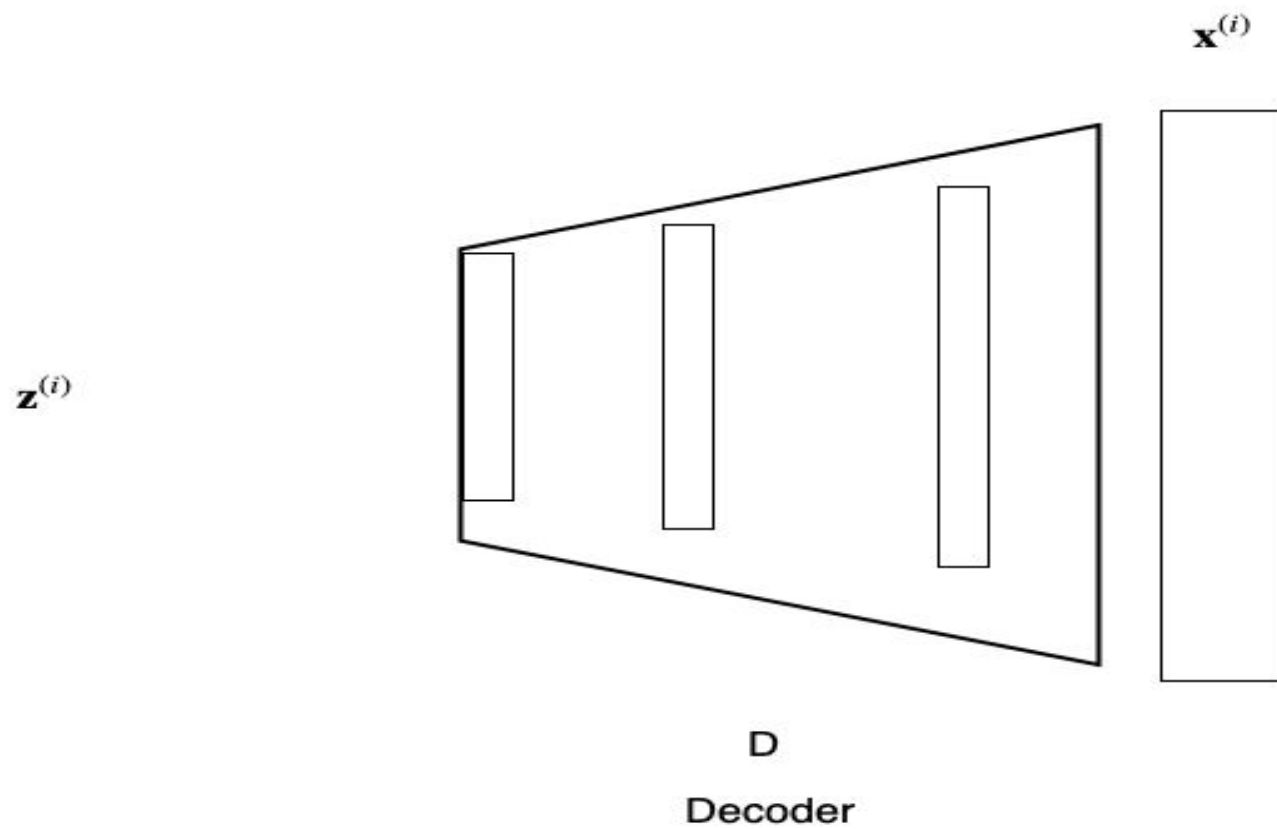
Most of the Machine Learning we have studied thus far is *discriminative*

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

We can use the Decoder on *arbitrary* \mathbf{z} to *generate* a completely new \mathbf{x} :

- $p(\mathbf{x}^{(i')} | \mathbf{z}^{(i')})$ for some i' not in training
- *generate* a new example i' , in the domain of \mathbf{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: $\mathbf{z}^{(i)} = \mathbf{E}(\mathbf{x}^{(i)})$
- Output of Decoder: $\tilde{\mathbf{x}}^{(i)} = \mathbf{D}(\mathbf{z}^{(i)})$
- "Target": $\mathbf{x}^{(i)}$

Both the Encoder and Decoder are parameterized (learnable parameters)

- Goal: find the parameters such that

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

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

Cost/Loss function

There are two obvious candidates for per-example loss

Mean Squared Error (MSE)

$$J^{(i)} = \sum_{j=1}^{|Z|} (\mathbf{x}_j^{(i)} - \tilde{\mathbf{x}}_j^{(i)})^2$$

Binary Cross Entropy

For the special case where *each* original feature is in the range $[0, 1]$

- e.g., pixel is on/off
- we can treat each original feature as a probability and use Binary Cross Entropy

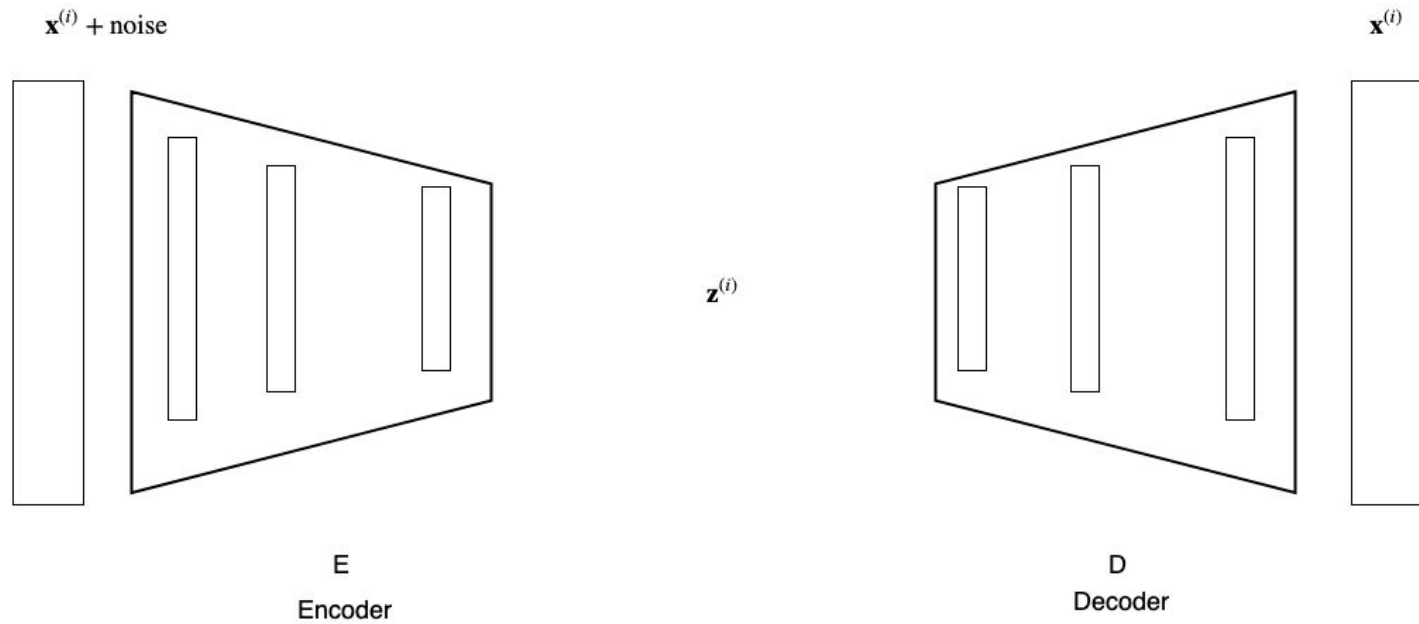
$$\mathcal{L}^{(i)} = \sum_{j=1}^{|z|} \left(x_j^{(i)} \log(\tilde{x}_j^{(i)}) + (1 - x_j^{(i)}) \log(1 - \tilde{x}_j^{(i)}) \right)$$

Autoencoder: extensions

There are some extensions of the "vanilla" Autoencoder we have described thus far.

De-noising Autoencoder

Denoising Autoencoder



Variational Autoencoder (VAE): Generative ML

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

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

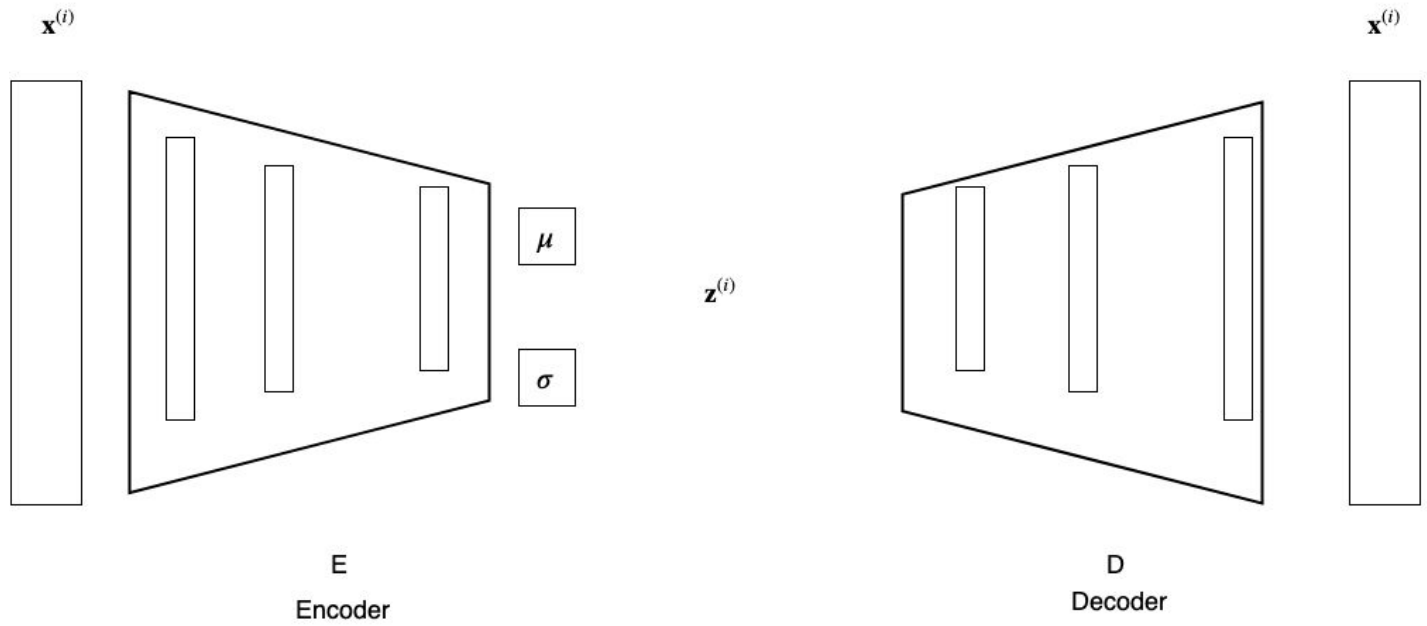
That is: the AE has not been trained to *extrapolate* beyond the training inputs.

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 \mathbf{X} .

Our goal is constructing a Decoder that can extrapolate.

Variational Autoencoder (VAE)



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

The difference is that \mathbf{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 $\mathbf{z}^{(i)}$
- It will compute *variables* $\mu^{(i)}$ and $\sigma^{(i)}$
 - $\mathbf{z}^{(i)}$ will be *sampled* from a distribution with mean $\mu^{(i)}$ and standard deviations $\sigma^{(i)}$

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

Note

μ and σ are computed values (and hence, functions of \mathbf{x}) and not parameters

- so training learns a *function* from $\mathbf{x}^{(i)}$ to $\mu^{(i)}$ and $\sigma^{(i)}$

To train a VAE:

- pass input $\mathbf{x}^{(i)}$ through the encoder, producing $\mu^{(i)}, \sigma^{(i)}$
 - use $\mu^{(i)}, \sigma^{(i)}$ to sample a latent representation $\mathbf{z}^{(i)}$ from the distribution
- pass the sampled $\mathbf{z}^{(i)}$ through the decoder, producing $\mathcal{D}(\mathbf{z}^{(i)})$
- measure the reconstruction error $\mathbf{x}^{(i)} - \mathcal{D}(\mathbf{z}^{(i)})$, just as in a vanilla AE
- backpropagate 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

$$E(\mathbf{x}) = q_{\phi}(\mathbf{z}|\mathbf{x}) \approx p_{\theta}(\mathbf{z}|\mathbf{x})$$

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

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 Cost functions that reflect our objectives

Conditional VAE

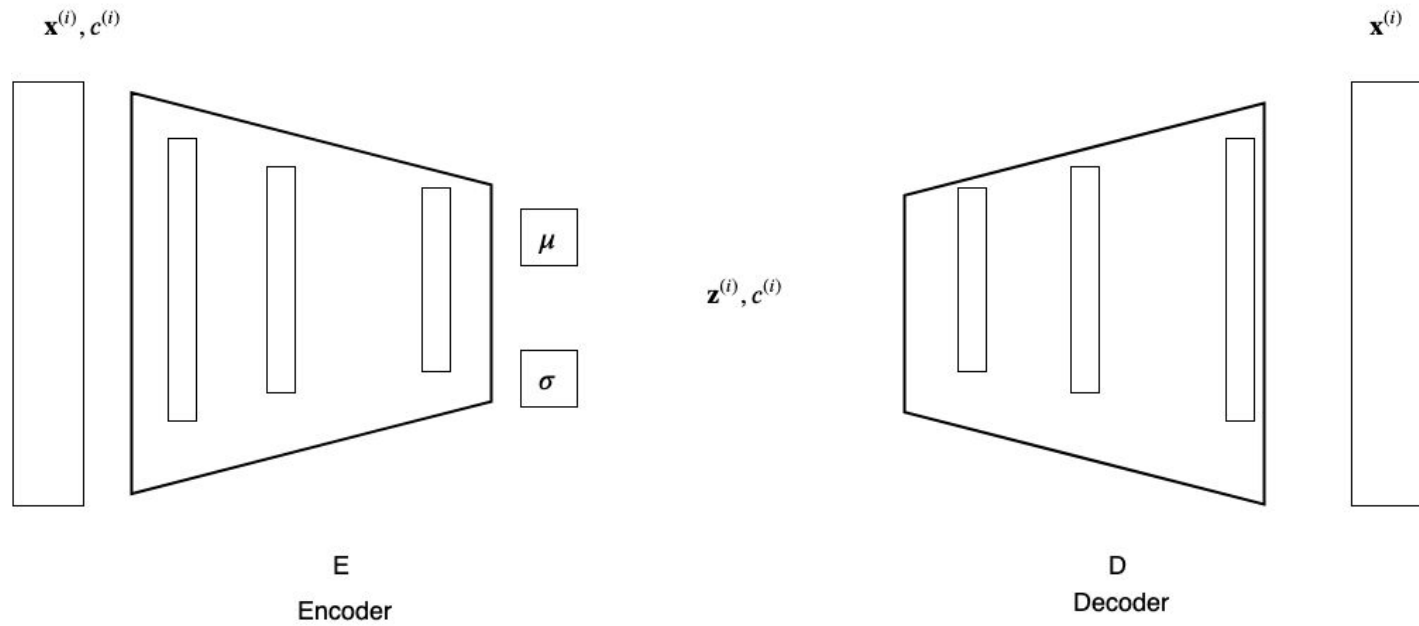
Once a VAE is trained, $\mathcal{D}(z)$ should produce a realistic output, for any z from the distribution.

However, if the distribution of \mathbf{X} includes examples from many classes

- Assuming we have labels as auxiliary 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{x} .

Conditional VAE (CVAE)



The class label c

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

$$\tilde{\mathbf{x}}^{(i)} = \mathcal{D}(\mathbf{z}^{(i)}, c)$$

So now we

- create a latent \mathbf{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 c : $q_{\phi}(\mathbf{z}|\mathbf{x}, c)$
- So is the decoding distribution $p_{\theta}(\mathbf{x}|\mathbf{z}, c)$

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

We will now switch to a notebook running on Google Colab [Autoencoder example from github](https://colab.research.google.com/github/kenperry-public/ML_Fall_2019/blob/master/Autoencoder_example.ipynb) (https://colab.research.google.com/github/kenperry-public/ML_Fall_2019/blob/master/Autoencoder_example.ipynb).

VAE: Probabilistic formulation

Note: Advanced material

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 [VAE tutorial](https://arxiv.org/pdf/1606.05908.pdf) (<https://arxiv.org/pdf/1606.05908.pdf>).

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 begin with a distribution $p(\mathbf{z})$ of the latent variables, and 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 .

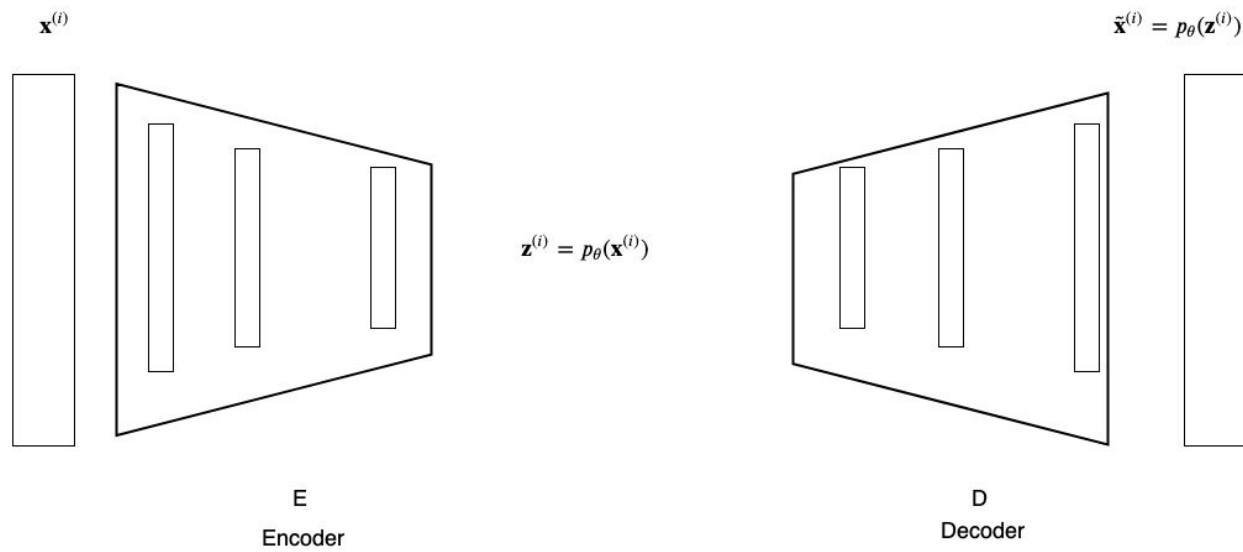
That is: we will train a model to learn θ .

So henceforth p will be subscripted with θ .

We motivate these distributions as they relate to the VAE:

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

VAE derivation: 1



From the joint distribution $p_{\theta}(\mathbf{x}, \mathbf{z})$ we can obtain

- $p_{\theta}(\mathbf{x}|\mathbf{z}) = \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{z})}$
 - the conditional distribution of \mathbf{x} given \mathbf{z}
 - this represents the output distribution of the decoder
- $p_{\theta}(\mathbf{z}|\mathbf{x}) = \frac{p_{\theta}(\mathbf{x}|\mathbf{z})p_{\theta}(\mathbf{z})}{p_{\theta}(\mathbf{x})}$ (by Bayes rule)
 - this represents the distribution for the encoder
- $p_{\theta}(\mathbf{x}) = \int_{\mathbf{z}} p_{\theta}(\mathbf{x}|\mathbf{z})p_{\theta}(\mathbf{z})$
 - the unconditional distribution of \mathbf{x} , the input space, by marginalizing \mathbf{z} .

The loss function: first attempt

Let's try to create a loss function, given that we are dealing with probabilities.

Our first attempt at loss is ℓ :

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

This is maximizing the likelihood of the training data.

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

Can you spot the problem ?

Recall that

$$p_{\theta}(\mathbf{x}) = \int_{\mathbf{z}} p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})$$

But this integral is problematic.

- p_θ does not have a closed form
- we have to approximate the integral numerically
- by iterating over every possible value of multi-dimensional \mathbf{z}

\mathbf{z} is multi-dimensional and to calculate the integral with respect to \mathbf{z} we have to integrate over the full range of each dimension.

As the dimension of \mathbf{z} becomes large, it is no longer computationally tractable to numerically evaluate the integral.

For the same reason $p_\theta(\mathbf{z}|\mathbf{x})$ is problematic since $p_\theta(\mathbf{x})$ appears in the denominator.

Avoiding intractability

The solution is change the objective of the encoder

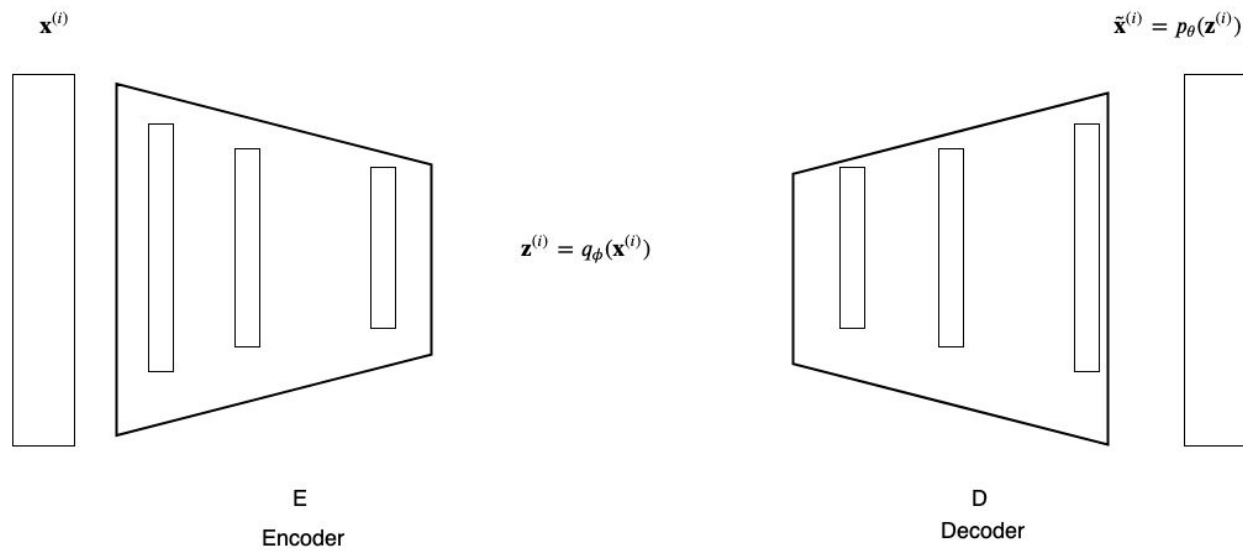
- from producing the intractable $p_{\theta}(z|x)$
- to producing an *approximation* $q_{\phi}(z|x)$
- that is both *tractable* and "close" (in distribution) to the intractable $p_{\theta}(z|x)$.

As usual, we use the KL divergence as a measure of similarity of two distributions:

$$\text{KL}(q_{\phi}(z|x) || p_{\theta}(z|x))$$

$q_{\phi}(z|x)$ will be implemented via a NN parameterized by ϕ .

VAE derivation: 2



Our new loss function

$$\begin{aligned}\mathcal{L} &= -\log(p_{\theta}(\mathbf{x})) + \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z}|\mathbf{x})) \\ &= \mathcal{L}_R + \mathcal{L}_D\end{aligned}$$

has two objectives

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

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

$$\mathcal{L}_D = \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z}|\mathbf{x}))$$

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

$$\mathcal{L} = -\mathbb{E}_z (\log(p_\theta(\mathbf{x}|z))) + \text{KL}(q_\phi(z|\mathbf{x}) \parallel p_\theta(z))$$

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

- Re-write
 $\log(p_\theta(\mathbf{x})) = \mathbb{E}_{z \sim p_\theta(z|\mathbf{x})} (\log(p_\theta(\mathbf{x}|z)))$
- the KL term becomes

$$\text{KL} (q_\phi(z|\mathbf{x}) \parallel p_\theta(z))$$

rather than the original

$$\text{KL} (q_\phi(z|\mathbf{x}) \parallel p_\theta(z|\mathbf{x}))$$

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

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

Advanced: Obtain \square by rewriting $\text{KL}(q_\phi(z|x) || p_\theta(z|x))$

Let's examine the discrepancy between the approximation $q_\phi(z|x)$ and $p_\theta(z|x)$

$$\begin{aligned}
 \text{KL}(q_\phi(z|x) || p_\theta(z|x)) &= \sum_z q_\phi(z|x) (\log(q_\phi(z|x)) - \log(p_\theta(z|x))) \\
 &= \mathbb{E}_z (\log(q_\phi(z|x)) - \log(p_\theta(z|x))) \\
 &= \mathbb{E}_z (\log(q_\phi(z|x)) \\
 &\quad - (\log(p_\theta(x|z)) + \log(p_\theta(z)) - \log(p_\theta(x))))
 \end{aligned}$$

$$\begin{aligned}
 \text{KL}(q_\phi(z|x) || p_\theta(z|x)) \\
 - \log(p_\theta(x)) &= \mathbb{E}_z (\log(q_\phi(z|x)) - (\log(p_\theta(x|z)) + \log(p_\theta(z)))) \\
 &= \mathbb{E}_z ((\log(q_\phi(z|x)) - \log(p_\theta(z))) - \log(p_\theta(x|z))) \\
 &= -\mathbb{E}_z (\log(p_\theta(x|z))) + \text{KL}(q_\phi(z|x) || p_\theta(z))
 \end{aligned}$$

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

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

So the RHS is a tractable form that is equivalent to the LHS and will serve as the loss function for the VAE.

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

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

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

One important consideration is that, since we learn by SGD, we need to be able to differentiate.

Another consideration is that the functional form of the distribution (i.e., an empirical distribution doesn't have a closed functional form) may simplify the math (e.g. normal).

To force the tractability of $q_{\phi}(\mathbf{z}|\mathbf{x})$ we will define *prior distribution* $p_{\theta}(\mathbf{z})$ to have a tractable, closed form (often a normal).

Loss function: discussion

The reconstruction loss should be familiar: it tries to force the decoded output to be "close" to the input.

What would happen if we omitted the KL divergence constraint \mathcal{D} from \mathcal{L} ?

Without it, the model could theoretically learn encodings $q_\phi(\mathbf{z}|\mathbf{x})$ whose distribution had near zero variance.

This would collapse the VAE into the vanilla AE.

So by choosing $p_\theta(\mathbf{z})$ with a non-zero variance, we force the encoder to be probabilistic.

Variational inference

The Cost function can be simplified quite a bit

- by choosing $p_{\theta}(\mathbf{z})$ to be a unit Normal
- by choosing $q_{\phi}(\mathbf{z}|\mathbf{x})$ to be Normal

This is beyond the scope of this talk but we refer the reader to [VAE tutorial](https://arxiv.org/pdf/1606.05908.pdf) (<https://arxiv.org/pdf/1606.05908.pdf>) (Also recommended by Geron in footnote 7, Chapt 15).

To summarize

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

Re-parameterization trick

There is still one more problem for training:

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

This is not a problem in the forward pass.

Optimization via back propagation requires the ability to take derivatives of the loss wrt the trainable parameters.

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

The trick is to re-express \mathbf{z} :

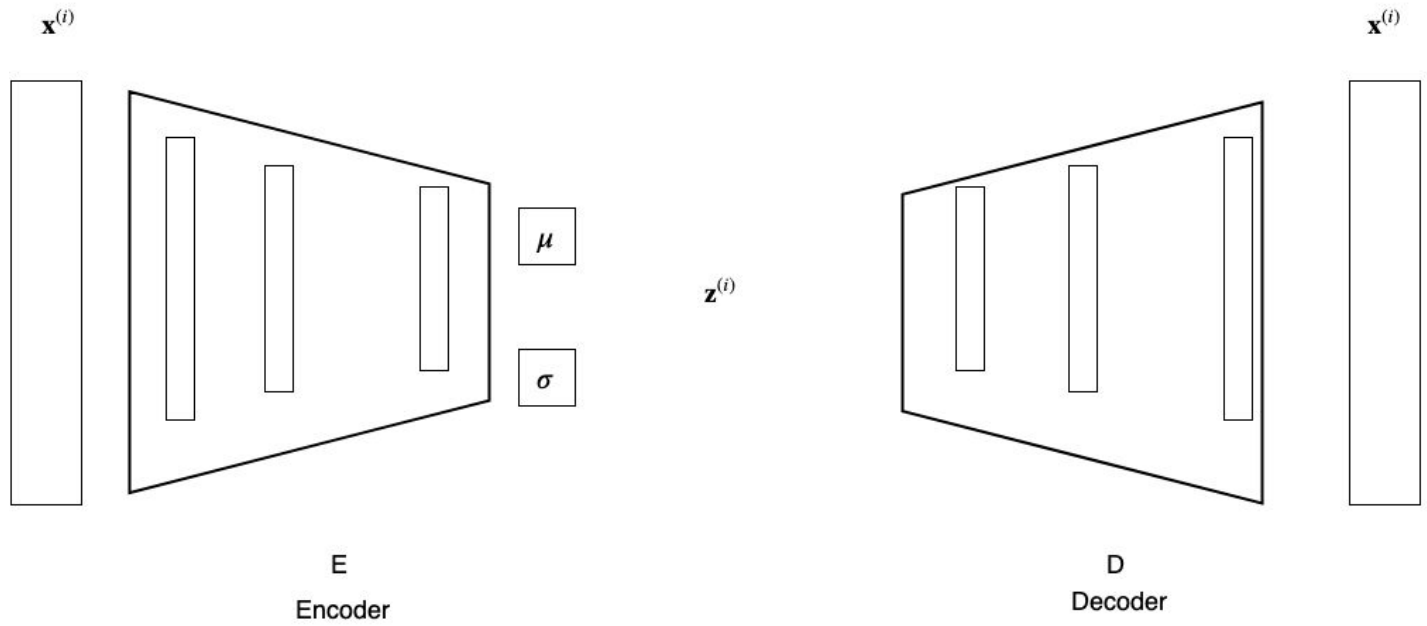
$$\begin{aligned}\mathbf{z} &= \boldsymbol{\mu}_{\theta}(\mathbf{x}) + \boldsymbol{\sigma}_{\theta}(\mathbf{x})\boldsymbol{\epsilon} \\ \boldsymbol{\epsilon} &\sim \hat{\mathcal{P}}(\mathbf{z})\end{aligned}$$

That is, we obtain \mathbf{z}

- By sampling an $\boldsymbol{\epsilon}$ from the constraining distribution $\hat{\mathcal{P}}(\mathbf{z})$
- Scaling the random $\boldsymbol{\epsilon}$ by variable (function of \mathbf{x}) $\boldsymbol{\sigma}$ and adding variable (function of \mathbf{x}) $\boldsymbol{\mu}$.

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

Variational Autoencoder (VAE)



We still can't take derivatives of \mathcal{L}_R with respect to ϵ , but we don't need to !

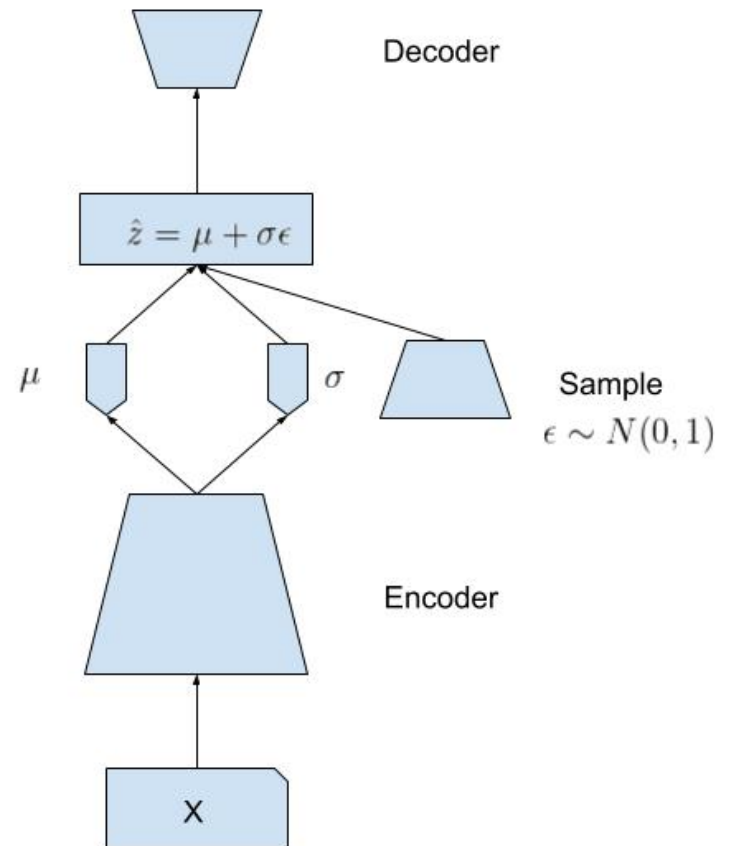
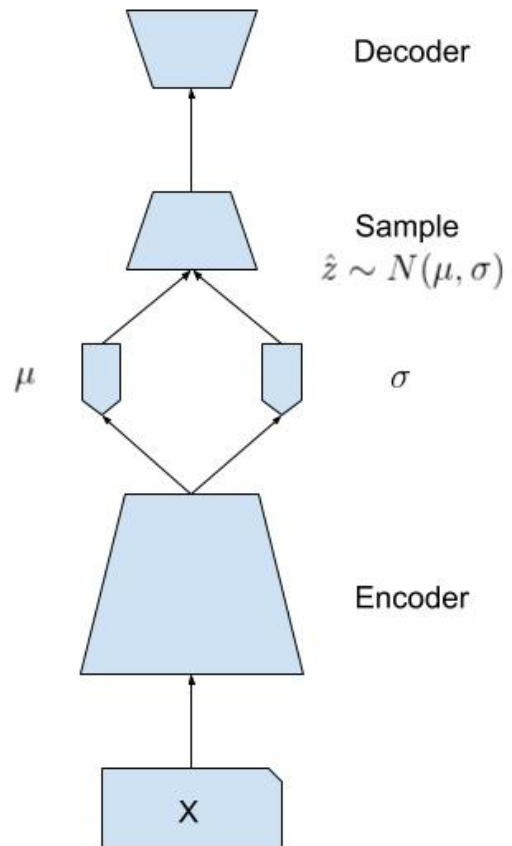
We only need to take derivatives of \mathcal{L}_R with respect to $\phi, \theta, \mu, \sigma$, which we can do.

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 ?

- but over a batch, the expected value over the drawn ϵ should be $E(\hat{p}(z))$ (which is $\mathbf{0}$ if we constrain \hat{p} to be $\mathbf{0}$ centered)

$$\mathcal{L}_D(\phi, \mu, \sigma, \mathbf{x}) = \mathcal{D}_{\text{KL}}(q_\phi(z|\mathbf{x}), \hat{p}_{\mu, \sigma}(z))$$



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

Done