

Autoencoder vs PCA

Autoencoder

encoder weight, bias  
 $Z = f(W_e \cdot x + b_e)$   
 low dimensionality representation  
 original I/P vector

decoder weight, bias  
 $X' = f(W_d \cdot z + b_d)$   
 reconstructed I/P vector

PCA

matrix of eigenvectors of covariance  
 $X = V \Sigma V^T$   
 original data matrix  
 diagonal matrix of variance explained by principle component.

can learn non-linear representation of data  
 (use non-linear activation in stacked autoencoder)

not interpretable

Linear technique.  
 Principal components are linear combinations of original variable.

linear combinations are interpretable.

Self-supervised/unsupervised.

Coming up w/ method to find implicit label  
 self-supervised.

Uses classification or regression loss to do the learning/predict.  
 Labels are implicit w/ data.

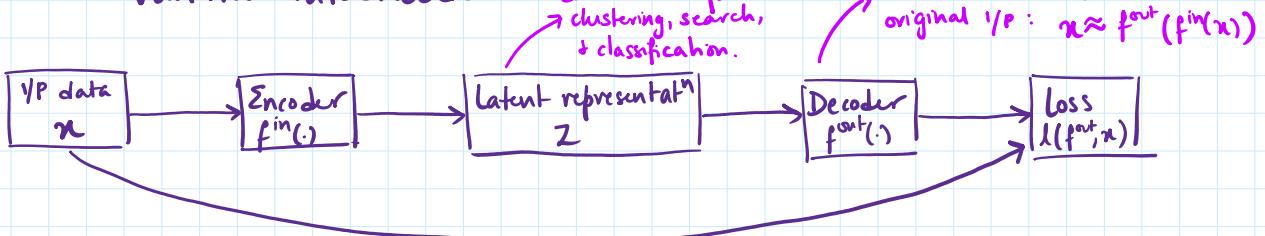
- autoencoder: predict I/P from I/P
- image inpainting: predict missing parts.
- image sorting: break image into parts, then sort.

Unsupervised

Uses of AE

- dimensionality reduction
- anomaly detection
- data compression
- image denoising.

## Vanilla Autoencoder.



- encoder compresses I/P into latent space representation
- decoder reconstructs I/P from  $z$ .
- objective: reconstruct I/P as accurately as possible.
- latent space dimensionality  $\leq$  input dimensionality

## Undercomplete Autoencoder.

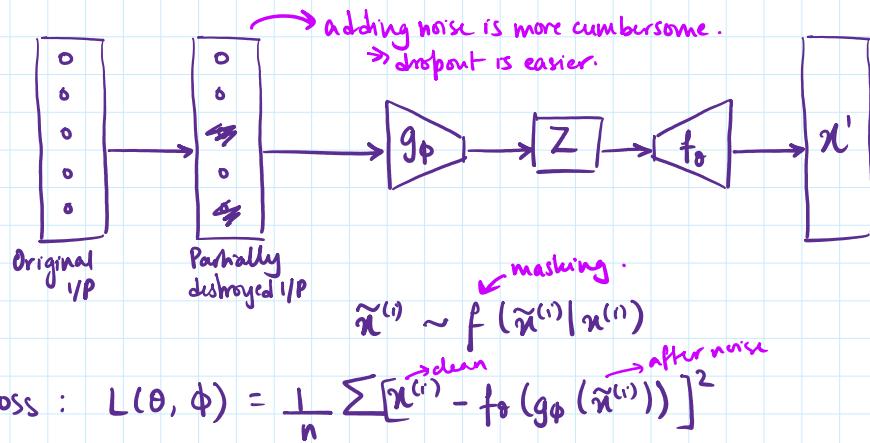


Cost function:  $L(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n [x^{(i)} - f_\phi(g_\phi(x^{(i)}))]^2$

$\downarrow$   
 $(y_{true} - y_{pred})^2$   
 basically MSE.

- $n$   $\downarrow$   
 $(y_{\text{true}} - y_{\text{pred}})^2$   
 basically MSE.
- objective: learn compact + meaningful representation of data.
  - latent space dimensionality  $< \text{I/P}$  dimensionality.  
should be small otherwise latent rep. may not be effective  
(it's just copying the I/P)

## Denoising Autoencoder.



We FINALLY get a denoised version of the original input.

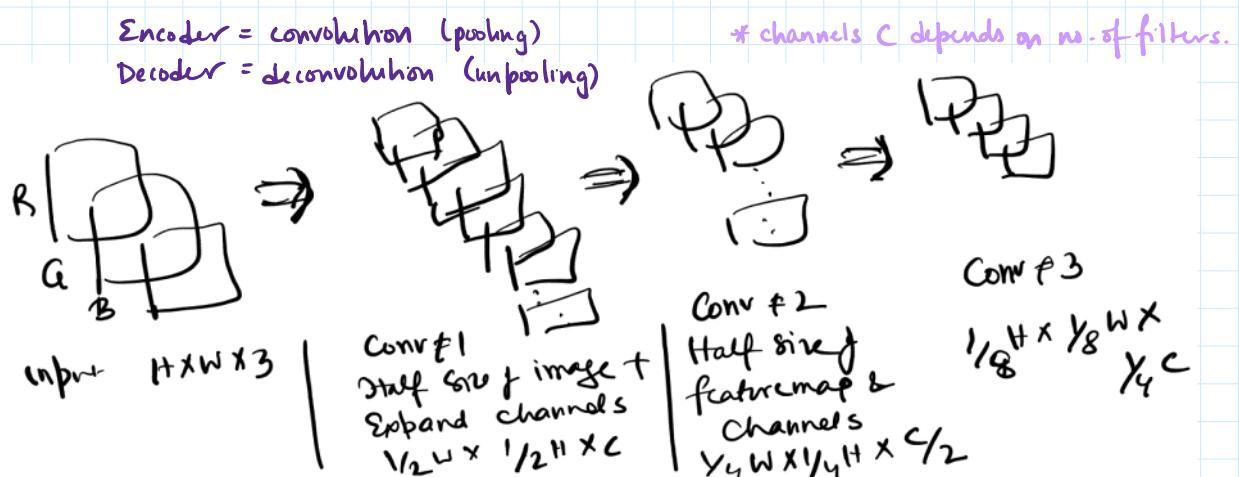
+ network learns to remove noise

= more robust model.

## Autoencoder in Anomaly Detection.

- ① Bottleneck FORCES autoencoder to learn only info from I/P that's impt for reconstruction
- ② If trained on normal data → reconstruction loss typical of non-anomalous data.
- ③ If trained on anomalous data → reconstruction loss super high.

## Autoencoders and Convolution.



Deconvolution :

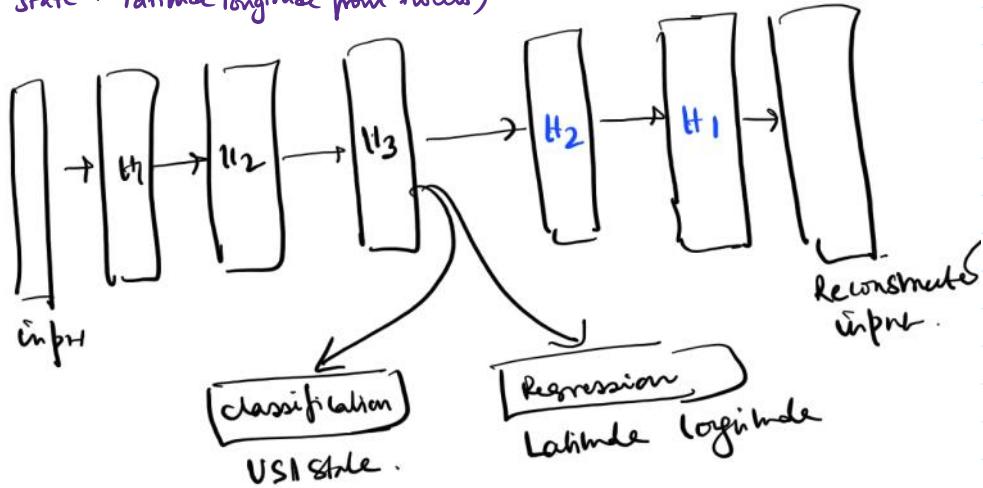
$$\frac{1}{2} H \times \frac{1}{2} W \times \frac{1}{2} C \rightarrow \frac{1}{4} H \times \frac{1}{4} W \times \frac{1}{2} C \rightarrow \frac{1}{8} H \times \frac{1}{8} W \times C \rightarrow H \times W \times 3$$

Deconvolution :

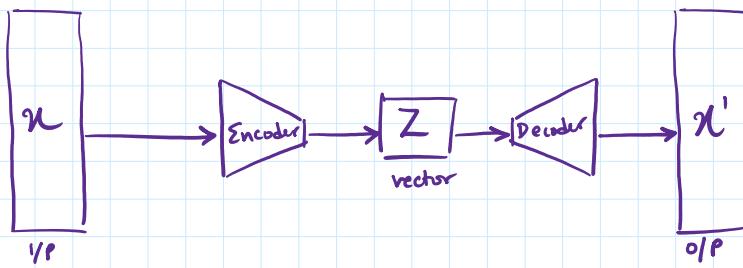
$$\frac{1}{8} H \times \frac{1}{8} W \times \frac{1}{4} C \rightarrow \frac{1}{4} H \times \frac{1}{4} W \times \frac{1}{2} C \rightarrow \frac{1}{2} H \times \frac{1}{2} W \times C \rightarrow H \times W \times 3$$

## Stacked Denoising Autoencoder.

(to infer state + latitude longitude from twists)



Is autoencoder generative? no.



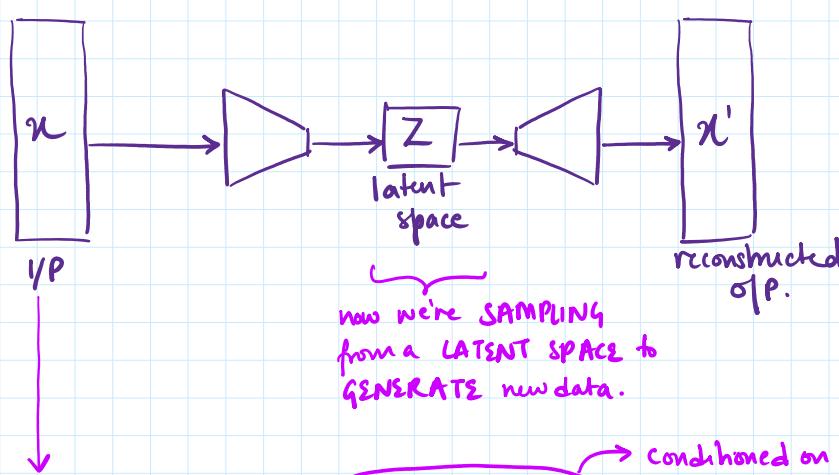
- ① z is a vector — doesn't capture semantic rel'n b/w data samples
- ② AE only learns MAPPING FUNCTION that does COMPRESSION into a vector with a good COMPRESSION LOSS.
- ③ So it can't really generate anything.  
It can only try to replicate the I/P.
- ④ Doesn't learn I/P distribution.

HOWEVER:

an autoencoder w/ generative ability is a variational autoencoder.

## Variational Autoencoder.

learns a latent space instead of learning a vector.



represents parameters of multivariate distribution, origin of I/P.

hence we can generate similar objects.

↓  
Model  $1/P$  as coming from random variable  $X$  → conditioned on other random variable  $Z$ .  
↓ latent (hidden)  
modelled as multivariate gaussian

Now, since  $X$  is conditioned on  $Z$ :

$$p(x) = \int p(x, z) dz$$

likelihood of our data

marginalize over joint probability WRT latent variable

not possible.  
Too computationally slow.

We can also write this as:  $p(z|x) = \frac{p(x,z)}{p(x)}$

or  $p(x) = \frac{p(x,z)}{p(z|x)}$

prob. of latent variable given data.  
Also unknown.

Issues:

① To have  $p(x)$ , we need  $p(z|x)$ .

Not computationally tractable.

② To have  $p(z|x)$  we need tractable  $p(x)$ .

$$\textcircled{3} \quad p(x) = \frac{p(x,z)}{p(z|x)} \quad p(z|x) = \frac{p(x,z)}{p(x)} \quad \rightarrow \text{So we approximate.}$$

Let true posterior to be found be  $p_0(z|x)$

Like any probability distribution:  $\int q_\psi(z|x) dz = 1$ .

Then a whole bunch of math.

- ① log likelihood
- ② Change approximation to expectation
- ③ Expand on chain rule of probability

$$\log p_0(x) = \text{ELBO} + D_{KL}(q_\psi(z|x) || p_0(z|x))$$

↑ evidence based lower bound.

↑ cuz  $\log p_0(x) \geq \text{ELBO}$

↑ maximize thus

↑ in order to maximize that.

↑ Maximize reconstruction likelihood of decoder

$\geq 0$  ↑ minimize distance of learnt distribution from prior belief of  $z$ .

\*  $L(\theta, \psi; x)$  optimize wrt  $\psi$  (variational param) and  $\theta$  (generational param) is a challenge.

→ cuz gradient loss bound wrt  $\psi$  is problematic due to stochastic nature.

Now,

→ VAE uses an external source of noise that's combined w/  $z$ .

↓  
backprop can't be done over stochastic  $z$ .  
→ reparameterization

→ Also, actual loss fn  $L(\theta, \psi; x) =$    
KL divergence of  $\text{z-space} + \text{learnt}$   $\text{z-space}$  + reconstruction MSE loss.

So, what is the net carry forward from this note?

- (a) AutoEncoder is not a generative model but VAE is
- (b)  $z$  is a vector  $y$ .  $z$  is latent space  $\mathbb{Z}$
- (c) Sample from observable data  $x$ , we are learning latent space  $z$  and generating  $x'$  w.r.t.  $z$ . Hence, it is generative.
- (d) Overall maximization of log likelihood of  $x$  i.e.  $\log p_{\theta}(x)$  is turned into a maximization of ELBO (that Proj is excluded) broken into 2 parts (that Proj is excluded)
  - (i) Reconstruction loss (NLL)
  - (ii) DKL between the prior  $z$  and learnt  $z$ .
- (e) Actual implementation of the "loss function" has a challenge as Gradient descent or stochastic  $\epsilon$  can be difficult (variance can be large) + by reparameterization trick, is done by taking out source of variance external variable!!

### The Trick

The reparameterization trick allows us to bypass this problem by transforming the sampling process into a deterministic one. Here's how it works:

1. **Latent Variables:** Instead of sampling directly from the latent distribution, we sample from a standard normal distribution (which is easy to work with).
2. **Transformation:** We then transform this sample using the mean and standard deviation learned by the encoder network. This transformation is deterministic and allows us to backpropagate gradients through the network.