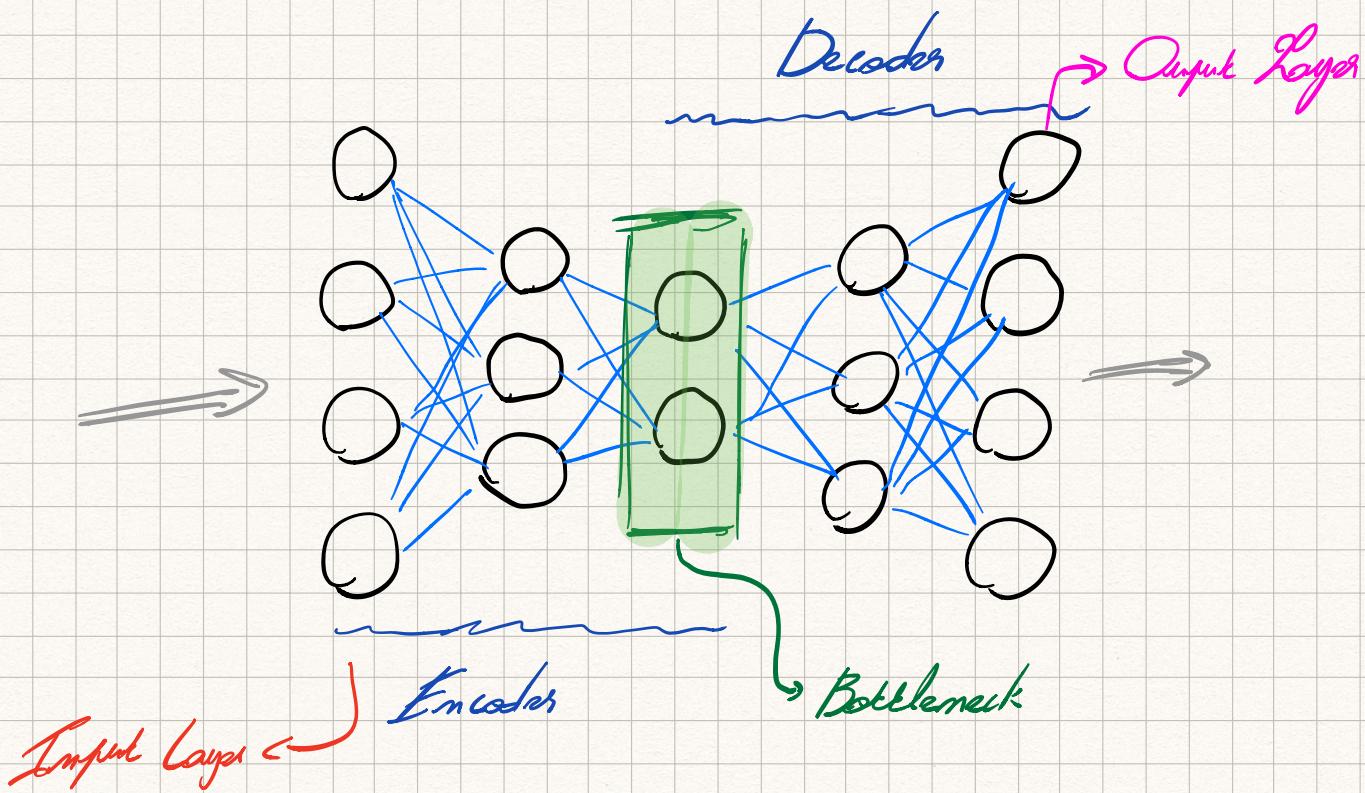
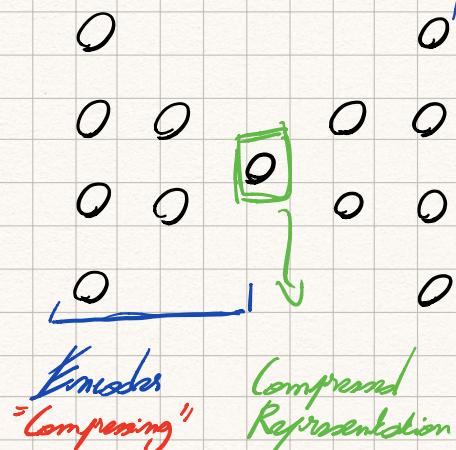


Autoencoder:



Decoder = "Decompressing"



Autoencoders are neural networks that encode try to learn = lower data dimension "representations". By training a double element network:

- Encoder
- Decoder (could be used to generate future = never seen on training "images")

⚠ It can not be used to generate new data, only decoding

↳ We do not know the = "compressed data representation", i.e. which distribution follows our "encoder" to generate numbers.

↳ VAE try to address this limitation.

The goal in autoencoders is to minimize the reconstruction error;

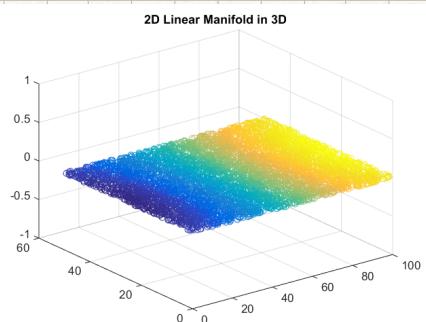
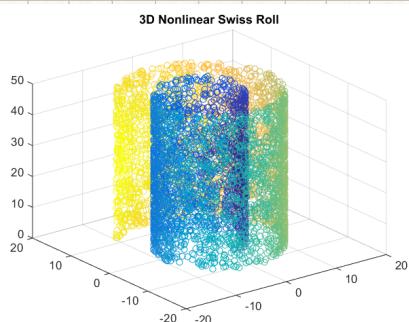
$$L(x, \hat{x}) = \text{loss of } \hat{x} \text{ with respect to } x.$$

- ⚠ If we create an AE with no activation functions (no non-linearity introduced) we would observe some kind of PCA (Principal Component Analysis)
- ②② PCA has, though, a difference, that is the imposed restriction of orthonormality, not present in AE. \Rightarrow PCA more efficient; easier to interpret.

Usually, in order to avoid network overfitting, we use regularization:

$$L(x, \hat{x}) + \text{regularizer} \rightarrow \text{The "Bottleneck" reduced dimension "ensures" it.}$$

②② Manifold \rightarrow continuous, non-intersecting surface.



Visual representation of a non-linear dimensionality reduction resulting in a 2D linear manifold.

Sparse Autoencoders:

Even if we restrict the "Bottleneck" layer enough, the neural network may still memorize the whole train set.

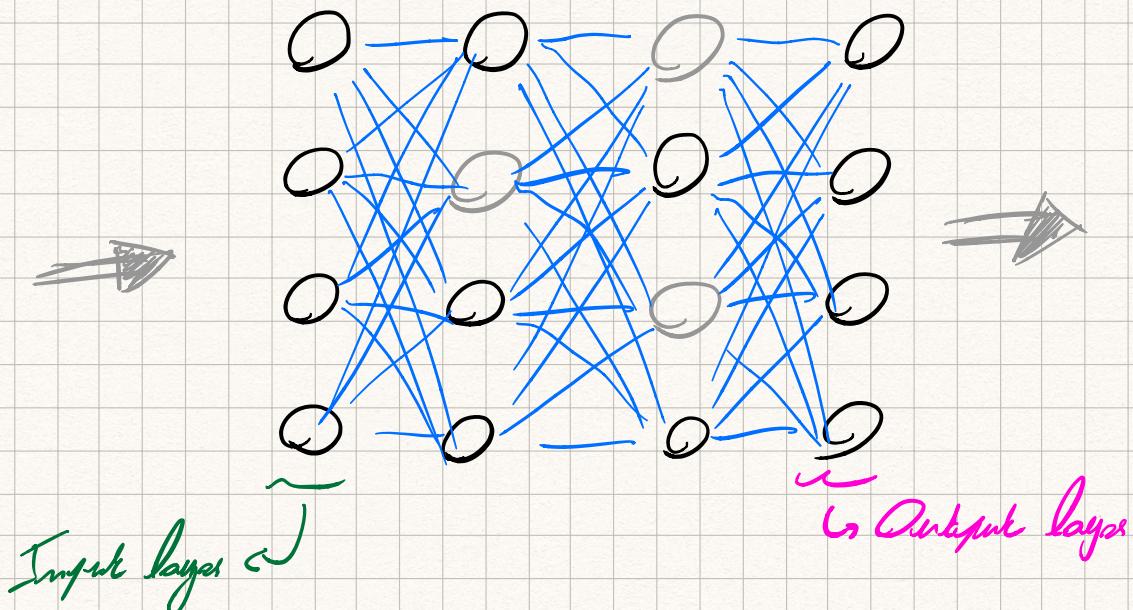
That's why "SAE" (Sparse Autoencoder) exist, we try to train a model that is "data-dependent" \Rightarrow Different neural nets will be created for different data.

AE

Uses full network activation

SAE

Allows to create specific paths depending on data input.



$O \approx$ Active Neuron

$O \approx$ Inactive Neuron

The ways to obtain are:

① L1-Regulation: also called "Lasso regression"

$$\lambda \sum_{j=1}^p |\beta_j| \Rightarrow \text{Shrinks some feature coefficients (neuron weight's) to } 0 \Rightarrow \text{Sparse connection.}$$

② L2-Regulation: Ridge Regression

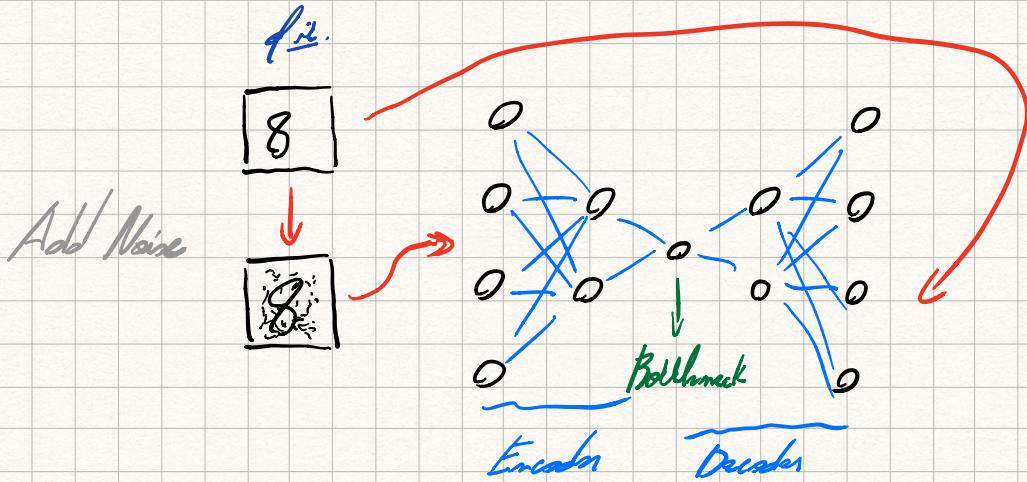
$$\lambda \sum_{j=1}^p \beta_j^2 \Rightarrow \text{Penalty = big values of features} \Rightarrow \text{Weights are smaller, i.e. weaker connection}$$

③ KL-Divergence: difference between 2 probability distributions.

$$L(x, z) + \sum_j KL(p || \hat{p}_j)$$

$$KL \rightarrow \sum_{j=1}^{J^{(h)}} p \log \frac{p}{\hat{p}_j} + (1-p) \log \frac{1-p}{1-\hat{p}_j} \quad (KL-Divergence)$$

Denoising autoencoder: At where the data has been corrupted at the start of the model, in order to "fix the model" to learn the true representation.



Contractive AE: At where we ensure that for similar input, the learned encoding should also be similar.

↳ At learns how to "contract" similar inputs in the same way in the embedding layer.

Constructing a loss function which penalizes = "large derivatives" of our hidden layers with respect to input examples.

↳ $\|A\|_F$ L2-Norm
↳ Squared Frobenius norm $\|A\|_F$ of the gradient matrix S of hidden layers activation with respect to input observation. ↳ First order partial derivatives

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$$