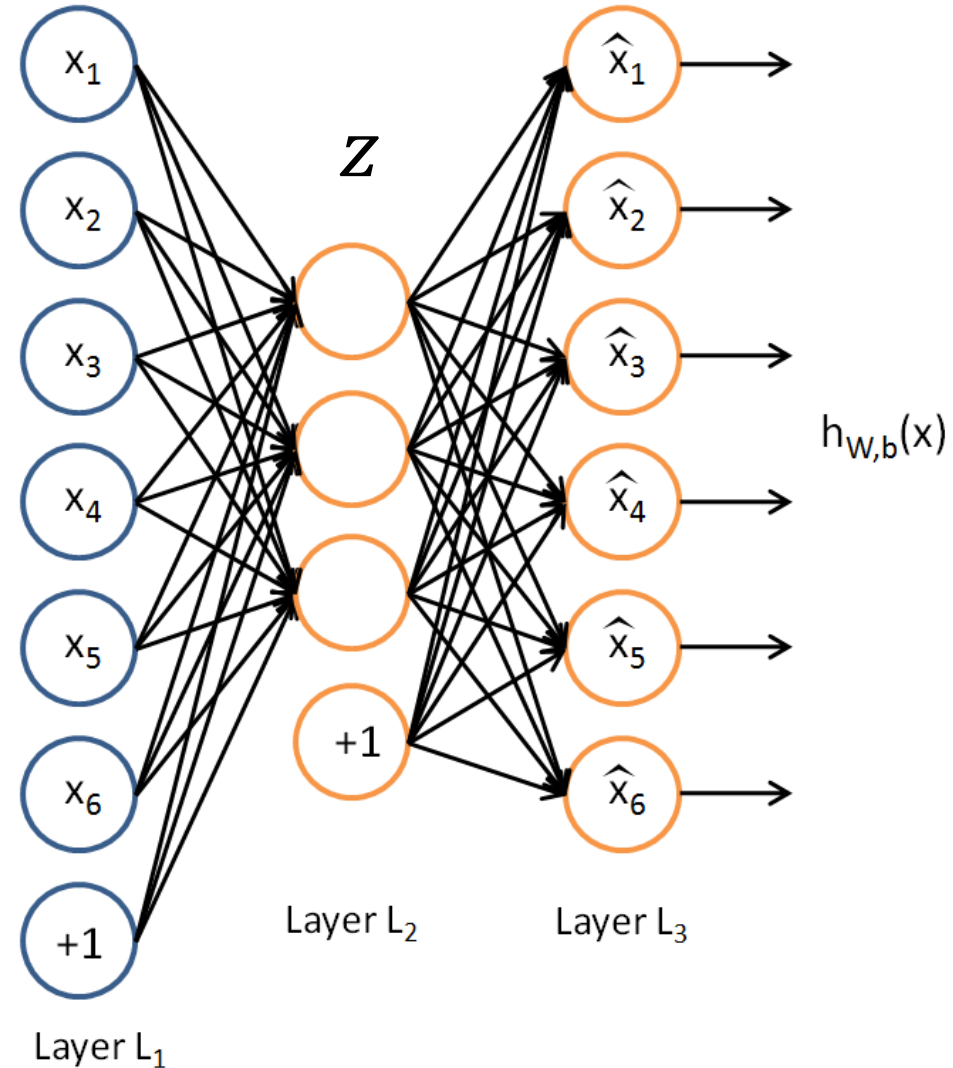# Encoder-decoder model

Biplab Banerjee

# Deep CNN based image segmentation
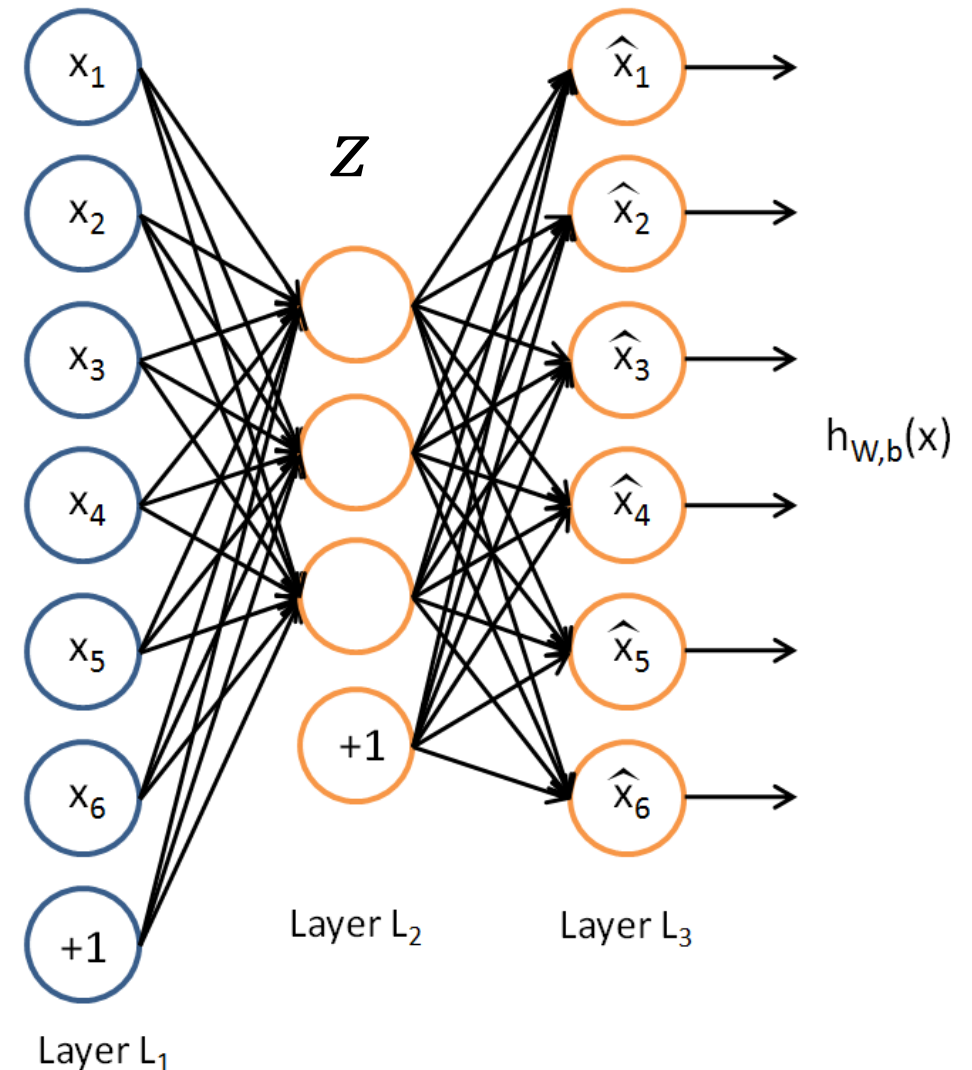
- Auto-encoder
- Regularized auto-encoder
- Class-encoder
- Image segmentation

# Traditional Autoencoder

$x_1$

$x_2$

$x_3$

$z$

$\widehat{x}_1$

$\widehat{x}_2$

$\widehat{x}_3$

$x_4$

$h_{W,b}(x)$

$\widehat{x}_4$

$x_5$

$x_6$

$+1$

$\widehat{x}_5$

$\widehat{x}_6$

$+1$

Layer $L_2$

Layer $L_3$

Layer $L_1$

# Traditional Autoencoder

- Unlike the **PCA** now we can use activation functions to achieve non-linearity.

- It has been shown that an AE without activation functions achieves the **PCA** capacity. (later)

# Uses

- The autoencoder idea was a part of NN history for decades (LeCun et al, 1987).

- Traditionally an autoencoder is used for dimensionality reduction and feature learning.
- Representation learning

# Simple Idea

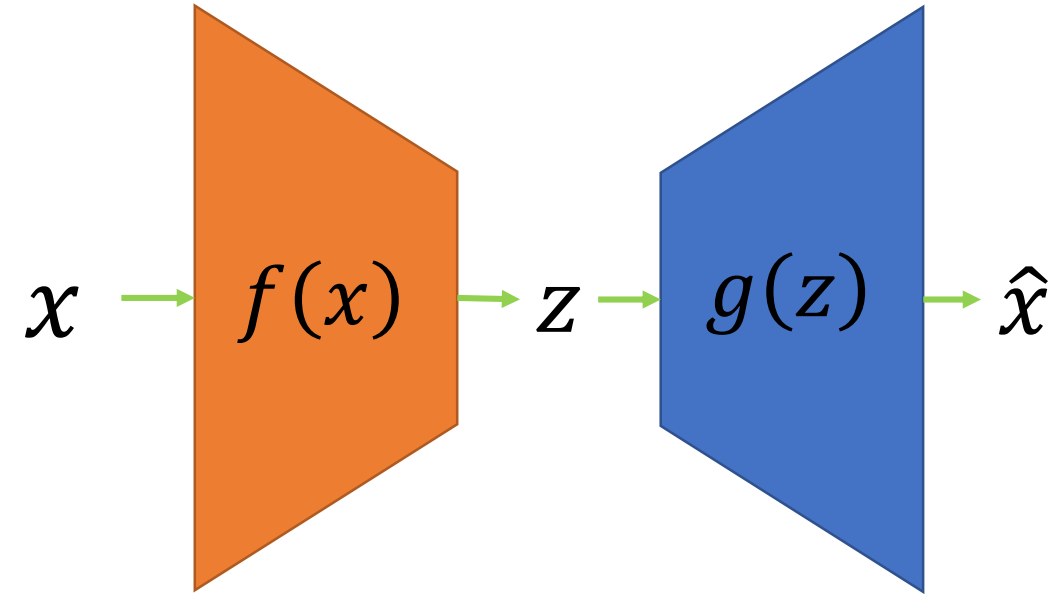- Given data $x$ (no labels) we would like to learn the functions $f$ (encoder) and $g$ (decoder) where:

$$f(x) = s(wx + b) = z$$

and

$$g(z) = s(w'z + b') = \hat{x}$$

s.t $h(x) = g(f(x)) = \hat{x}$

where $h$ is an **approximation** of the identity function.



$x \to f(x) \to z \to g(z) \to \hat{x}$

($z$ is some **latent** representation or **code** and $s$ is a non-linearity such as the sigmoid)

($\hat{x}$ is $x$'s reconstruction)

# Training the AE

Using **Gradient Descent** we can simply train the model as any other FC NN with:
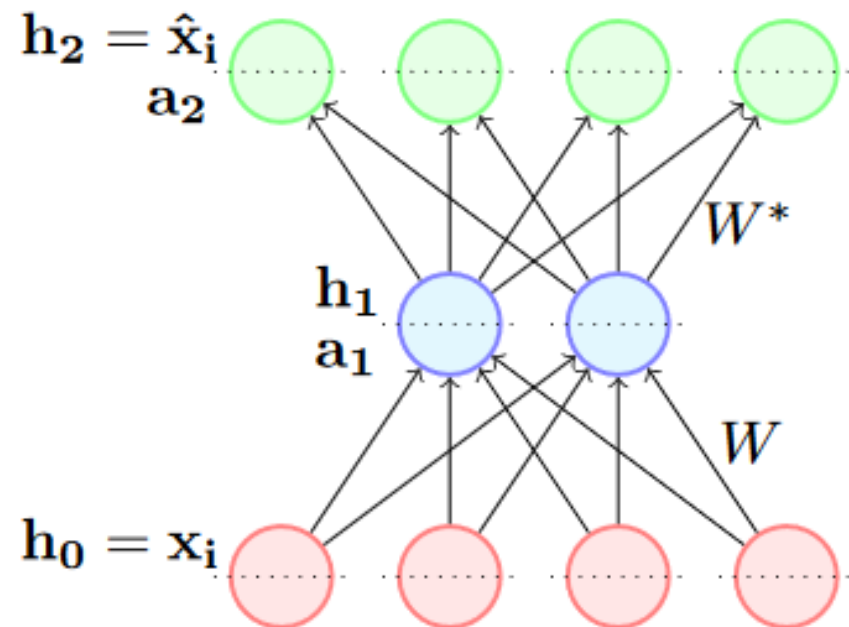
- Traditionally with *squared error* loss function

$$L(x, \hat{x}) = \|x - \hat{x}\|^2$$

- If our input is interpreted as bit vectors or vectors of bit probabilities the *cross entropy* can be used

$$H(p, q) = -\sum_x p(x) \log q(x)$$

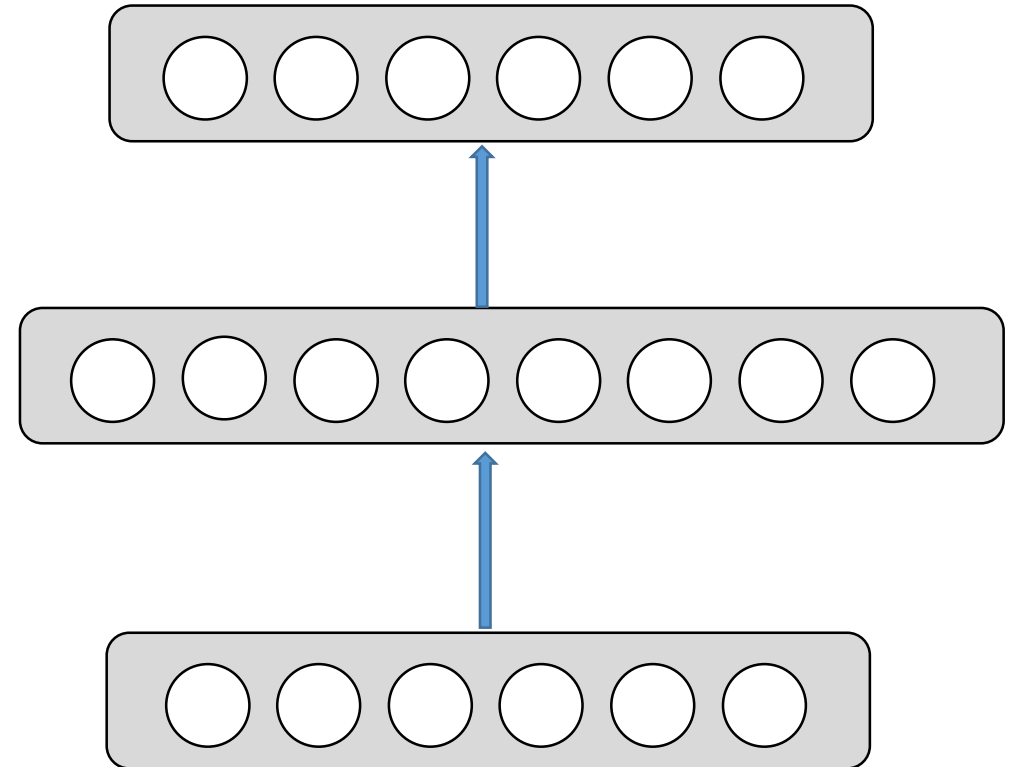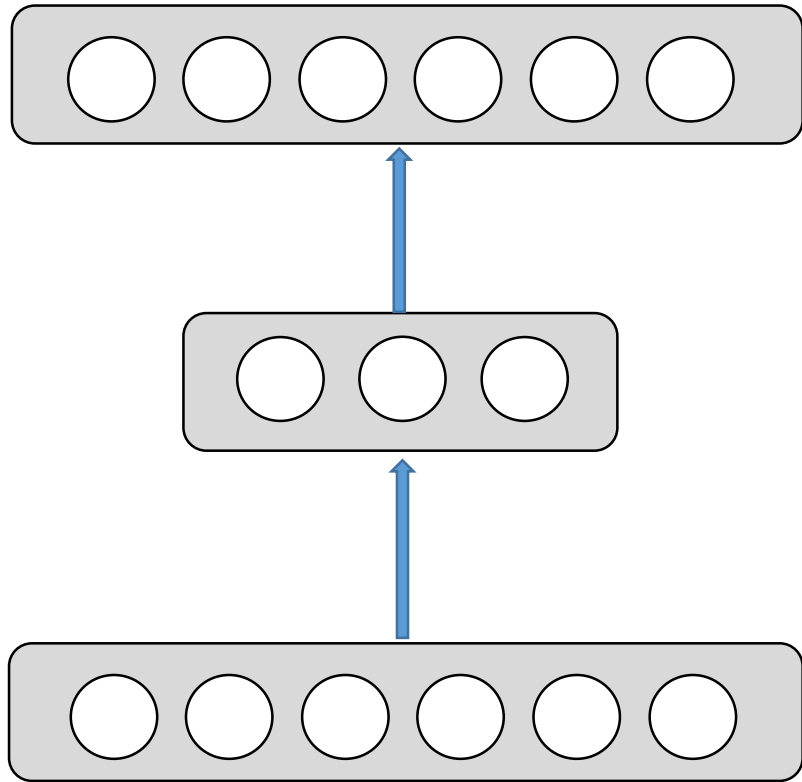$$\mathscr{L}(\theta) = (\hat{\mathbf{x}}_i - \mathbf{x}_i)^T (\hat{\mathbf{x}}_i - \mathbf{x}_i)$$

$\mathbf{h}_2 = \hat{\mathbf{x}}_i$

$\mathbf{a}_2$

$\mathbf{h}_1$

$\mathbf{a}_1$

$W^*$

$W$

$\mathbf{h}_0 = \mathbf{x}_i$

- $\dfrac{\partial \mathscr{L}(\theta)}{\partial W^*} = \dfrac{\partial \mathscr{L}(\theta)}{\partial \mathbf{h}_2} \boxed{\dfrac{\partial \mathbf{h}_2}{\partial \mathbf{a}_2} \dfrac{\partial \mathbf{a}_2}{\partial W^*}}$

- $\dfrac{\partial \mathscr{L}(\theta)}{\partial W} = \dfrac{\partial \mathscr{L}(\theta)}{\partial \mathbf{h}_2} \boxed{\dfrac{\partial \mathbf{h}_2}{\partial \mathbf{a}_2} \dfrac{\partial \mathbf{a}_2}{\partial \mathbf{h}_1} \dfrac{\partial \mathbf{h}_1}{\partial \mathbf{a}_1} \dfrac{\partial \mathbf{a}_1}{\partial W}}$
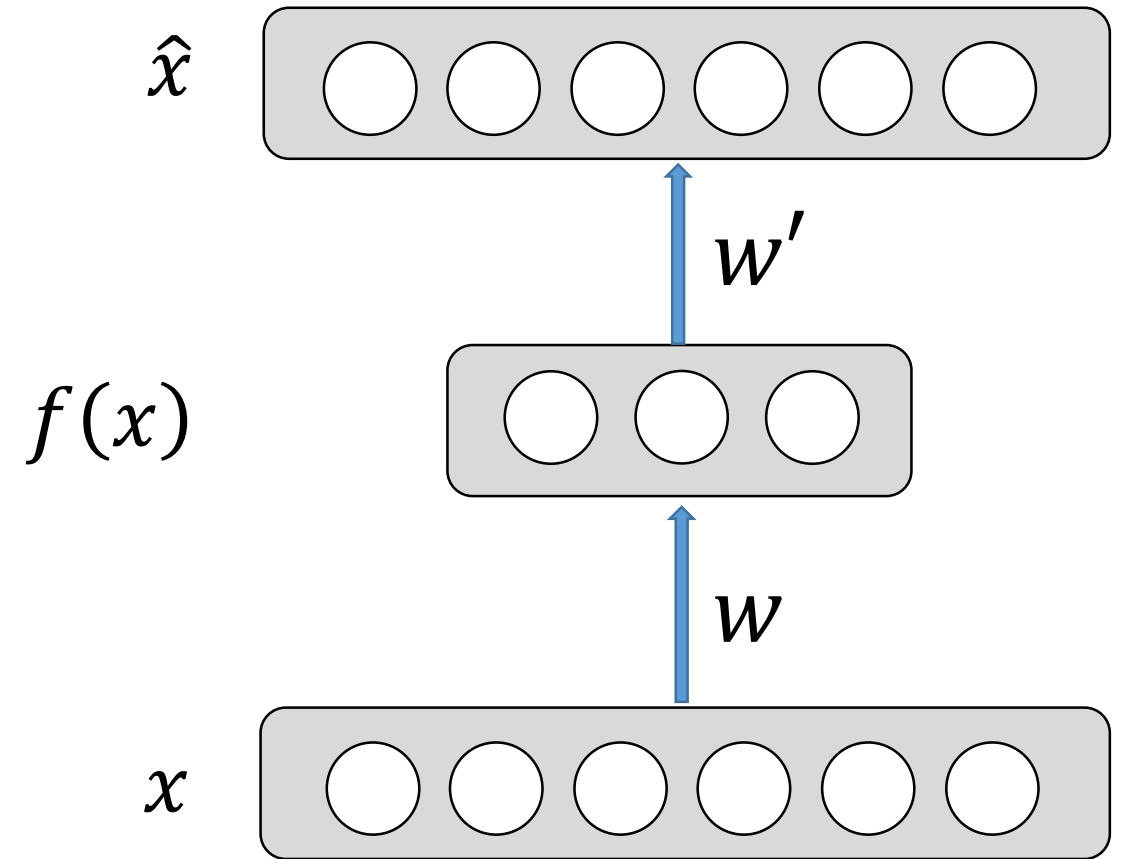
# Undercomplete AE VS overcomplete AE

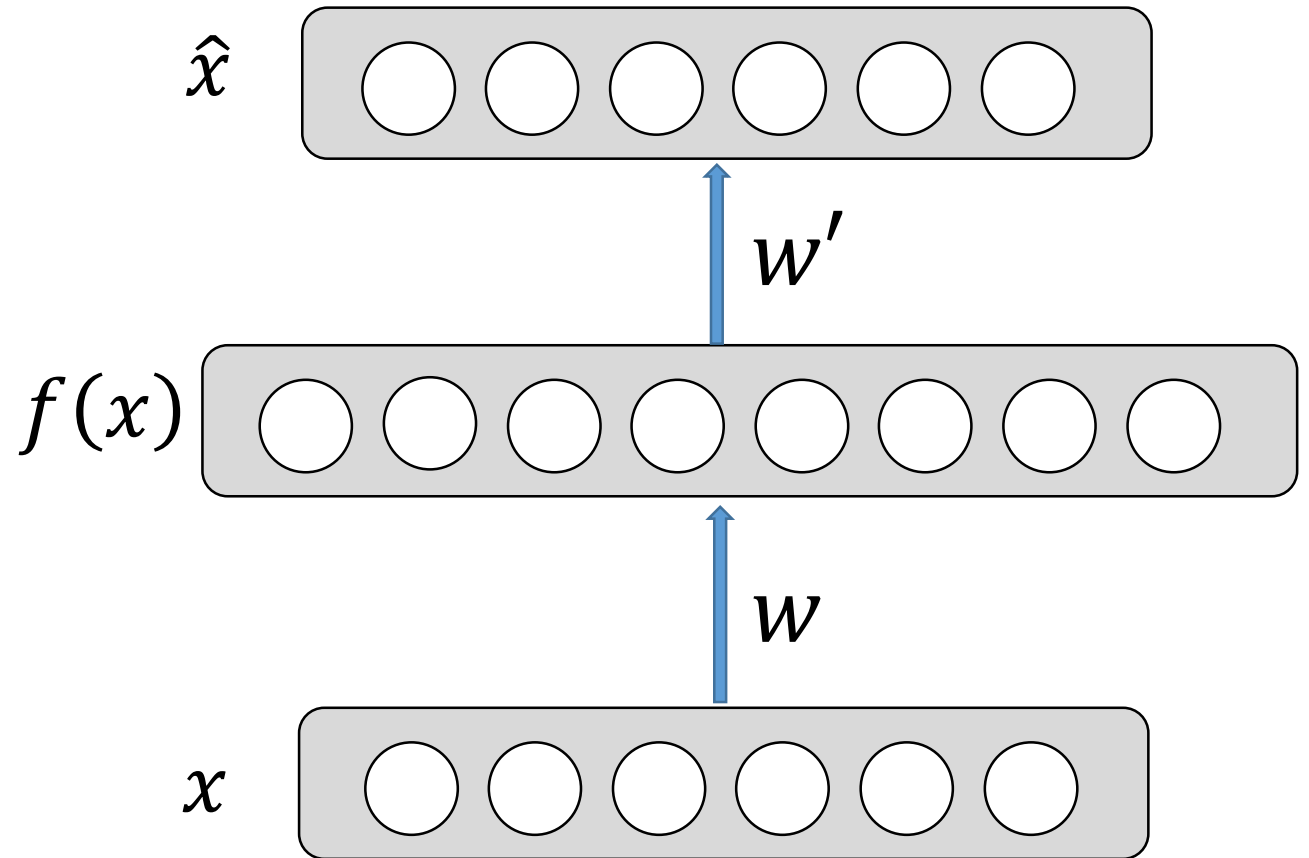We distinguish between two types of AE structures:

# Undercomplete AE

- Hidden layer is **Undercomplete** if smaller than the input layer
  - ❑ Compresses the input
  - ❑ Compresses well only for the training dist.

- Hidden nodes will be
  - ❑ Good features for the training distribution.
  - ❑ Bad for other types on input
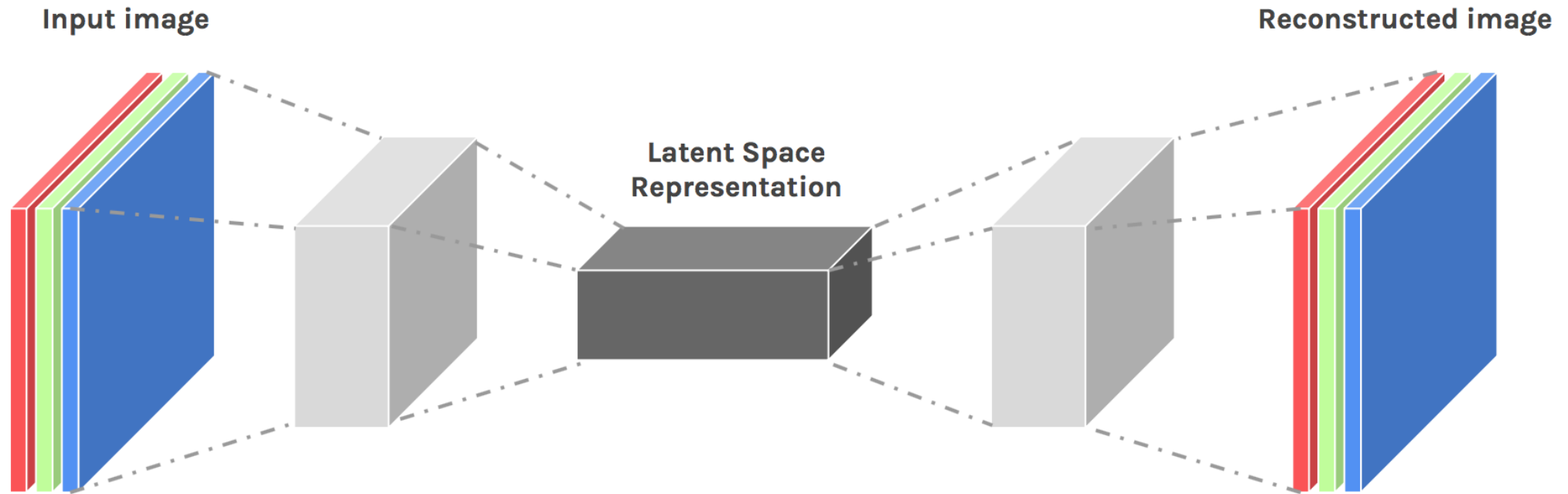
$\hat{x}$

$w'$

$f(x)$

$w$

$x$

# Overcomplete AE

- Hidden layer is **Overcomplete** if greater than the input layer
  - ❑ No compression in hidden layer.
  - ❑ Each hidden unit could copy a different input component.

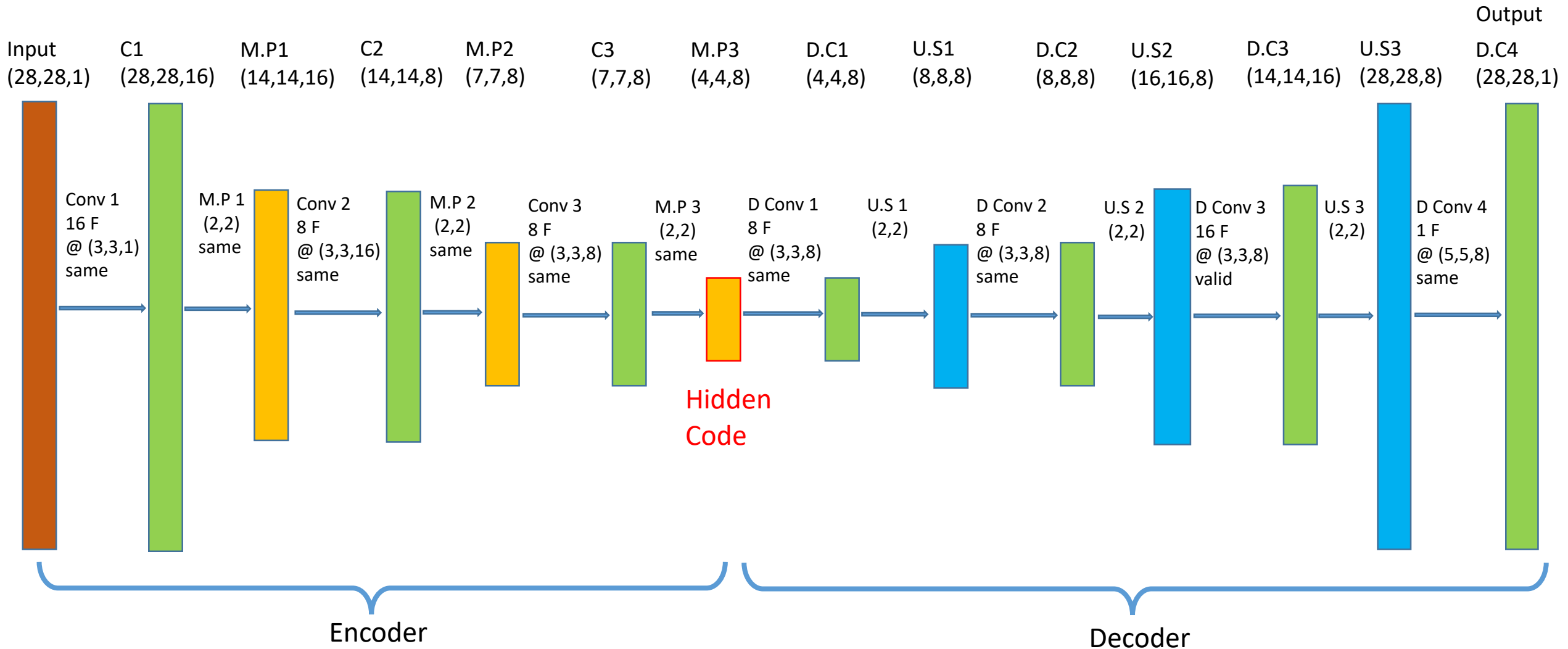- No guarantee that the hidden units will extract meaningful structure.

$\hat{x}$

$W'$

$f(x)$

$W$

$x$

# Convolutional AE



Input image

Latent Space Representation

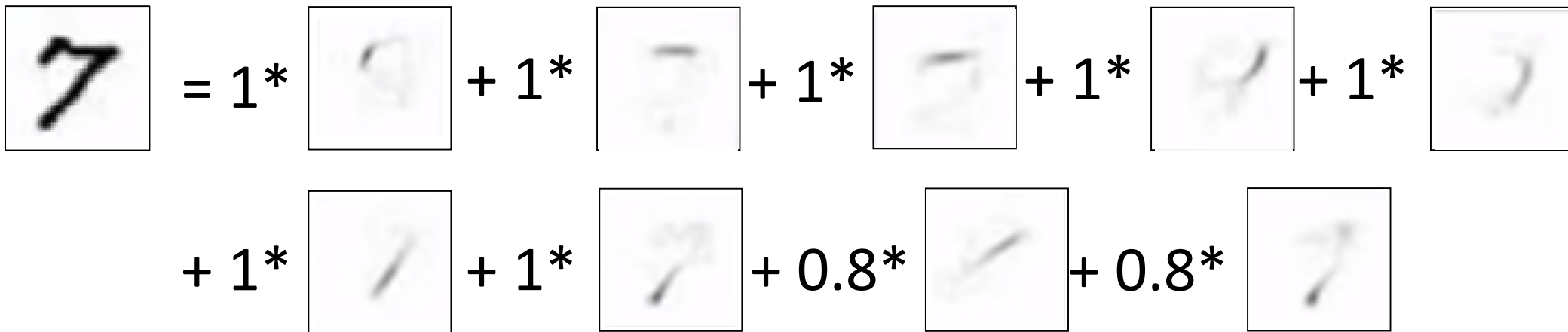Reconstructed image

Convolutional AE

# Regularization

Motivation:
- We would like to learn meaningful features **without** altering the code's dimensions (Overcomplete or Undercomplete).

The solution: imposing other constraints on the network.

# Sparsely Regulated Autoencoders

- We want our learned features to be as **sparse** as possible.
- With sparse features we can generalize better.

# Sparsely Regulated Autoencoders

$a_j$ is defined to be the activation of the $j$th hidden unit (bottleneck) of the autoencoder.

Let $a_j\ (x)$ be the activation of this specific node on a given input $x$.

# Sparsely Regulated Autoencoders

Further let,

$$\hat{\rho}_j = \frac{1}{m} \sum_{i=1}^{m} \left[ a_j \left( x^{(i)} \right) \right]$$

be the average activation of hidden unit $j$ (over the training set).

Thus we would like to force the constraint:

$$\hat{\rho}_j = \rho$$

where $\rho$ is a "sparsity parameter", typically small. In other words, we want the average activation of each neuron $j$ to be close to $\rho$.

# Sparsely Regulated Autoencoders

- We need to penalize $\hat{\rho}_j$ for deviating from $\rho$.
- Many choices of the penalty term will give reasonable results.

- For example: $$\sum_{j=1}^{Bn} KL\left(\rho | \hat{\rho}_j\right)$$
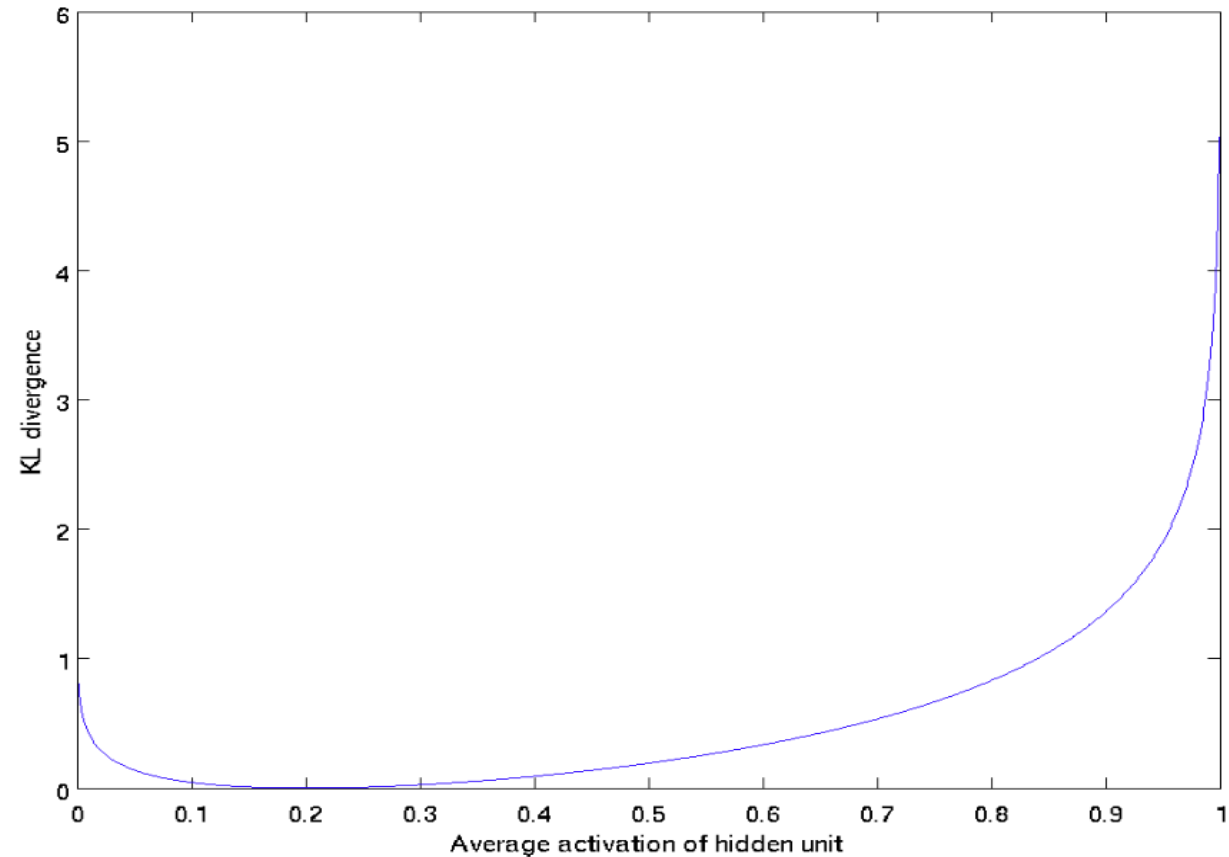
where $KL\left(\rho | \hat{\rho}_j\right)$ is a Kullback-Leibler divergence function.

# Sparsely Regulated Autoencoders

- A reminder:
  - KL is a standard function for measuring how different two distributions are, which has the properties:

$$KL\left(\rho | \hat{\rho}_j\right) = 0 \text{ if } \hat{\rho}_j = \rho$$

otherwise it is increased monotonically.



$\rho = 0.2$

# Sparsely Regulated Autoencoders

- Our overall cost functions is now:

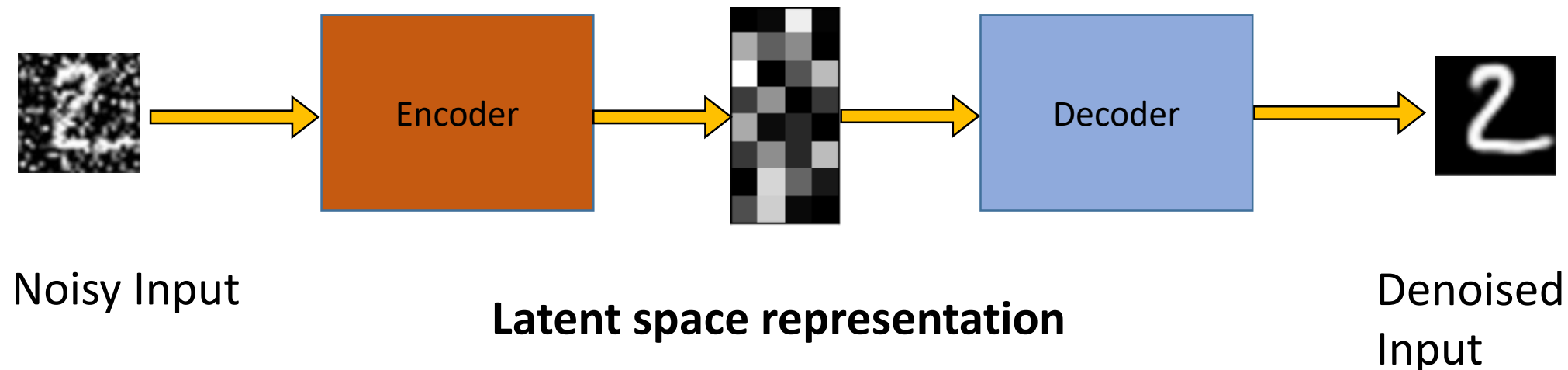$$J_S(W, b) = J(W, b) + \beta \sum_{j=1}^{Bn} KL(p|\hat{\rho}_j)$$

*Note: We need to know $\hat{\rho}_j$ before hand,
 so we have to compute a forward pass on all the training set.

# Denoising Autoencoders

**Intuition:**

- We still aim to encode the input and to NOT mimic the identity function.
- We try to undo the effect of *corruption* process stochastically applied to the input.
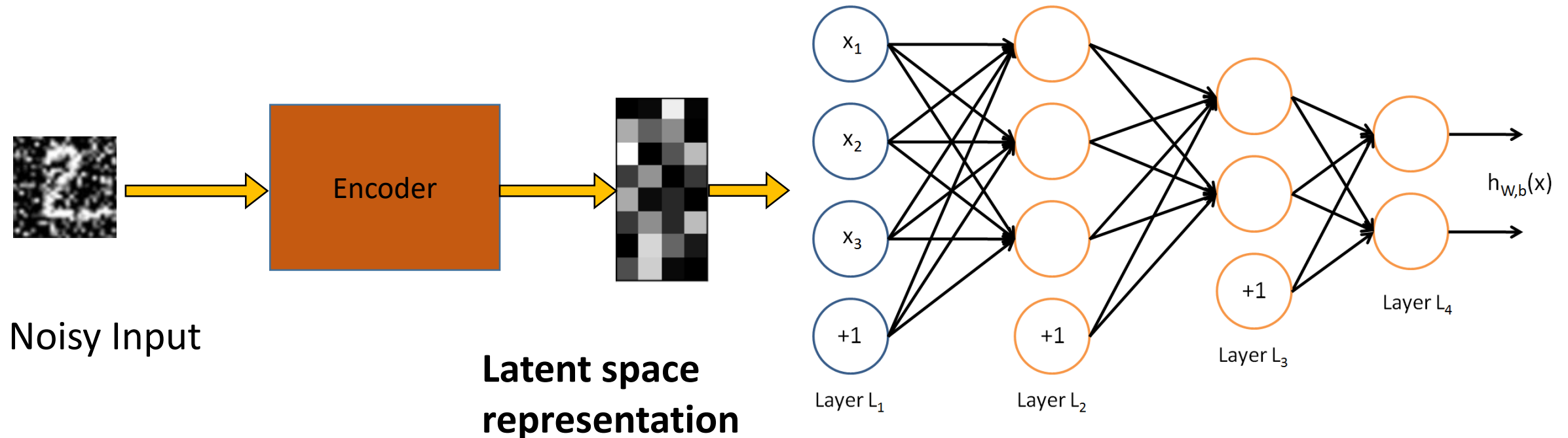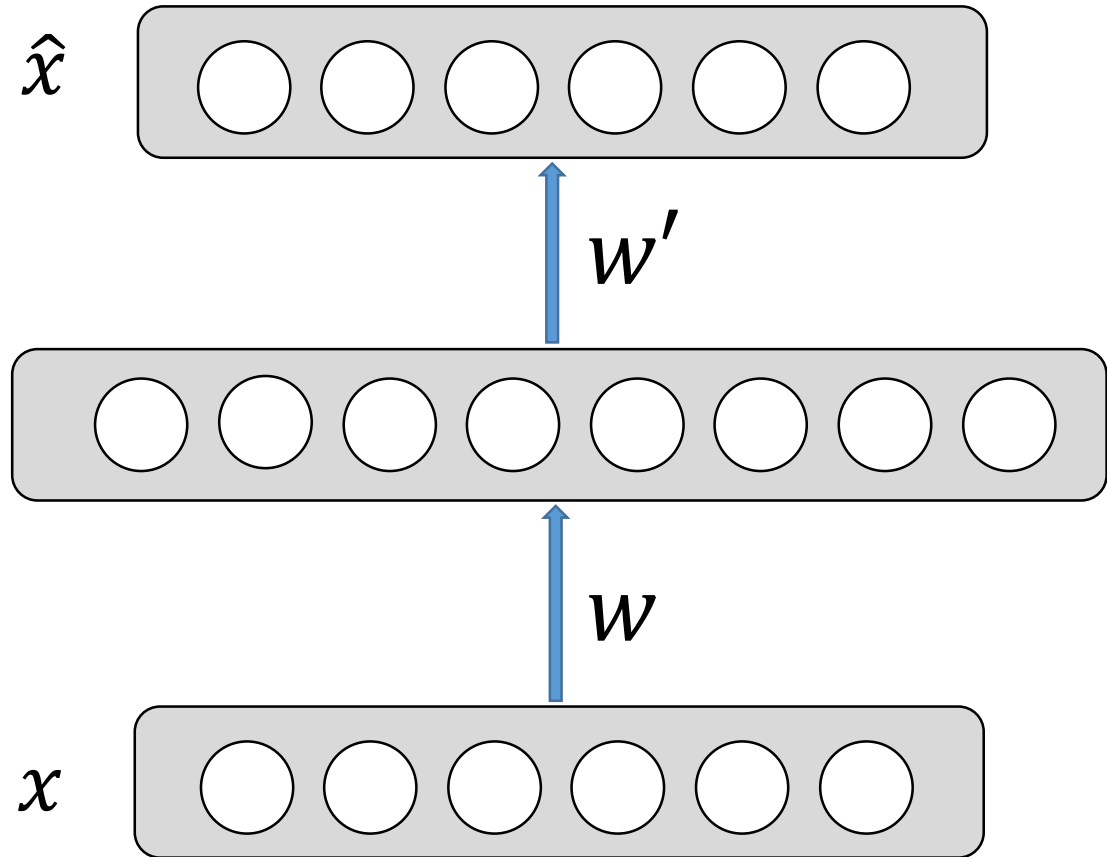
**A more robust model**



Noisy Input

**Latent space representation**

Denoised Input

# Denoising Autoencoders

Use Case:

- Extract robust representation for a NN classifier.



Noisy Input

**Latent space representation**

$x_1$

$x_2$

$x_3$

+1

Layer $L_1$

+1

Layer $L_2$

+1

Layer $L_3$

Layer $L_4$

$h_{W,b}(x)$

Encoder

# Contractive autoencoders

- We wish to extract features that **only** reflect variations observed in the training set. We would like to be invariant to the other variations.

- Points close to each other in the input space should maintain that property in the latent space.



$\hat{x}$

$w'$

$w$

$x$

# Contractive autoencoders

- Definitions and reminders:
- - Frobenius norm (L2): $\|A\|_F = \sqrt{\Sigma_{i,j}|a_{ij}|^2}$

- - Jacobian Matrix: $J_f(x) = \dfrac{\partial f(x)}{\partial x} = \begin{bmatrix} \dfrac{\partial f(x)_1}{\partial x_1} & \cdots & \dfrac{\partial f(x)_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \dfrac{\partial f(x)_m}{\partial x_1} & \cdots & \dfrac{\partial f(x)m}{\partial x_n} \end{bmatrix}$

# Contractive autoencoders

- Our new loss function would be:
- $$L^*(x) = L(x) + \lambda\Omega(x)$$

- where $\Omega(x) = \left\|J_f(x)\right\|_F^2$ or simply: $\sum_{i,j} \left(\frac{\partial f(x)_j}{\partial x_i}\right)^2$

and where $\lambda$ controls the balance of our reconstruction objective and the hidden layer "flatness".

$$Z_j = W_i X_i$$

$$h_j = \phi(Z_j)$$

$$\frac{\partial h_j}{\partial X_i} = \frac{\partial \phi(Z_j)}{\partial X_i}$$

$$= \frac{\partial \phi(W_i X_i)}{\partial W_i X_i} \frac{\partial W_i X_i}{\partial X_i}$$
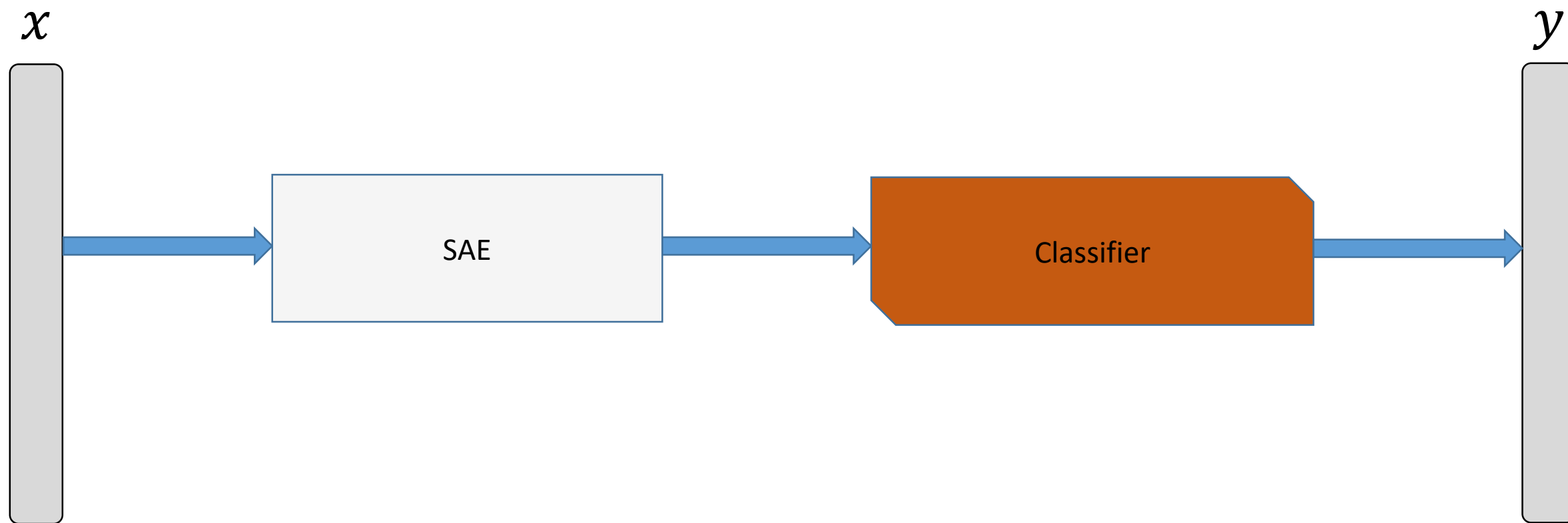
$$= \left[ \phi(W_i X_i)(1 - \phi(W_i X_i)) \right] W_i$$

$$= \left[ h_j(1 - h_j) \right] W_i$$

$$\frac{\partial h}{\partial X} = diag[h(1 - h)] W^T$$

$$\|J_h(X)\|_F^2 = \sum_{ij} \left( \frac{\partial h_j}{\partial X_i} \right)^2$$

$$= \sum_i \sum_j [h_j(1 - h_j)]^2 (W_{ji}^T)^2$$

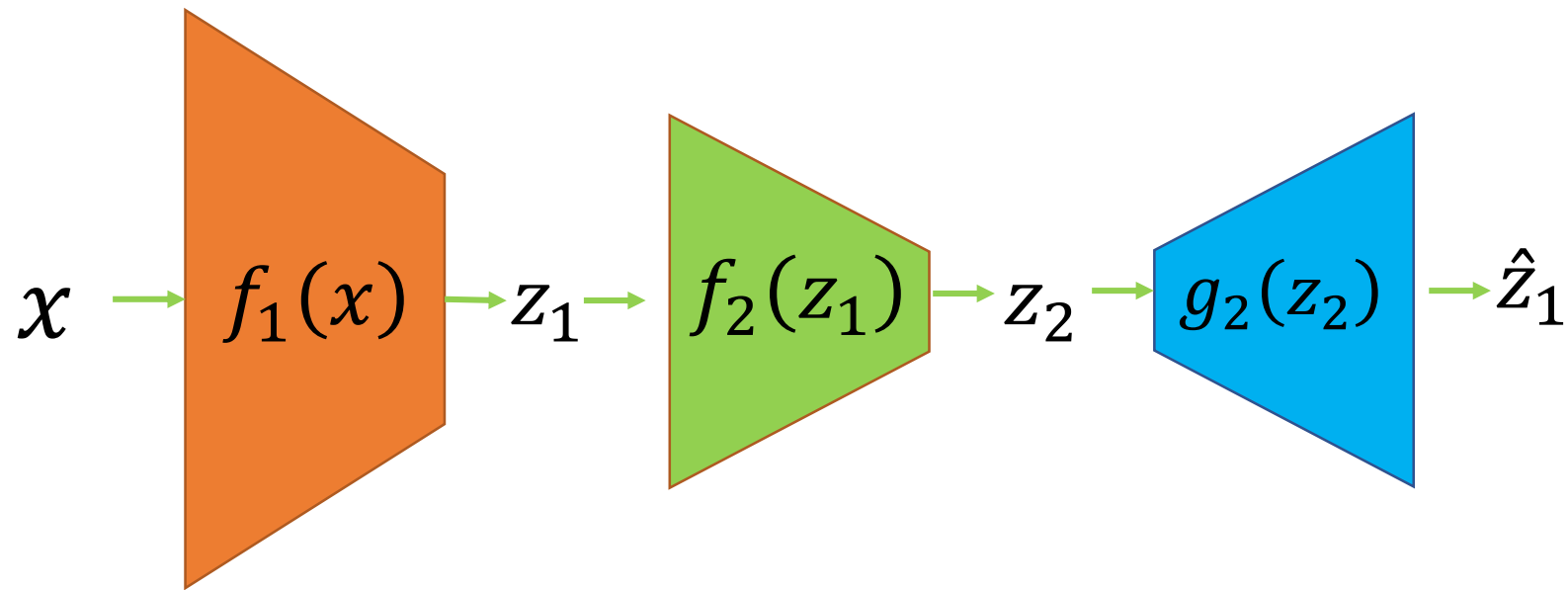$$= \sum_j [h_j(1 - h_j)]^2 \sum_i (W_{ji}^T)^2$$

# Stacked AE

$x$

$y$



SAE

Classifier

# Stacked AE – train process

First Layer Training (AE 1)

$$x \rightarrow f_1(x) \rightarrow z_1 \rightarrow g_1(z_1) \rightarrow \hat{x}$$

# Stacked AE – train process

Second Layer Training (AE 2)

$x \rightarrow \boxed{f_1(x)} \rightarrow z_1 \rightarrow \boxed{f_2(z_1)} \rightarrow z_2 \rightarrow \boxed{g_2(z_2)} \rightarrow \hat{z}_1$

# Stacked AE – train process

Add any classifier

$x$ → $f_1(x)$ → $z_1$ → $f_2(z_1)$ → $z_2$ → Classifier → Output

# Image segmentation - segnet

# Image segmentation U-net

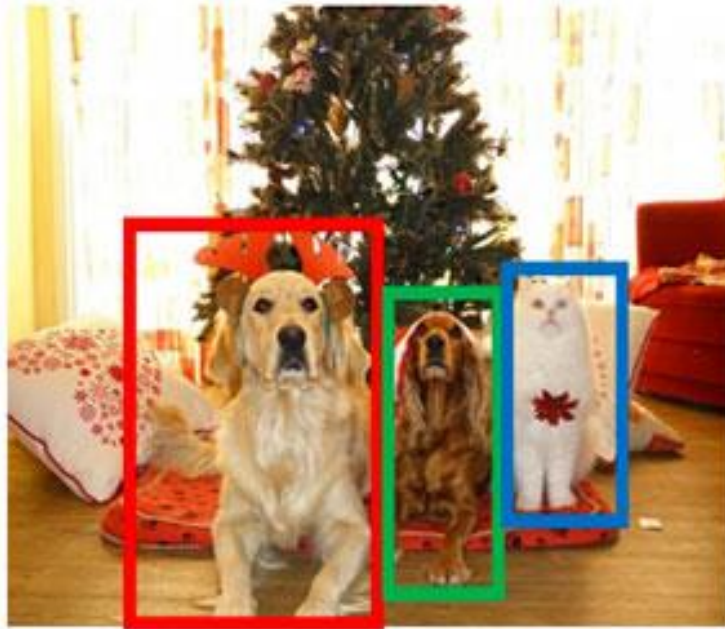# Image segmentation FCN

# Instance segmentation