# Autoencoders



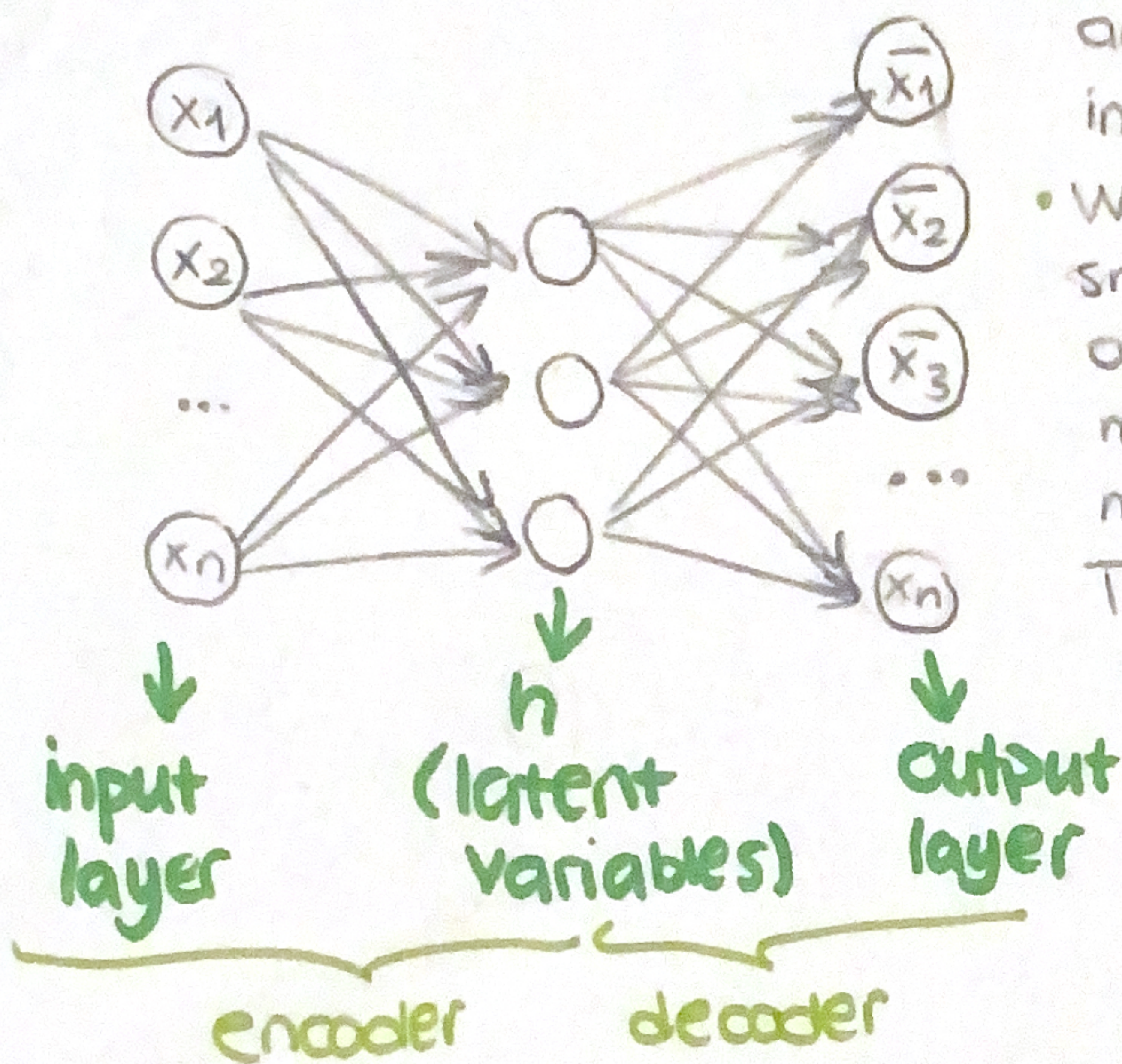**Idea:** Autoencoders are trained to reconstruct input.

- When hidden layer is smaller than input & output layers, the model only learns the most salient features. This is called undercomplete autoencoder.
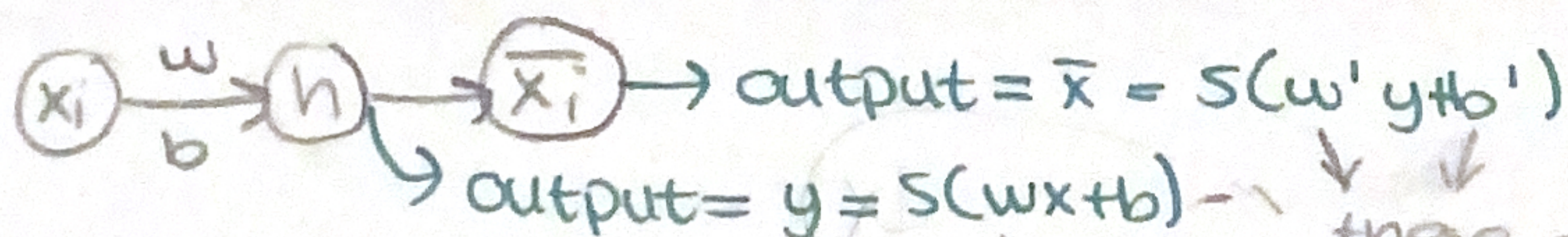
- Autoencoders learn unsupervised.

- $h_{w,b}(x) \approx x$
  $\downarrow$
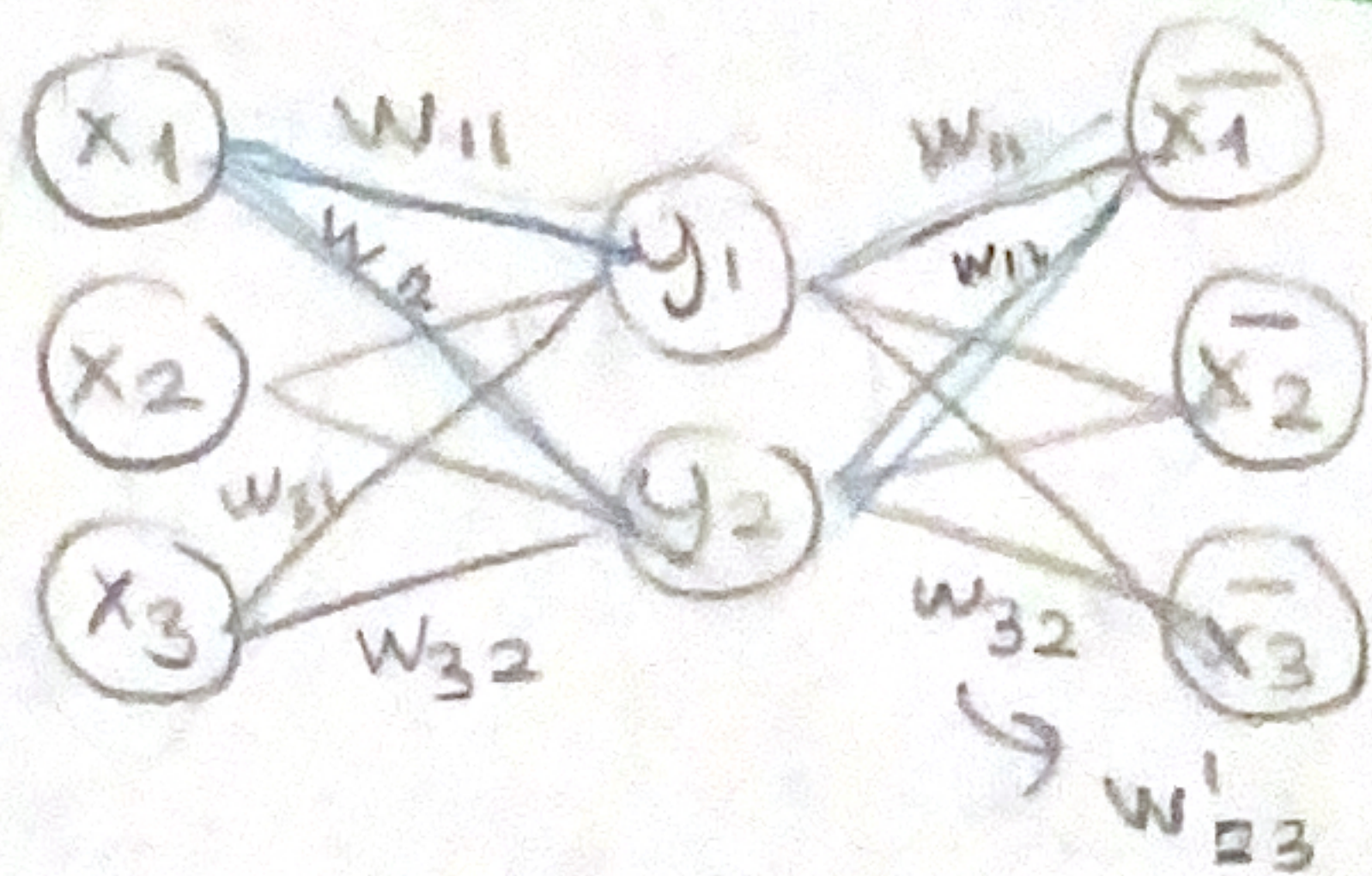  identity function (we're trying to learn this)

input layer  ·  h (latent variables)  ·  output layer

encoder  ·  decoder

$$LOSS \rightarrow (\bar{x_i} - x_i)$$

$x_i \xrightarrow[b]{w} h \rightarrow \bar{x_i} \rightarrow$ output $= \bar{x} = S(w'y + b')$

output $= y = S(wx + b)$

$L(x - J(h_{w,b}$

these are weights and biases of decoder

weights and biases of encoder

one is transpose of another!

$$w' = w^T$$

$$W = \begin{array}{c} \quad x_1 \quad\; x_2 \quad\; x_3 \\ \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \end{bmatrix} \end{array}$$

$$w' = w^T \rightarrow w' = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix}$$



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