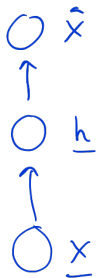


Representation Learning

Tirtharaj Dash

Dept. of CS & IS and APPCAIR
BITS Pilani, Goa Campus

November 27, 2021



maximising $P(\underline{x} | \underline{h})$

maximising $P(\underline{h} | \underline{x})$

> minimise the
reconstruction
loss

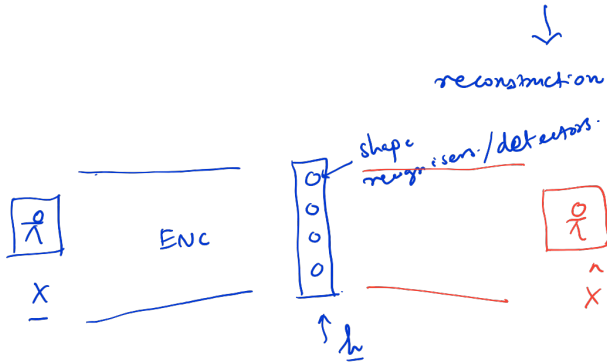
$$L(\underline{x}, \underline{\hat{x}})$$

$$\underline{x} \approx \underline{\hat{x}}$$

① MSE

② CE / NLL loss

Hinton "To recognise shapes, ^{first} learn to generate images."

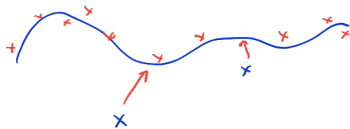


→ (1) sparse autoencoder.

$L + \lambda$ sparsity penalty \rightarrow penalise h

→ (2) Robust to noise: Denoising autoencoder.

Denoising AE



X

\downarrow

$X + \text{noise}$

\downarrow

$\sim \mathcal{N}(0, \sigma^2)$
||
 λ^2

Regularisation for Representation learning:

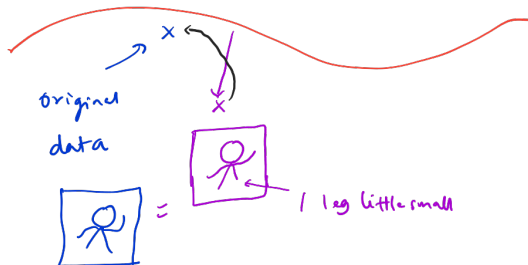
① , ②

③ Contractive Autoencoder : Penalise the gradient.

Contractive AE:

→ Method to penalise a sparse representation.

→ Robustness to perturbation in data.

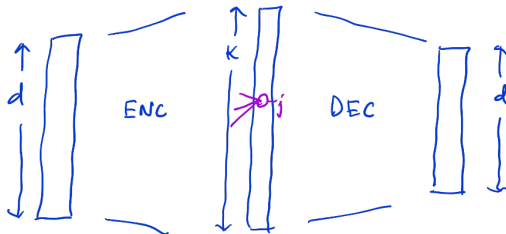


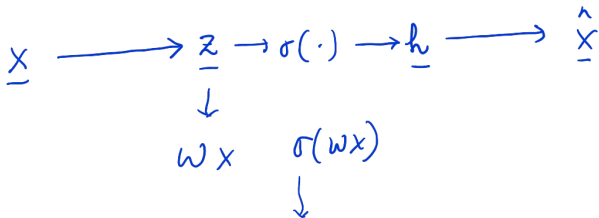
Let our reconstruction loss be:

$$L(\underline{x}, \underline{\hat{x}}) = \sum_{i=1}^d (x_i - \hat{x}_i)^2 \quad \left| \begin{array}{l} \underline{x} \in \mathbb{R}^d \\ \underline{\hat{x}} \in \mathbb{R}^d \\ \underline{h} \in \mathbb{R}^k \end{array} \right.$$

In Contrastive AE:

$$R = \frac{1}{2} \sum_{i=1}^d \sum_{j=1}^k \left(\frac{\partial h_j}{\partial x_i} \right)^2$$





$$\frac{\partial h_j}{\partial x_i} = w_{ji} \overbrace{h_j (1 - h_j)}^{\sigma'(w x)} \quad \forall i, j$$

Total loss

$$L_{\text{Total}} = \sum_{i=1}^d (x_i - \hat{x}_i)^2 + \frac{\lambda}{2} \sum_{j=1}^k h_j^2 (1 - h_j)^2 \sum_{i=1}^d w_{ji}^2$$

after skipping some steps

[HOMEWORK :-)]

$$\frac{\partial L_{\text{Total}}}{\partial w_{ji}} = x_i \frac{\partial L_{\text{Total}}}{\partial z_j} + \lambda w_{ji} h_j^2 (1 - h_j)^2$$

Task-specific
loss

shrinking the hidden rep.

[VERIFY]

(reconstruction
error should be
 $\rightarrow 0$)

L_2 -regularised AE

Q. Let hidden activation
is linear ?

$\lambda w_{ji} \underbrace{h_j^2}_{\text{circled}} \quad \|h\|^2$

Relationship between Denoising AE and Contractive AE:

achieves robustness stochastically.

$$\sim \mathcal{N}(0, \lambda^2)$$

↓
Gaussian

adding contractive regularization term.

↓
hidden activation is linear

↓
L₂-reg. AE.

$$|h| > |x|$$

Greedy, layerwise Pretraining using AE:

Without pretraining.

minimizing loss.



maximize the likelihood?

$$p(h|D) \propto p(D|h) p(h)$$

$$\approx \boxed{p(D|\underline{w})} p(\underline{w})$$

likelihood. prior

↓

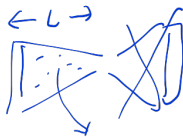
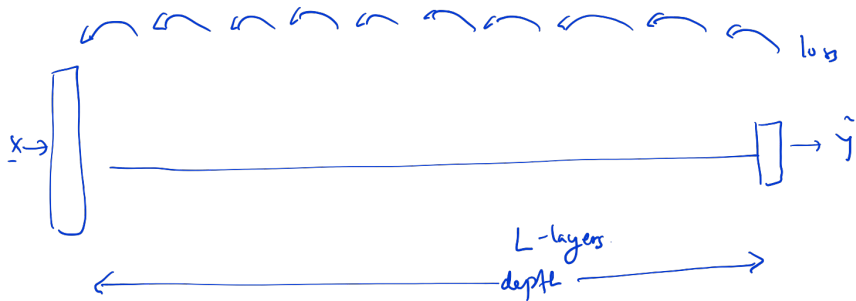
With Pretraining

$$p(D|\underline{w}) \boxed{p(\underline{w})}$$

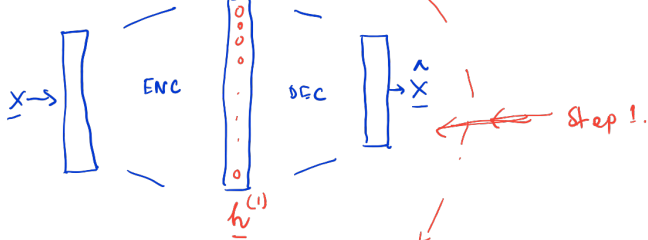
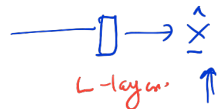
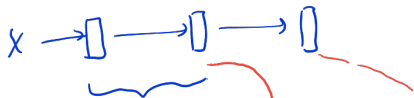


Known data.

↓
= source
=

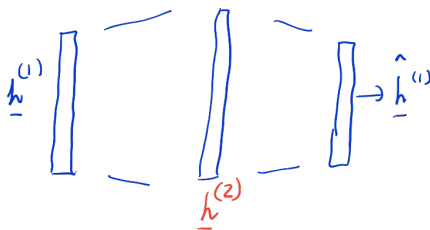


- ① resource
- ② computational issue
 - ① variable grad.
 - ⋮



$$\underline{\underline{p(\underline{w})}}$$

Training unsupervised
each layer
in a
greedy manner



Step 2

Next class:

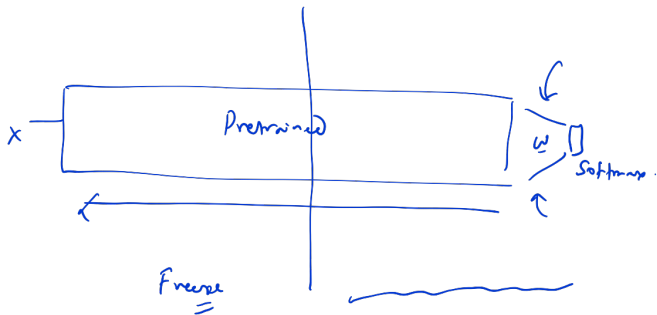
Learning "hidden"
"latent"

Probabilistic structure from data.

(VAE)

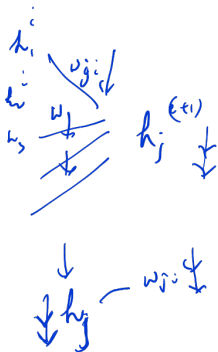
$p(\underline{x})$

(GAN)



$$\frac{\partial L_{\text{Total}}}{\partial w_{ji}} = \square + \lambda w_{ji} \frac{h_j^2 (1-h_j)^2}{\dots}$$

$$w_{ji} \leftarrow w_{ji} - \eta \left(\frac{\partial L_{\text{Total}}}{\partial w_{ji}} \right)$$



$$\leftarrow w_{ji} - \eta \left[\dots \right] \frac{h \cdot (1-h)}{\dots}$$

$$- \eta \left[\lambda w_{ji} \frac{h_j^2 (1-h_j)^2}{\text{large}} \right]$$