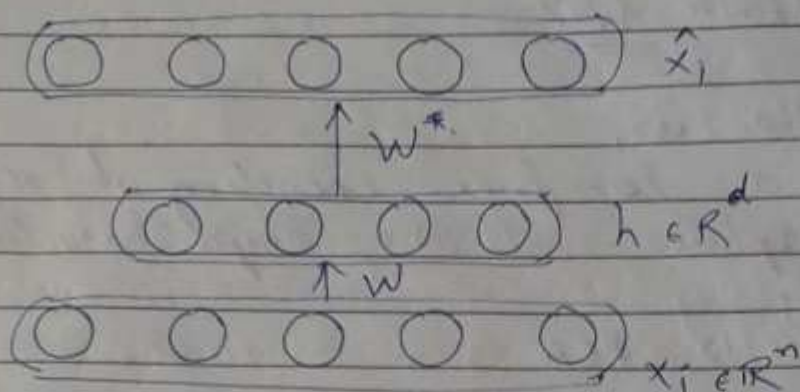


AutoEncoder

PAGE No.	
DATE	/ /



→ It's like a feed-forward NN.

① → It encodes its input x_i into a hidden representation h .

$$h = g(\underbrace{wx_i + b}_{\text{Linear transformation}})$$

can be any kind of sigmoid fn. (Linear or non-linear transformation)

② Decodes - the ip again from this hidden representation. & reconstructs the ip.

$$\hat{x}_i = f(w^*h + c)$$

→ x_i should be similar to \hat{x}_i → very similar to PCA.

1 Case

→ Let us consider the case when $\dim(h) < \dim(x_i)$

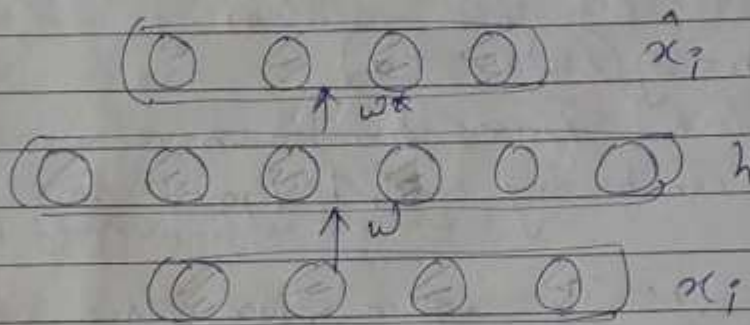
→ If we are still able to reconstruct

\hat{x}_i perfectly from h , then what does it say about h ?

- It tells us, " h is a loss-free encoding of x_i ."
- It captures all the imp. characteristics of x_i .
- Similar to PCA.

→ An encoder where, $[\dim(h) < \dim(x_i)]$ is called an under complete autoencoder.

Loss :- Consider, $\dim(h) \geq \dim(x_i)$



- In such a case, the autoencoder could learn a trivial encoding by simply copying x_i into h then copying h into \hat{x}_i .

eg

Compression → we have 10 bits
 ↓
 4 bits → (encoding info. present)

But 10 bits
 ↓
 16 bits (other 6 bits blank)
 ↓
 10 bits

PCA linearly transforms data into new coordinate sys^t along most of the variation in the data can be described with fewer dimensions.

PAGE NO.	
DATE	/ /

- ⇒ Such an identity encoding is useless in practice as it does not really tell us anything about the important characteristics of the data.
- ⇒ Such as autoencoder is ~~an~~ where, $\dim(h) \gg \dim(x_i)$ is called an over-complete autoencoder (hidden layer more neurons than i/p)
- E.g. whether a person is likely to get a disease or not
 - we want parameters like, height, wt & BMI
- we only computed BMI.
- we conclude at the end, that the person has high likelihood of getting diabetes depending on BMI.
- we want to know, whether it was the height or whether it was the wt which was responsible for this.
- So original i/p, features are entangled & we want to deentangle it.
(i.e. we want to go from smaller feature space to larger feature space)
- ⇒ so here, there is a need to prevent the identity fn(representation)

Regularization in Autoencoders

→ why we need Regularization?
to avoid overfitting or
to Enable Generalization

→ In case of over complete autoencoder
overfitting is likely.

→ Not generalized to unseen data

→ Overfitting - happens for more no. of
parameters

overcomplete autoencoders have
more no. of parameters

so we need Regularization

→ It can happen in undercomplete
autoencoders also.

→ As Eg. $\begin{bmatrix} 3 & 3 \\ 3 & 3 \\ 3 & 3 \end{bmatrix}$
 $100 \times 100 = 10K$

Lot of data is not imp., ~~here~~

After shrinking we are still having
large no. of parameters as, even

→ after shrinking also there is no
much redundancy in the itp that
shrinking still leads to a large
no. of parameters \therefore we could still
overfit. so we can need regularization

→ In case of over-complete Autoencoder

Solution:

The simplest solution is to add a "L2-regularization" term to the objective fn.

$$\min_{\theta, w, w^*, b, c} \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n (\hat{x}_{ij} - x_{ij})^2 + \lambda \|\theta\|^2$$

\Rightarrow all the parameters are weights & put it in large vector

$$\theta = (w_1, w_2, w_3, \dots)$$

$$\frac{\partial L(\theta)}{\partial \theta}$$

- This is easy to implement & just add a term λw to the gradient $\frac{\partial L(\theta)}{\partial w}$ (& similarly for other parameters)

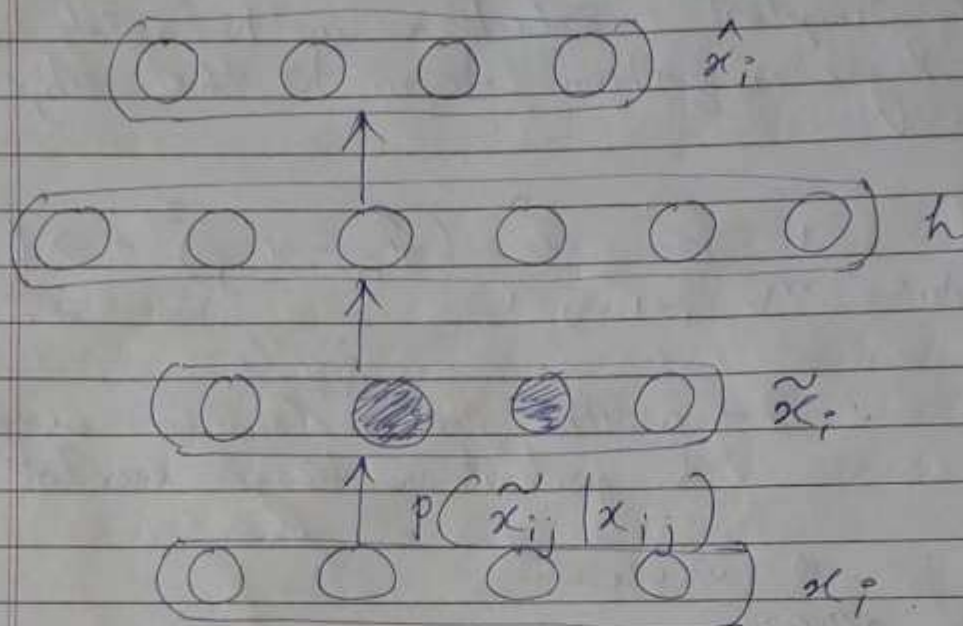


Another trick is to 'tie the weights' of the encoder & decoder i.e. $w^* = w^T$.
 \Rightarrow tying the wts - means - it reduces the parameter space of auto encoder. So it reduces the complexity of the model which releases overfitting.

\rightarrow parameter sharing: - we force certain weights to share the same value

- we divide the parameters or weights of a model into groups by leveraging prior knowledge, & all parameters in each group are constrained to take the same value.

Denoising Autoencoders (DAE)



→ A denoising encoder simply corrupts the input data using a probabilistic process ($P(\tilde{x}_{ij} | x_{ij})$) before feeding it to the network.

→ A simple $P(\tilde{x}_{ij} | x_{ij})$ used in practice is the following:-

$$P(\tilde{x}_{ij} = 0 | x_{ij}) = q$$

ie.

with probability q , it will set $\tilde{x}_{ij} = 0$ (corrupting) & with probability $(1-q)$ it is kept as it is (retaining data).
(Assume i/p are binary 0 & 1)

⇒ But we have corrupted data as
 i.e. Δ getting off as \hat{x}_i similar to x_i

→ i.e. the objective fn. is still to
 reconstruct the original (uncorrupted)
 x_p

$$\arg \min_{\theta} \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n (\hat{x}_{ij} - x_{ij})^2$$

Eg.



we corrupt, still it should
 predict it as '3'

So it should be more robust
 So it should not just memorize the
 data so to learn better we can
 feeding corrupted data

for eg:

it will have to learn to reconstruct a
 corrupted x_{ij} correctly by relying on
 its interactions with other elements of x_p

Eg.

Hand-written digit Recognition

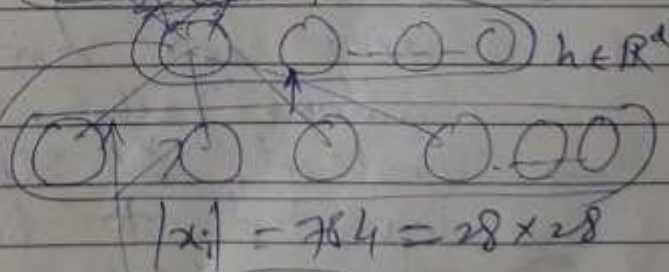
$\hat{x}_i \in \mathbb{R}^{784}$

3 1 8 5

8 4 1 5

8 0 1 9

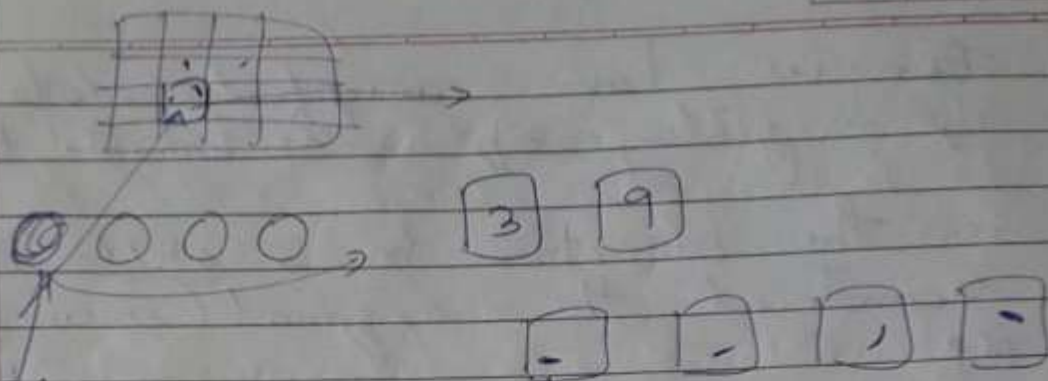
from MNIST Data



$$|x_i| = 784 = 28 \times 28$$

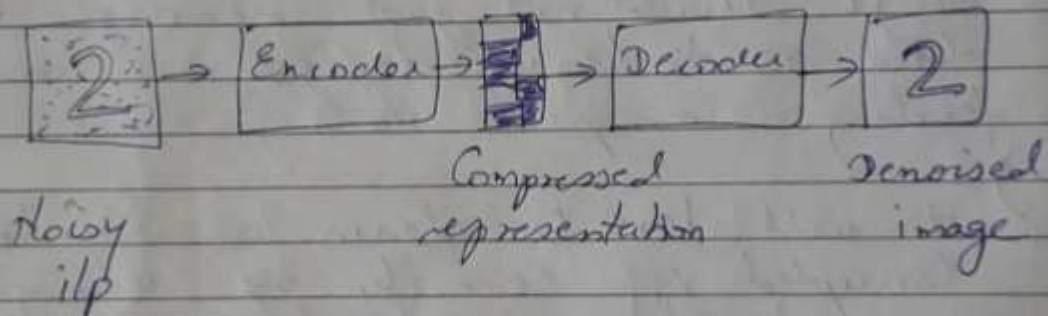
Captures interaction
 betw diff neurons





This neuron fires for this stroke,
other neurons for other strokes.
So together all neurons can decide
that the combination of these
strokes is 3 or 9.

★



- Denoising autoencoders are an extension of simple autoencoders; however it's worth noting that denoising autoencoders were not originally meant to automatically denoise an image.
- Instead, this procedure was invented to help :-
- ① The hidden layers of the autoencoder learn more robust filters.
 - ② Reduce the risk of overfitting in the autoencoder.

③ prevent the autoencoder from learning a simple identity function.

→ noise was stochastically (i.e. randomly) added to the input data & then the autoencoder was trained to recover the original, nonperturbed signal.

Eg. preprocessing an image to improve the accuracy of an Optical Character recognition (OCR) algorithm.

② Ensures - that it is actually learning latent representations sparse Auto Encoders: instead of redundant info. in our inp data.

- A hidden neuron with sigmoid activation will have values between 0 & 1.
- we say that neuron is activated when its o/p is close to 1 & not activated when o/p is 0.

* A sparse Encoder tries to ensure the neuron is inactive most of the times (i.e. it is close to 0 for most of the ilps).

i.e. the Average activation of a neuron is close to 0.

→ If the neuron l is sparse (i.e. mostly inactive) then $\hat{p}_l \rightarrow 0$.

- The 'average value of the activation' of a neuron l is given by

$$\hat{p}_l = \frac{1}{n} \sum_{i=1}^n h(x_i)_l$$

2 diff ways to constraint sparsity penalty:
 ① L_1 -regularization (Lasso) \rightarrow puts each feature to 0
 ② KL-divergence

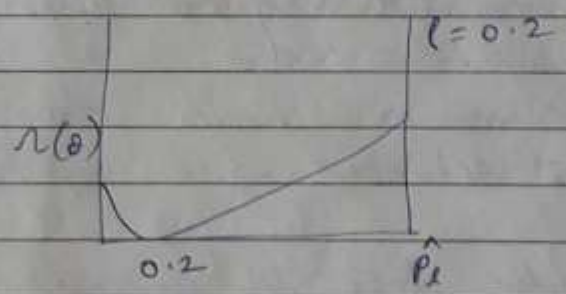
\Rightarrow A sparse autoencoder uses a sparsity parameter ρ (rate) (typically \approx very close to 0, say 0.005) & tries to enforce the constraint $\hat{\rho} = \rho$.

\Rightarrow To ensure $\hat{\rho} = \rho$, we add the term

$$r(\theta) = \sum_{k=1}^K \left[l \log \frac{l}{\rho_k} + (1-l) \log \frac{1-l}{1-\rho_k} \right]$$

$$L'(\theta) = L(\theta) + r(\theta)$$

\Rightarrow this takes its minimum value when $\hat{\rho}_k = \rho$



- for a single k , we take $\rho = 0.2$
 we reach to zero only when $\hat{\rho}_k = \rho$

$\left\{ \begin{array}{l} - L(\theta) \text{ is the squared error loss or cross entropy loss} \\ \& r(\theta) \text{ is the sparsity constraint} \end{array} \right.$

Finally, $\frac{\partial L'(\theta)}{\partial w} = \frac{\partial L(\theta)}{\partial w} + \frac{\partial r(\theta)}{\partial w}$

③ Contractive Autoencoders:-

- A contractive autoencoder also tries to prevent an overcomplete autoencoder from learning the identity function.
- It does so, by adding the following regularization term to the loss fun;

$$\alpha(\theta) = \|J_x(h)\|_F^2$$

- where, $J_x(h)$ is the Jacobian of the encoder.
- If the ip has n-dimensions & the hidden layer has k dimensions then, the Jacobian Matrix is,

$$J_x(h) = \begin{bmatrix} \frac{\partial h_1}{\partial x_1} & \dots & \dots & \dots & \frac{\partial h_1}{\partial x_n} \\ \frac{\partial h_2}{\partial x_1} & & & & \frac{\partial h_2}{\partial x_n} \\ \vdots & & & & \vdots \\ \frac{\partial h_k}{\partial x_1} & & & & \frac{\partial h_k}{\partial x_n} \end{bmatrix}$$

- * its a partial derivative of every neuron in the 1st hidden layer w.r.t the 1st ip.
- ie. Jacobian is derivative of h w.r.t x
- ⇒ Each J_x^{th} entry shows:- how much does the change with a small change in x_k

Similar to Euclidean norm

- Frobinious Norm -

Contracting Autoregressive

$$\left\{ \|I_x(h)\|_F^2 = \sum_{j=1}^n \sum_{k=1}^K \left(\frac{\partial h_j}{\partial x_k} \right)^2 \right\} \text{ Should be '0'}$$

it is sum of squares of all elements of a matrix

we want to Minimize $L(\theta) + \underbrace{r(\theta)}_{\text{frobinius norm} = 0}$ if $r(\theta) = 0$, $\text{rank}(\theta)$ will be very high

Consider, $\frac{\partial h_j}{\partial x_k} = 0$

→ It means that this neuron is not very sensitive to variations in the ip x_k

But, doesn't this contradict our other goal of minimizing $L(\theta)$ which requires h to capture variations in the ip.

→ $L(\theta)$ - wants to capture the variations

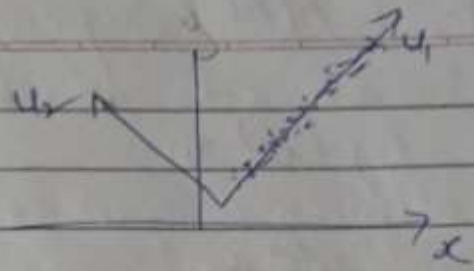
⇒ By putting these 2 contradicting objectives against each other we ensure that h is sensitive to only very important variations as observed in the training data.

→ $L(\theta)$ - Capture important variations in data

→ $r(\theta)$ - do not capture variations in data

⇒ Tradeoff - Capture only very important variations in the data

Eg.



- u_1 & u_2 are the dimensions.
 - u_1 is important - as it has lots of variance
 - $u_2 \rightarrow$ variance is out of line which doesn't seem important
- \rightarrow Its similar to PCA \rightarrow try to capture the variations across the important dimensions but not across non-important ones