# Lecture 7

## 7.1

Introduction to Autoencoders



$\hat{x_i}$ (output)

$w^* \uparrow$

$h$ (hidden)

$w \uparrow$

$x_i$ (input)

Special type of feed forward $\not\!f\!f$ neural network

Encodes i'th input $x_i$ into a hidden representation $h$

Encoder function

$$h = g(Wx_i + b)$$

$x_i \in \mathbb{R}^n \qquad h \in \mathbb{R}^d \qquad W \in \mathbb{R}^{d \times n}$

$b \in \mathbb{R}^{d \times n}$

$h \in \mathbb{R}^d$

$$\hat{x_i} = f(W^* h + c) \qquad W \in \mathbb{R}^{n \times d}$$

$x_i \in \mathbb{R}^n$

Decoder, decodes the input again from the hidden representation

The model is trained to minimize a certain loss function which will ensure that $\hat{x_i}$ is close to $x_i$

if $dim(h) < dim(x_i)$

if we are able to reconstruct $\hat{x_i}$ perfectly from $x_i$, $h$, then $h$ is a loss-free encoding of $x_i$. It captures all the important characteristics of $x_i$.

Analogy with PCA.

$dim(h) \geq dim(x_i)$

Identity encoding is useless.

→ Choice of $f(x_i)$ & $g(x_i)$

→ Choice of loss function

If i/p binary → function for decoder is logistic function

if i/p real number,
the function will be linear

$g$ is typically chosen as sigmoid.

## Loss function

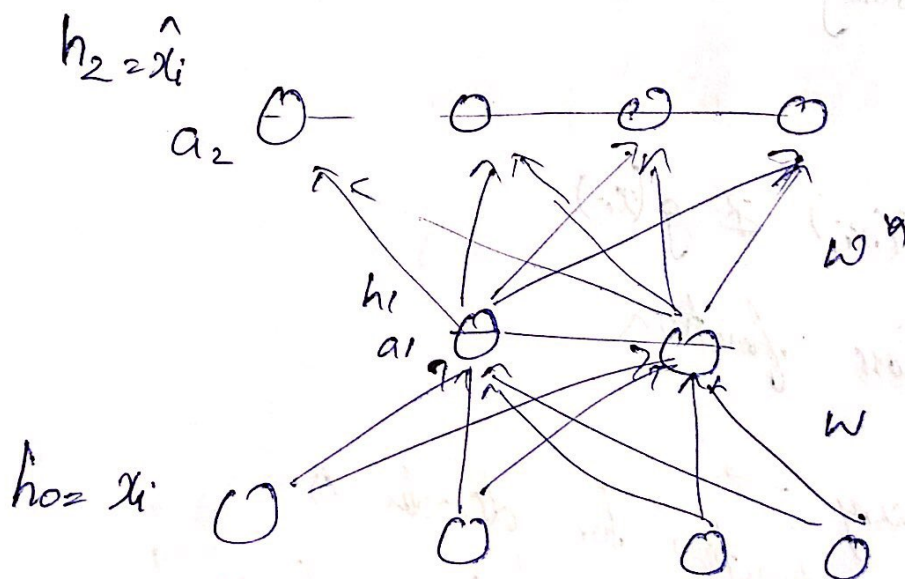i/p are real-valued.

Squared error loss

We can train the autoencoder just
like a regular feedforward network.

$$\frac{\partial \mathcal{L}(\theta)}{\partial W^*} \quad \& \quad \frac{\partial \mathcal{L}(\theta)}{\partial W}$$

$$\mathcal{L}(\theta) = (\hat{x_i} - x_i)^T (\hat{x_i} - x_i)$$

$h_2 = \hat{x_i}$

$$\frac{\partial L(o)}{\partial w^{v}} = \frac{\partial L(o)}{\partial h_2} \left[ \frac{\partial h_2}{\partial a_2} \cdot \frac{\partial a_2}{\partial w^a} \right]$$

$$\frac{\partial L(o)}{\partial w} = \frac{\partial L(o)}{\partial h_2} \left[ \frac{\partial h_2}{\partial a_2} \frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a_1} \frac{\partial a_1}{\partial w} \right]$$

$$\frac{\partial L(o)}{\partial h_2} = \frac{\partial L(o)}{\partial \hat{x_i}}$$

$$= \nabla_{\hat{x_i}} \left\{ (\hat{x_i} - x_i)^T (\hat{x_i} - x_i) \right\}$$

$$= 2 (\hat{x_i} - x_i)$$

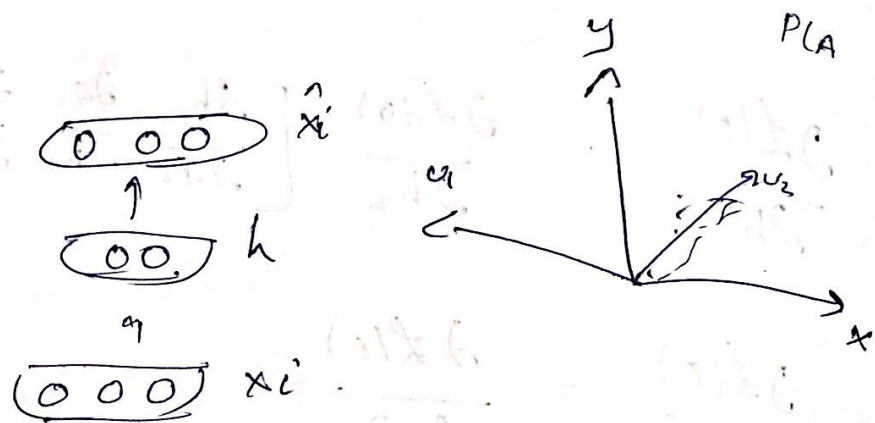Input — Binary    Logistic

→ Sigmoid deroduce output b/w 0 & 1

$$\min = \left\{ - \sum_{j=1}^{n} (x_{ij} \log \hat{x}_{ij} + (1 - x_{ij}) \log (1 - \hat{x}_{ij})) \right\}$$

$$P = [1, 0]$$

$$q = [0.8, 0.2] \qquad -\sum_{i=0} P_i \log q_i$$

$$q_0 \quad q_1$$

Cross Entropy Function
for input is binary

## 7.2 Link b/w PCA & Auto Encoder



The encoder part of an autoencoder is equivalent to PCA if we

→ use a linear encoder

→ use a linear decoder

→ using squared error loss

→ normalize the inputs to

$$\hat{x}_{ij} = \frac{1}{\sqrt{m}} \left( x_{ij} - \frac{1}{m} \sum_{k=1}^{M} x_{kj} \right)$$

$$(x^T) x = \frac{1}{m} (x')^T x'$$

~~Optimal~~

Linear Decode & ^Error loss function

Squared

the the optimal solution to the following
objective function

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

when we use of a linear encoder

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

This is equivalent to

$$\min_{W\Theta H} \left( \| X - HW^* \|_F \right)^2$$

$$HW^* = U_{,\leq k}$$

H is actually a linear encoding & find
an expression for the encoder weights W.

$$H = U_{.,\leq k} \Sigma_{k,k} \doteq (XX^T)(XX^T)^{-1} U_{,\leq k} \Sigma_{k,k}$$

$$= XX I$$

7.3 Regularization in Autoencoders

to generalization

~~training time~~

the

if less generalization (then overfitting)

Regularization is done to avoid overfitting

if the no of parameters are high,
it leads to overfitting

Over complete autoencoder a regularization
is required as the number of parameters are very
high) . But there are situations where we need
regularization for under complete auto encoder

Simplest solution add L2 regularization

$$\left[ \min_{\theta, w, w^*, b, c} \frac{1}{m} \sum_{i=1}^{m} \sum_{j=1}^{n} (\hat{x}_{ij} - x_{ij})^2 \right] + \lambda \|\theta\|^2 \quad (L2)$$
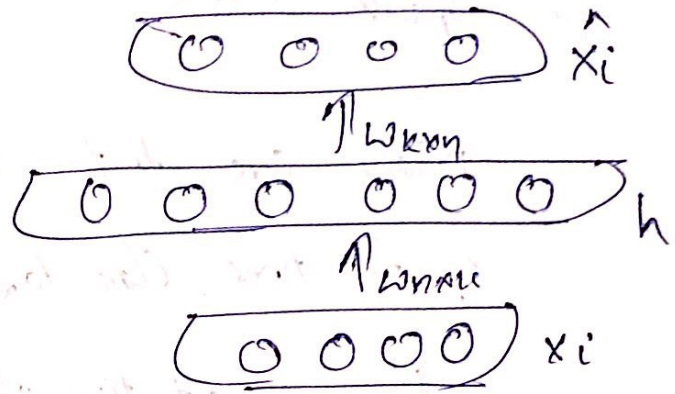
objective function $L'(\theta)$ $\quad \Omega(\theta)$

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

$$\frac{\partial L'(\theta)}{\partial w_1} \quad \text{adds the term } \alpha \lambda w$$
$$\text{to the gradient}$$
$$\partial \Omega(\theta)$$

$$W^* = W^T$$



Network diagram with layers: top layer $\hat{x_i}$, middle layer $h$ with $W_{key}$ and $W_{linear}$ connections, bottom layer $x_i$.
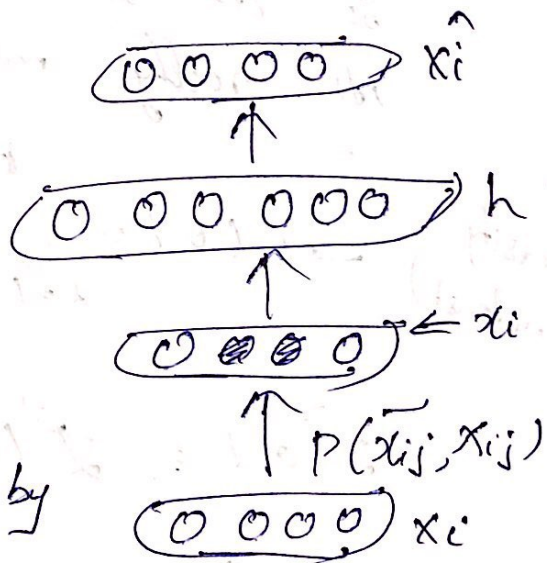
## 7.4 Denoising Autoencoders

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

[ Robust to change at test times ]



Network diagram: top layer $\hat{x_i}$, layer $h$, corrupted layer $\leftarrow x_i$ with $P(\widetilde{x_{ij}}, x_{ij})$, bottom layer $x_i$.

No longer get always by memorising the data..

• Or add Gaussian Noise
$$\widetilde{x_{ij}} = x_{ij} + N(0,1)$$

Gaussian Noise is better than $L2$ regularisation

## 7.5 Sparse Encoder

Make sure that neuron is inactive most of the time close to zero

Average activation of neuron is close to zero.

Average value of activation of neuron $\Rightarrow$ $P_i^\wedge = \dfrac{1}{m} \sum_{i=1}^{m} h(x_i) d e$

This is kind of regularization, So its like whenever its active, its really meaningful patterns.

(Sparsity constant) $\Omega(\theta) = \sum_{l=1}^{k} p \log \dfrac{P}{P_i^\wedge} + (1-p) \log \dfrac{1-p}{1-P_i^\wedge}$

$$\mathcal{L}(\theta) = \mathcal{L}'(\theta) + \Omega(\theta)$$

$$\Omega(\theta) = \sum_{l=1}^{k} p \log \left(\dfrac{p}{P_i^\wedge}\right) + (1-p) \log \dfrac{(1-p)}{(1-P_i^\wedge)}$$

# 7.6

## Contractive Autoencoders

Some kind of regularization

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

(Jacobian) — $\mathbb{R}^k$

— $\mathbb{R}^n$

$$\| J_x(h) \|_F^2 = \sum_{j=1}^{n} \sum_{\ell=1}^{k} \left( \frac{\partial h_\ell}{\partial x_j} \right)^2$$

Not very sensitive to variations.

### Regularization

Weight decay

$$\Omega(\theta) = \lambda \| \theta \|^2$$

Sparse

$$\Omega(\theta) = \sum_{\ell=1}^{k} p \log \frac{p}{\hat{p}_\ell} + (1-p) \log \frac{1-p}{1-\hat{p}_\ell}$$

Contractive

$$\Omega(\theta) = \sum_{j=1}^{n} \sum_{\ell=1}^{k} \left( \frac{\partial h_1}{\partial x_j} \right)^2$$