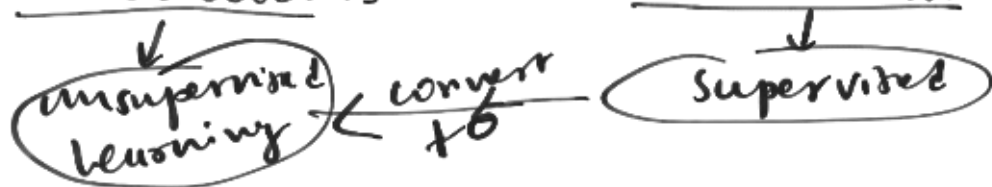


Sparse Autoencoders

Neural Networks



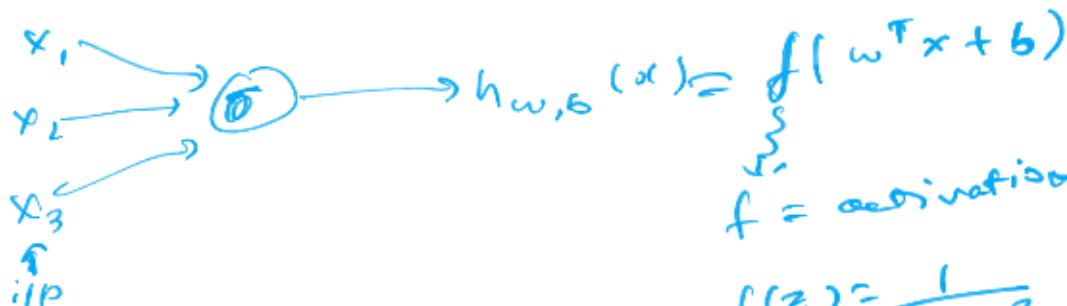
Neural Netw:

$(x^{(i)}, y^{(i)})$

Hypothesis: $h_{w,b}(x)$

parameters: w, b

"Neuron"



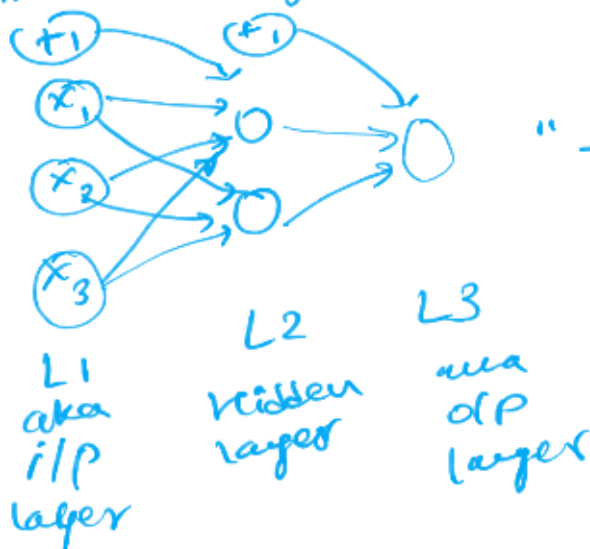
f = activation function

$$f(z) = \frac{1}{1 + e^{-z}} \Rightarrow \text{sigmoid AF.}$$



$$f(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Combine lots of these neurons to make a neural network



" $+1$ " \Rightarrow bias units ' b '

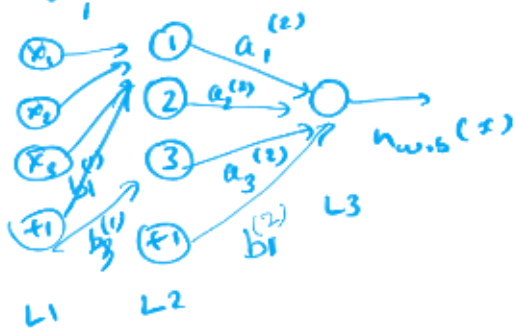
... : w, b

parameters: w, b

$w_{ij}^{(l)}$ = weight for connection b/w unit j in layer l and unit i in layer $l+1$

$b_i^{(l)}$ = bias associated with unit i in layer $l+1$

$a_i^{(l)}$ = activation (O/P value) of unit i in layer l



$$a_1^{(2)} = f(w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2 + w_{13}^{(1)} x_3 + b_1^{(1)})$$

$$a_2^{(2)} = f(w_{21}^{(1)} x_1 + w_{22}^{(1)} x_2 + w_{23}^{(1)} x_3 + b_2^{(1)})$$

$$a_3^{(2)} = f(\dots)$$

$$h_{w,b}(x) = f(w_{11}^{(2)} a_1^{(2)} + w_{12}^{(2)} a_2^{(2)} + w_{13}^{(2)} a_3^{(2)} + b_1^{(2)})$$

$h_{w,b}(x) \Rightarrow$ Parameters w, b

Online learning: $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots \uparrow$ no end

Stochastic gradient descent:-

For $i = 1, 2, 3, \dots$

Get example $(x^{(i)}, y^{(i)})$

update $w_{jk}^{(l)} := w_{jk}^{(l)} - \alpha \frac{\partial}{\partial w_{jk}^{(l)}} J(w, b; x^{(i)}, y^{(i)})$

$b_i^{(l)} := b_i^{(l)} - \alpha [\dots]$

$$b_j := v_j$$

Cost Function:- $J(w, b; x, y)$

$$= \underbrace{\frac{1}{2} \| h_{w,b}(x) - y \|^2}_{\text{Squared error}} + \underbrace{\frac{\lambda}{2} \sum_l \sum_i \sum_j (w_{ij}^{(l)})^2}_{\text{weighted decay}}$$

regularization term to prevent overfitting

λ = weight decay parameter

If $(x^{(i)}, y^{(i)}) \sim \mathcal{D}(\text{IID})$ \hookrightarrow Identically independently distributed.

$$\min_{w, b} E[J(w, b; x, y)]$$

For a NN it is NOT ok to initialize all parameters with zeroes

So instead we initialize with random values

$$w_{ij}^{(l)} \sim \mathcal{N}(0, \xi^2) \quad \{ \xi \}$$

\hookrightarrow Normal distribution

But if we do so, all hidden units will be computationally exactly the same because in each iteration i.e., $a_1^{(2)} = a_2^{(2)} = a_3^{(2)} = \dots$

\therefore weights randomly for symmetry

Breuninger

So we init
Bias values can be initialized with 0s or
randomly.

Backpropagation :-

Intuition: For node i in layer l
compute $\delta_i^{(l)}$ that measure how
responsible that node is for errors
in o/p.

\Rightarrow For o/p node i :-
(in layer n_l)
compute $\delta_i^{(l)}$ as a function of
 $h_{wi}^{(l)}(x)$ & y

\Rightarrow For hidden units i :- can't observe the o/p directly
 $\delta_i^{(l)} = \left(\sum_j w_{ji}^{(l)} \delta_j^{(l+1)} \right) \cdot f'(z_i^{(l)})$

Algorithm:-

- (1) Feedforward pass to compute activations of all units
- (2) For each unit in o/p layer (layer n_l)
compute $\delta_i^{(n_l)} = -(y_i - a_i^{(n_l)}) \cdot f'(z_i^{(n_l)})$
- (3) For each layer $l = n_l - 1, n_l - 2, \dots, 2$
 $\delta_i^{(l)} = \left(\sum_j w_{ji}^{(l)} \delta_j^{(l+1)} \right) \cdot f'(z_i^{(l)})$

(4) Update parameters:-

$$w_{ij}^{(l)} := w_{ij}^{(l)} - \alpha (a_j^{(l)} \delta_i^{(l+1)} + \lambda w_{ij}^{(l)})$$

$$b_i^{(l)} := b_i^{(l)} - \alpha \delta_i^{(l+1)}$$

Unsupervised learning:-

$x^{(1)}, x^{(2)}, x^{(3)}, \dots$

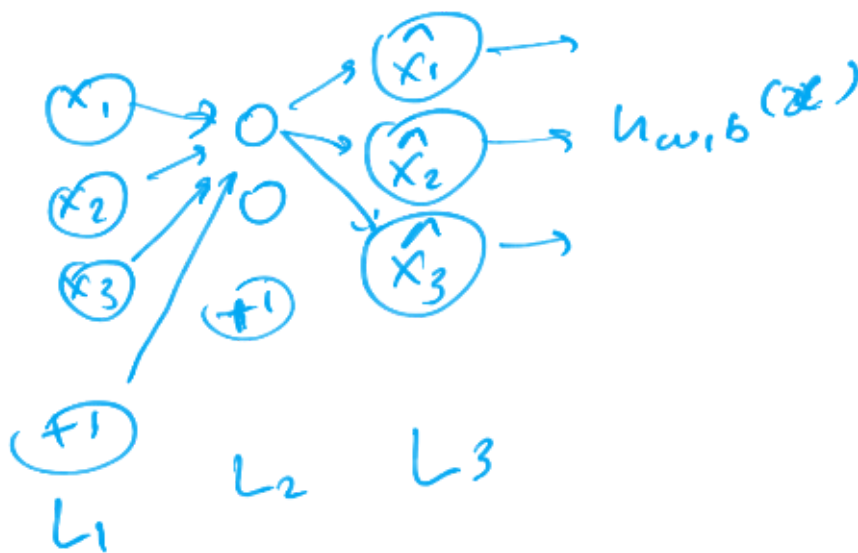
create $y^{(i)} = x^{(i)}$

so that we can give it to neural net. NN need labelled examples.

Autoencoder:-

$$h_{w,b}(x) \approx x$$

learn identity function.
approximation to ID. function.



eg: if 210×10 image

$$\Rightarrow x \in \mathbb{R}^{100}$$

50 hidden units (say)

100 O/P units

In a real image all pixels are not completely independent (i.e., not i.i.d.)

... to compress 100 pixels image do a

If we draw 50 units
 50 units hidden layer then are well be
 learning many redundancies, so are more
 in 10 units ~~not~~ not a meaningful
 compression.

→ limiting the size of hidden layer is
 one way to force autoencoders to learn an
 interesting function.

→ another constraint \rightarrow sparsity

Sparsity (Informally)

constraint $a_i^{(2)} \approx -1 \quad \forall i$

Neuron is "active" if $a_i^{(2)} \approx +1$

"inactive" if $a_i^{(2)} \approx -1$

\Rightarrow most neurons should be inactive
 most of the times.

Sparsity (Formal defn)

Imagine $x^{(1)}, x^{(2)}, x^{(3)}, \dots$ are drawn IID from \mathcal{D}

constraint: $E_{x \sim \mathcal{D}}[a_i^{(2)}] = \rho \quad \forall i$

where $\rho = -0.9$ (say) {i.e., close to -1}

Algo :-

(1) keep running estimate of $\hat{\rho}_i$ of $E_{x \sim \mathcal{D}}[a_i^{(2)}]$

(2) On each iteration of gradient descent, update

parameters to make expectation closer to $\rho = -0.9$ (say)

Explanation:-

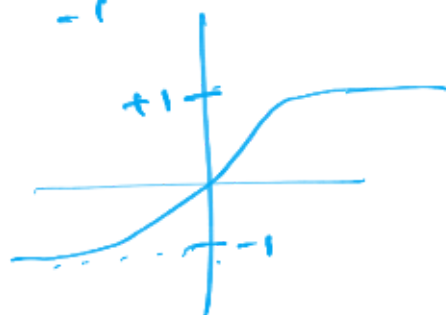
On each iteration, update
 (1) $\hat{p}_i := 0.999 \cdot \hat{p}_i + 0.001 \cdot a_i^{(2)}$ {for each i }
 Initialize $\hat{p}_i = 0$ (say) (eq 1)

(2) Recall $a_i^{(2)} = f\left(\sum_{j^{(1)}} w_{ij}^{(1)} x_j + b_i^{(1)}\right)$

If $\hat{p}_i > p \Rightarrow$ we want hidden unit i to be less active
 eg $(-0.7) > (-0.9)$ (want o/p values closer to -1)

So decrease $b_i^{(1)}$

bcz decreasing $b_i^{(1)}$ will move value of 'f'
 closer to -1



If $\hat{p}_i < p$

\downarrow
 we want hidden unit i to be more active
 So increase $b_i^{(1)}$

putting both cases together

update rule:

$$b_i^{(1)} := b_i^{(1)} - \underbrace{\alpha}_{\text{learning rate}} \underbrace{\beta (\hat{p}_i - p)}_{\text{additional learning rate}} \quad \forall i \quad \text{(eq 2)}$$

Summarised algo. for enforcing sparsity constraint
→ On each iteration (get new (x, y)):

(1) Run forward pass to compute all activations.
(2) Use backprop to perform 1 iteration of
stochastic gradient descent.

(3) Use (eq1) to update \hat{f}_i and
(eq2) to update $b_{\#}^{(i)}$