

MLFA-Neural-Net-2

07-Mar-2023 at 6:22 AM

Deep Learning

Generalized Backpropagation:

Earlier, we have talked about back propagation based learning for MLP with one input, one hidden and one output layer. In doing so, we have used chain rule for computing derivatives of the error function with respect to the parameters.

However, for a deep network, the number of hidden layers is arbitrary. This requires a general approach towards computing the gradient (e.g., performing backpropagation).

A feedforward neural network is a computational graph whose nodes are computing units and whose directed edges transmit numerical information from node to node. Such a network represents a **chain of function compositions** that transform an input to an output pattern. The learning problem consists of finding the optimal combination of weights so that the network function Φ approximates a given function f . Difficulty is that we are not given with the function f explicitly but only given with a set of examples.

For a given training set $X = \{x^t, r^t\}$, let us consider the error function:

$$E = \frac{1}{2} \sum_t (r_t - y_t)^2$$

The weights or parameters are initialized randomly. The backpropagation algorithm is used to find a **local minima** of the error function. The gradient of the error function is calculated and used to update the weights of the previous iteration. To accommodate arbitrary no of hidden layers the gradients have to be computed **recursively**.

Let us consider the weights of the network to be w_1, w_2, \dots, w_l . We need to calculate

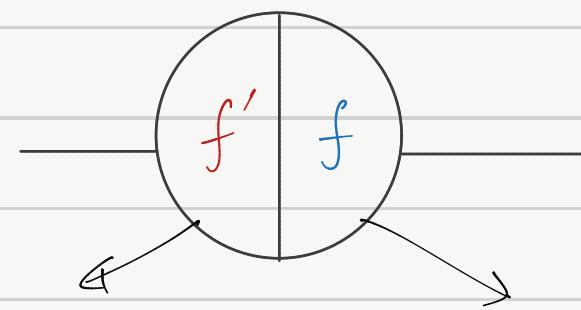
$$\nabla E = \left(\frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_l} \right)$$

The update in weights are

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} \quad i=1, \dots, l$$

Derivatives of Network Functions:

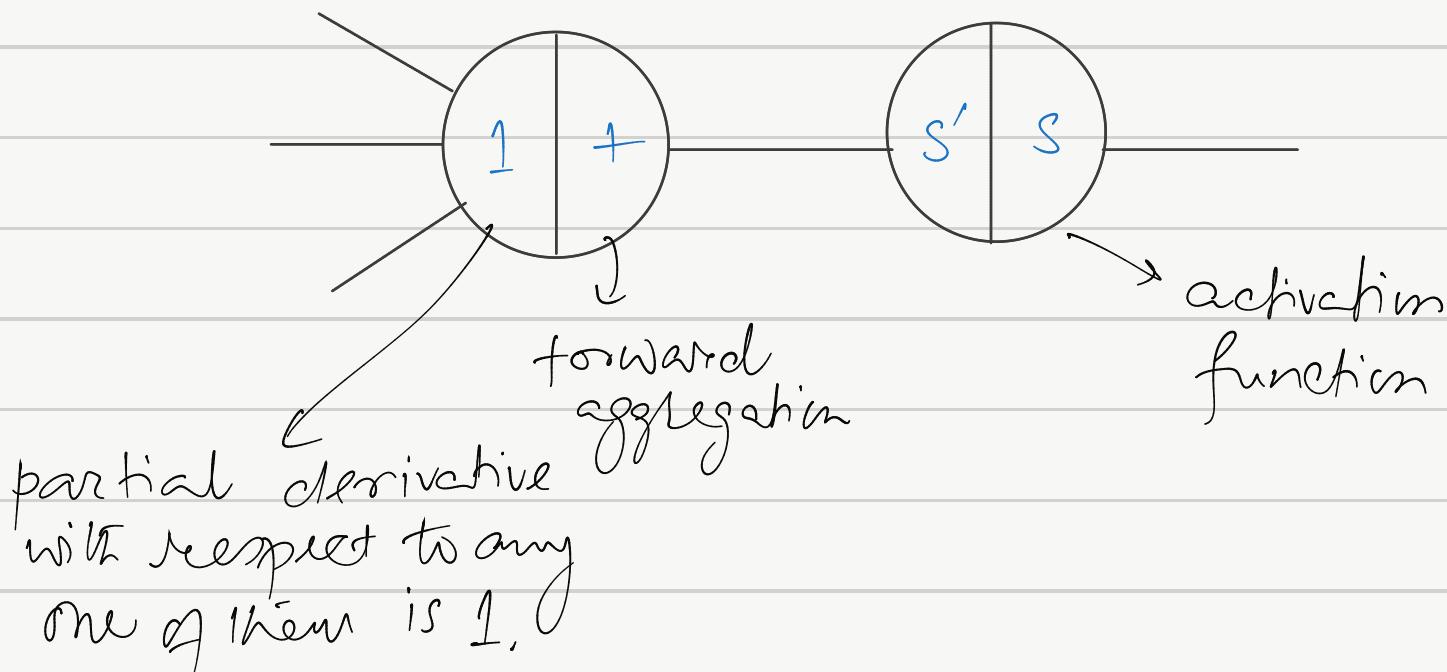
Let us consider a one dimensional network function (having one input and one output). This network is a complex chain of function composition. we represent the nodes as composite structure to record both forward and gradient computation



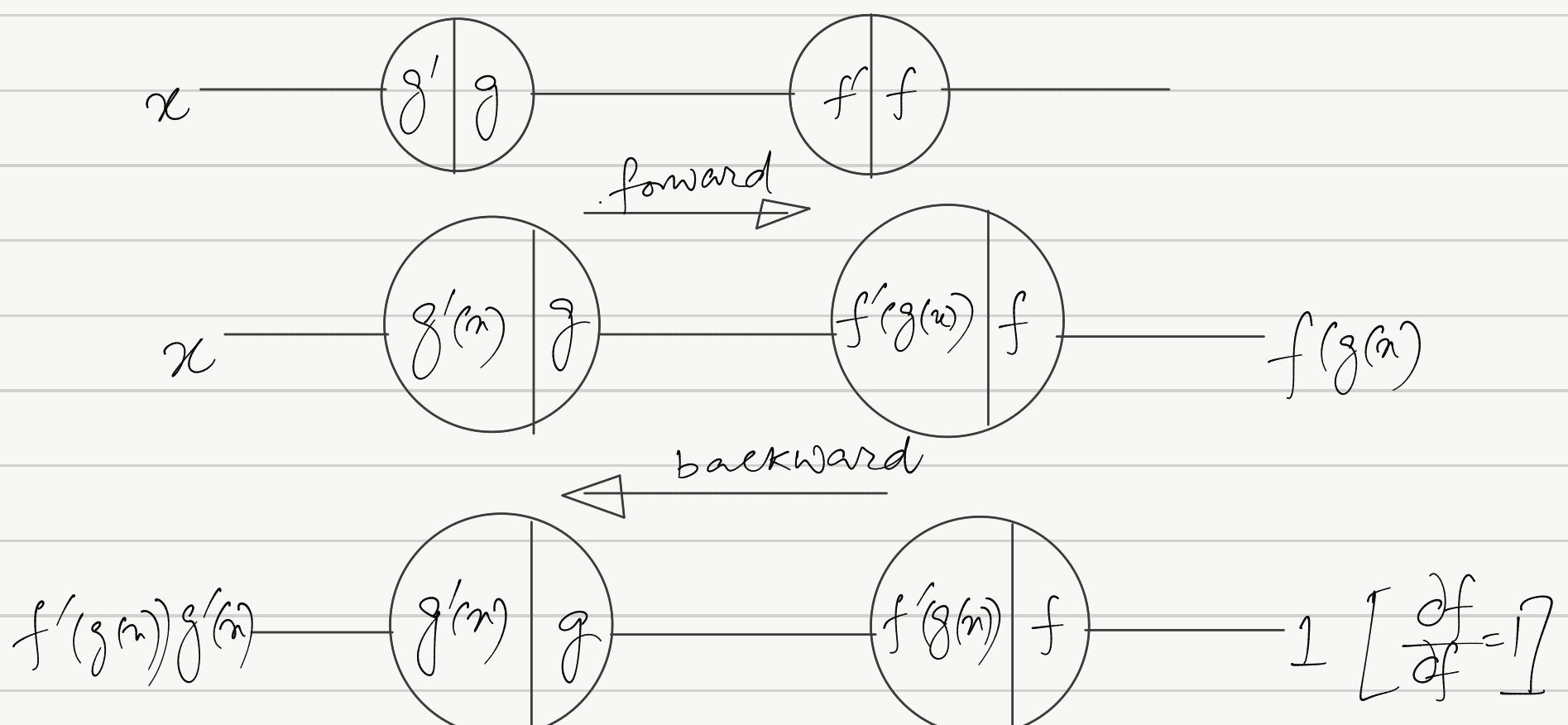
stores in local derivative

stores in forward computation

This can be further detailed as



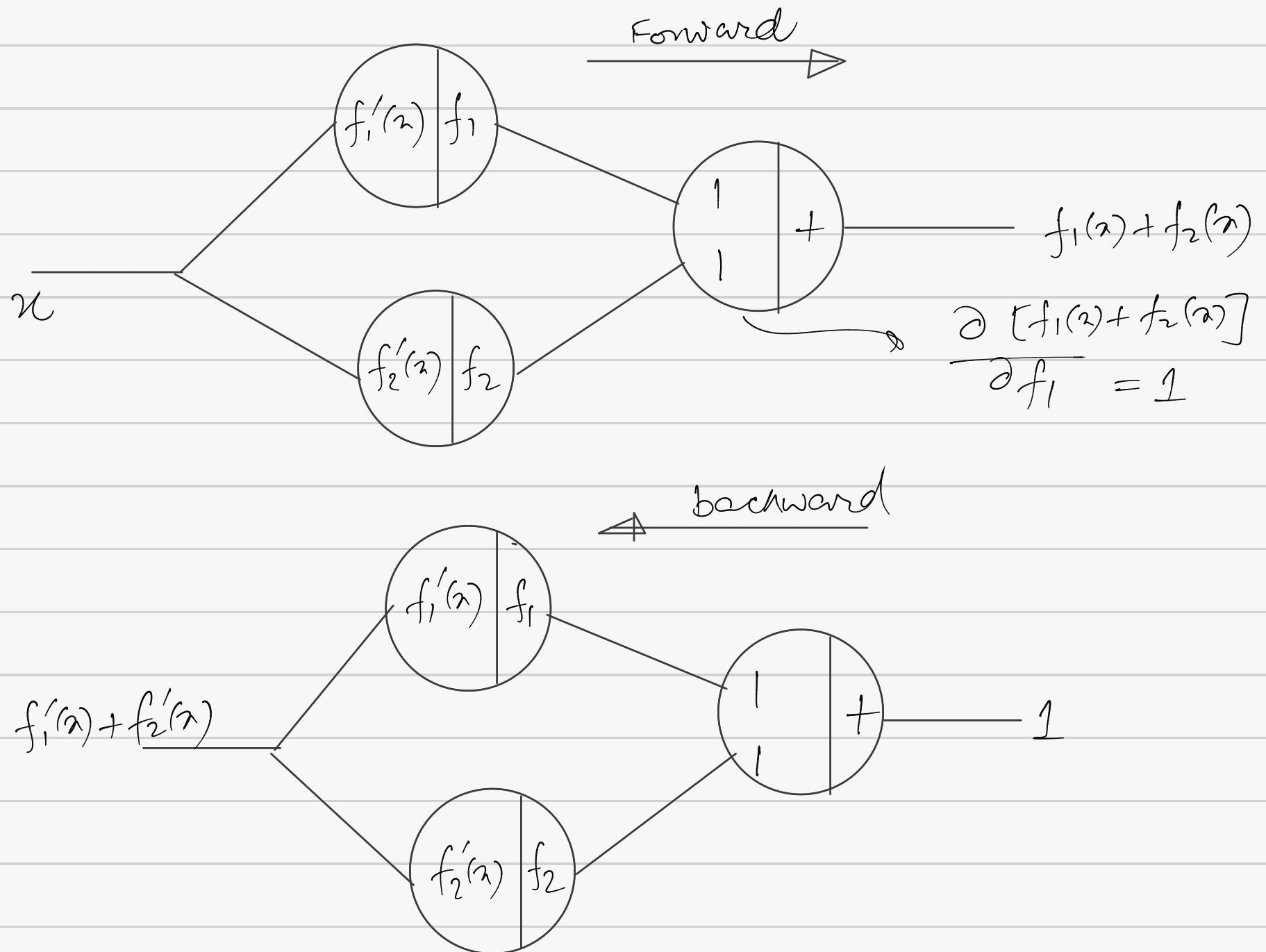
Case 1 : Function Composition $[f(g(x))]$



$$\text{Let } g(x) = y \quad \frac{\partial f(g(x))}{\partial x} = \frac{\partial f}{\partial y} \cdot \frac{\partial y}{\partial x} = f'(g(x)) \cdot g'(x)$$

- \Rightarrow The computation in the network is performed backwards with the input x
 \Rightarrow At each node the product with the value in the left side of the node is computed

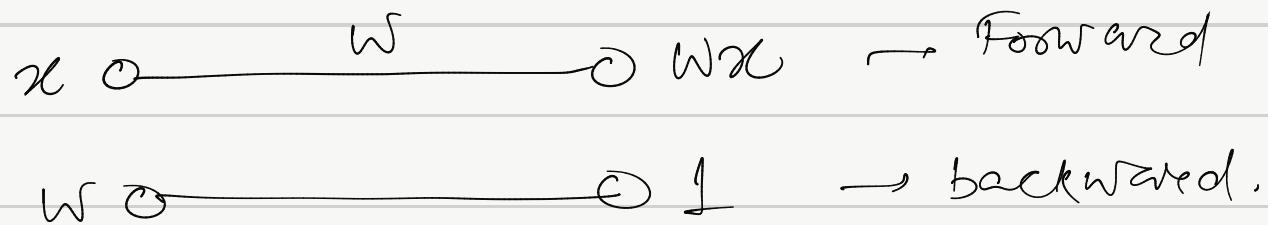
Case 2: Function Addition $[f_1(x) + f_2(x)]$



Let $y_1 = f_1(x)$ and $y_2 = f_2(x)$ $z = f_1(x) + f_2(x)$

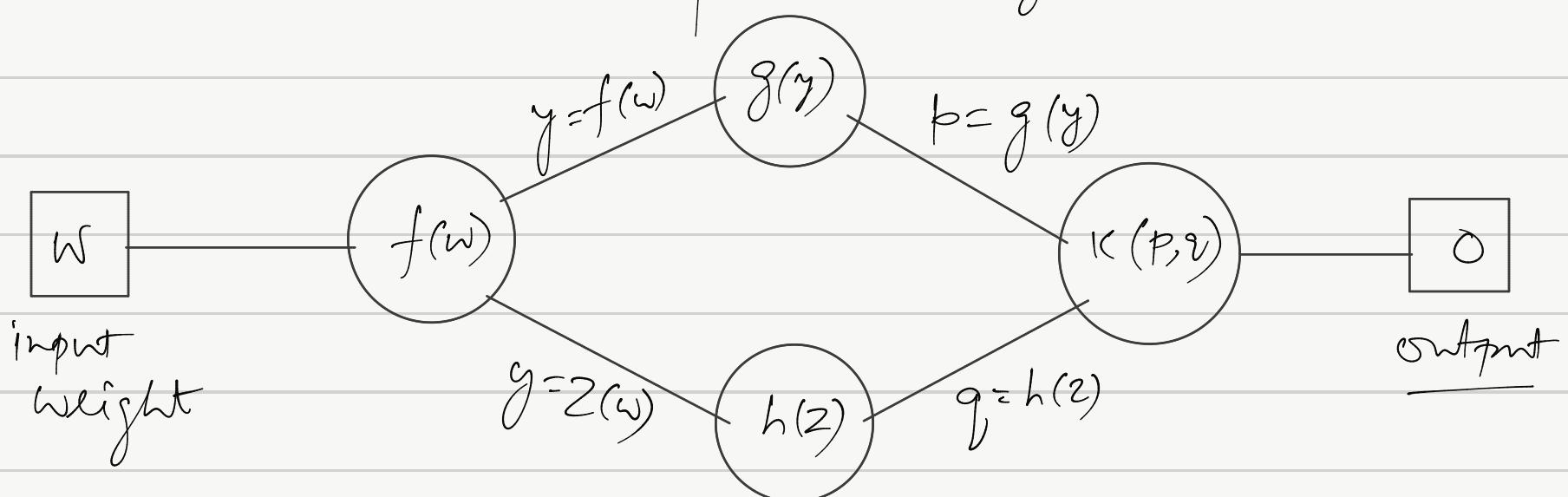
$$\frac{\partial z}{\partial x} = \frac{\partial y_1}{\partial x} + \frac{\partial y_2}{\partial x} = f'_1(x) + f'_2(x)$$

Case 3: Weighted Edges:



Derivative of Function Composition with one weight

- * From now on in derivatives in each node will be dropped from the diagram.



$$o = k(p, q) = k(g(f(w)), h(f(w)))$$

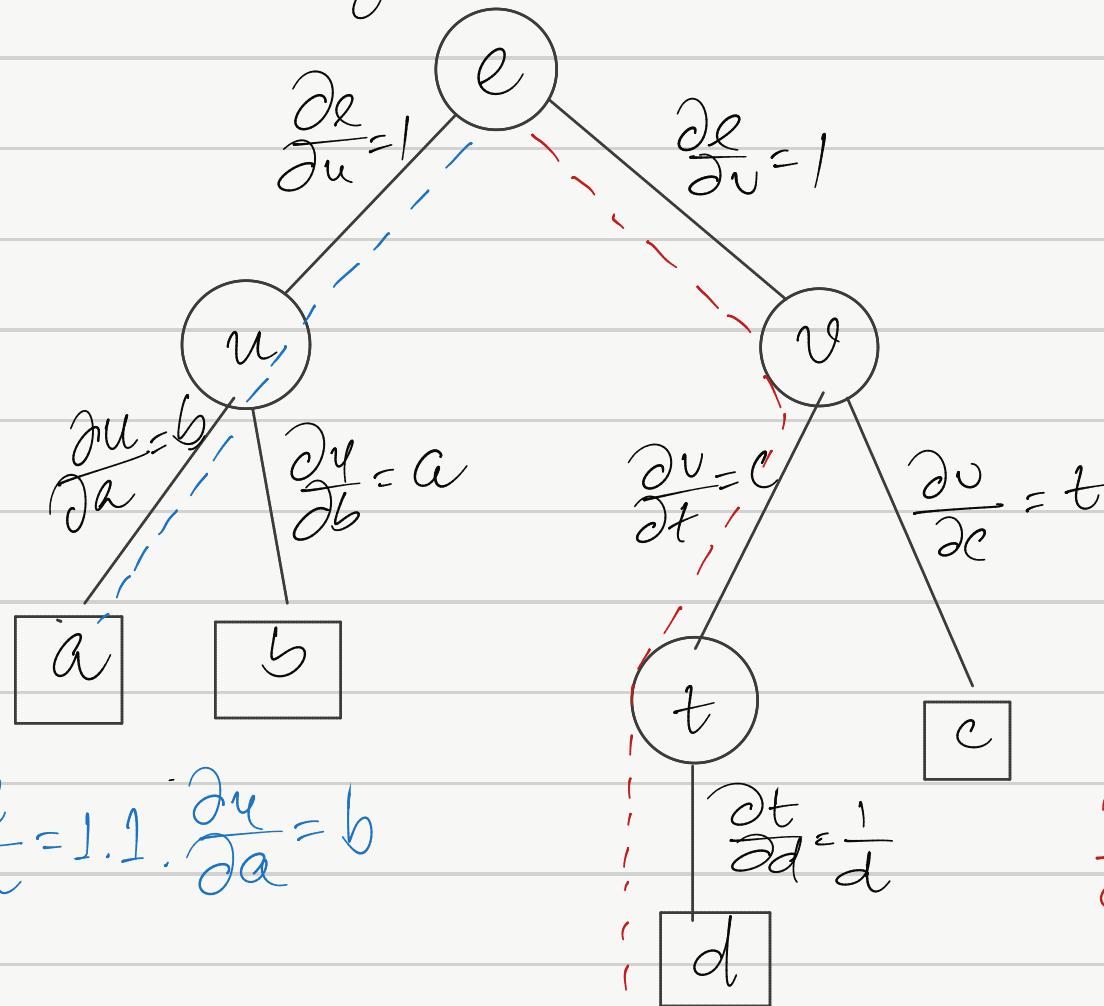
$$\begin{aligned} \frac{\partial o}{\partial w} &= \frac{\partial o}{\partial p} \cdot \frac{\partial p}{\partial y} \cdot \frac{\partial y}{\partial w} + \frac{\partial o}{\partial q} \cdot \frac{\partial q}{\partial z} \cdot \frac{\partial z}{\partial w} \\ &= \frac{\partial k(p, q)}{\partial p} \cdot g'(y) \cdot f'(w) + \frac{\partial k(p, q)}{\partial q} \cdot h'(z) \cdot f'(w) \end{aligned}$$

Derivative of a composite function with 4 inputs

$$e = a \cdot b + c \log d$$

$$\frac{\partial e}{\partial a} = b, \quad \frac{\partial e}{\partial b} = a, \quad \frac{\partial e}{\partial c} = \log d, \quad \frac{\partial e}{\partial d} = \frac{c}{d}$$

Computation graph:



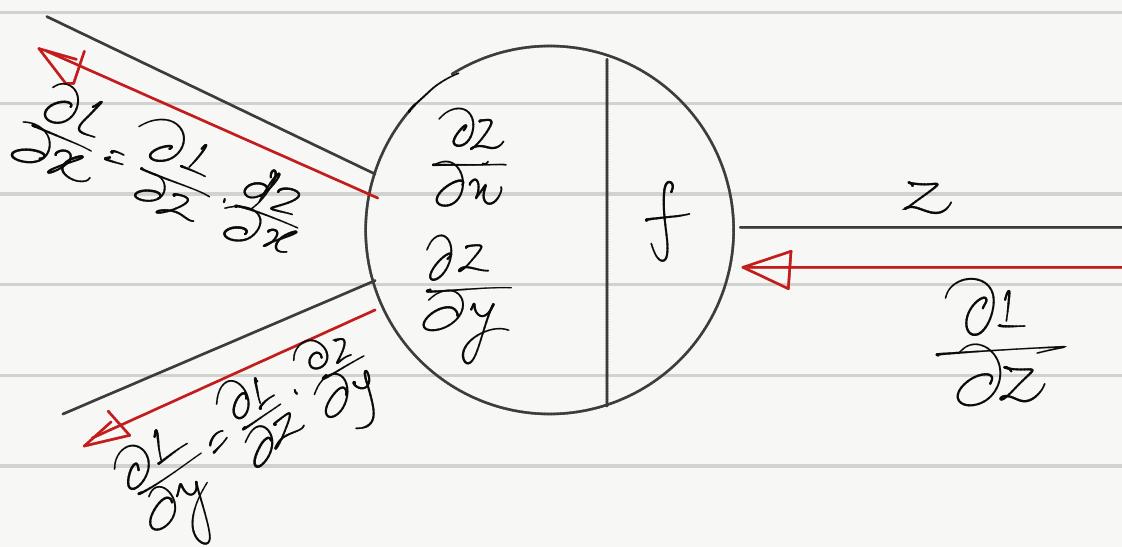
$$l = u + v$$

$$u = a \cdot b$$

$$v = c \cdot t$$

$$t = \log d$$

In general derivative for a neuron $Z = f(x, y)$

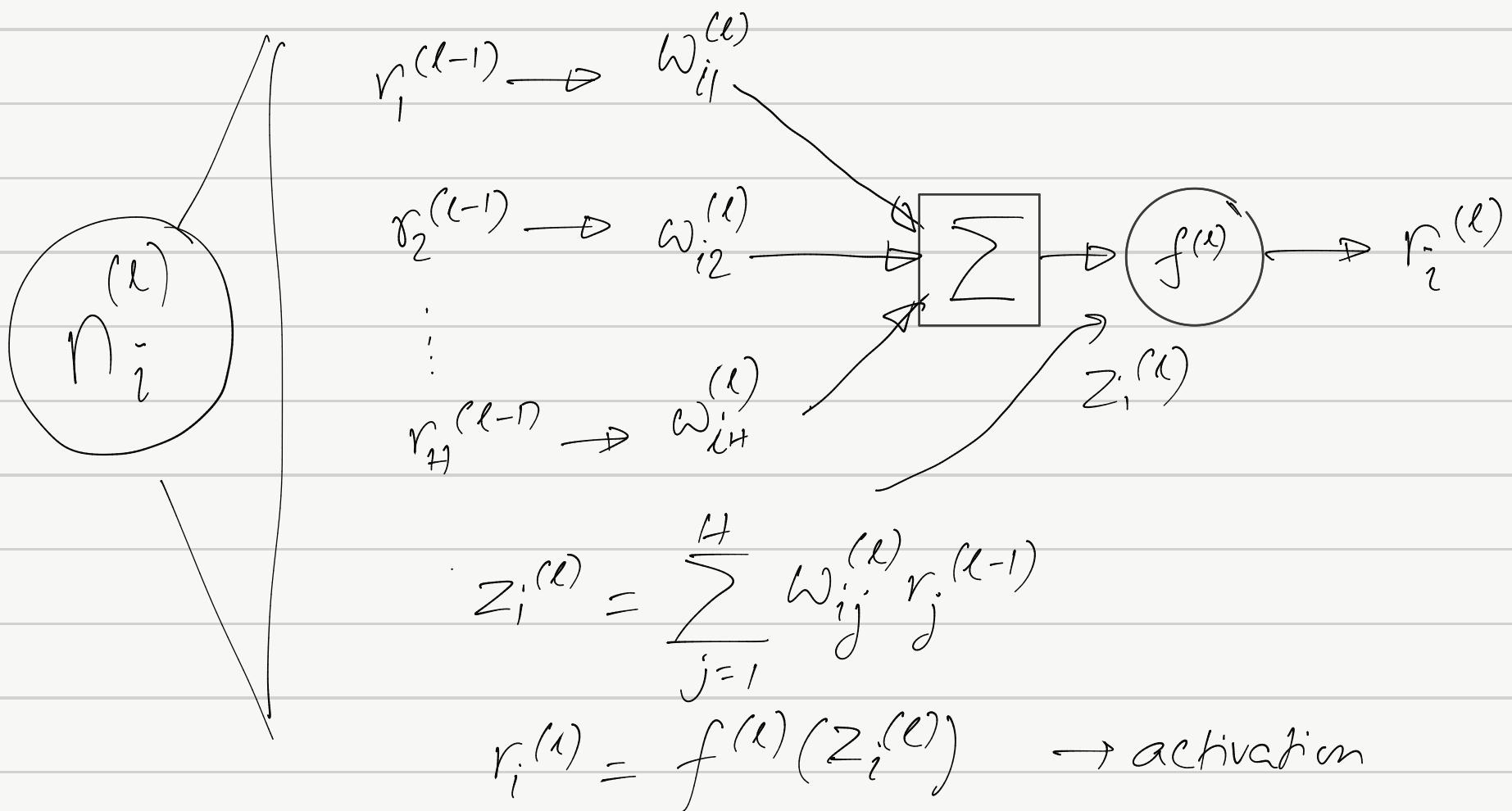


Alg:

- Choose random weights for the network
- Feed an example and do forward Computation
- calculate an error for each node
 - starting from the last stage and propagating error backwards
- update weights
- Repeat with other examples

Derivation of Backpropagation Rule :

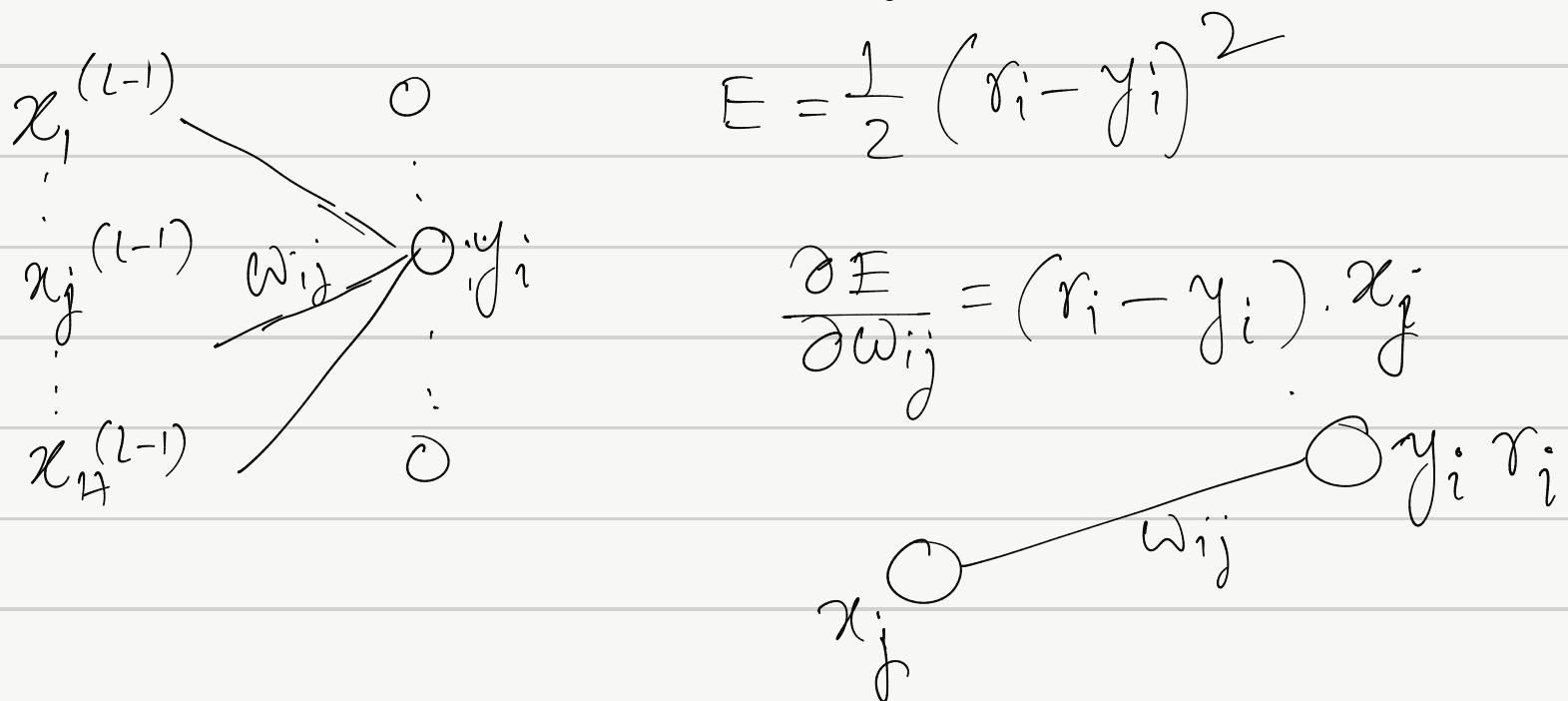
Computation in i^{th} neuron of ℓ^{th} layer :



Update Rule :

$$\Delta w_{ij}^{(\ell)} = -\eta \frac{\partial E}{\partial w_{ij}^{(\ell)}}$$

Consider in simplest case (the output layer)



$$\frac{\partial E}{\partial w_{ij}} = \delta_j x_j$$

In general, for in ℓ^{th} layer

$$\frac{\partial E}{\partial w_{ij}^{(\ell)}} = \frac{\partial E}{\partial z_i^{(\ell)}} \cdot \frac{\partial z_i^{(\ell)}}{\partial w_{ij}^{(\ell)}}$$

$\delta_i^{(\ell)}$ $\rightarrow z_i^{(\ell)} = \sum_{h=1}^H w_{ih}^{(\ell)} x_h^{(\ell-1)}$

local gradient $\frac{\partial z_i^{(\ell)}}{\partial w_{ij}^{(\ell)}} = r_j^{(\ell-1)}$

$$\frac{\partial E}{\partial w_{ij}^{(\ell)}} = \delta_i^{(\ell)} \cdot r_j^{(\ell-1)}$$

How to calculate $\delta_i^{(\ell)}$?

$$\delta_i^{(\ell)} = \frac{\partial E}{\partial z_i^{(\ell)}} = \frac{\partial E}{\partial r_i^{(\ell)}} \cdot \frac{\partial r_i^{(\ell)}}{\partial z_i^{(\ell)}}$$

$$r_i^{(\ell)} = f^{(\ell)}(z_i^{(\ell)}) \quad \frac{\partial r_i^{(\ell)}}{\partial z_i^{(\ell)}} = f'^{(\ell)}$$

$$\delta_i^{(\ell)} = \frac{\partial E}{\partial r_i^{(\ell)}} \cdot f'^{(\ell)}$$

$$\frac{\partial E}{\partial r_i^{(l)}}$$

How does $r_i^{(l)}$ contributes to the error fn. E ?

$\hookrightarrow r_i^{(l)}$ is distributed to all the neurons in the next layer.

\hookrightarrow contribution of $r_i^{(l)}$ to the error function is through all the neurons that it is connected in the next layer.

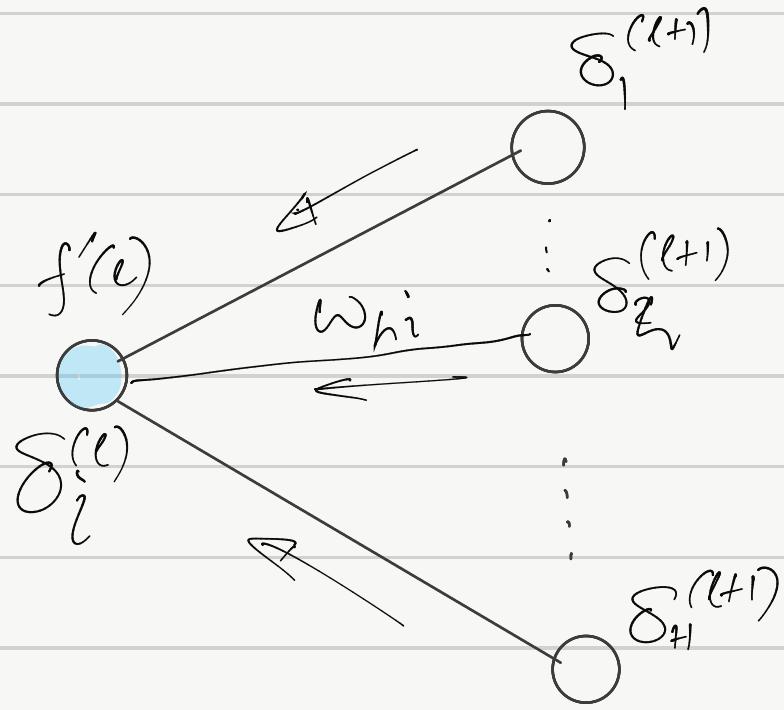
\hookrightarrow Thus variables of the error fn. that are of interest to $r_i^{(l)}$ are $z_1^{(l+1)}, \dots, z_H^{(l+1)}$

$$\frac{\partial E}{\partial r_i^{(l)}} = \frac{\partial E(z_1^{(l+1)}, \dots, z_H^{(l+1)})}{\partial r_i^{(l)}}$$

$$= \sum_{h=1}^H \frac{\partial E}{\partial z_h^{(l+1)}} \cdot \frac{\partial z_h^{(l+1)}}{\partial r_i^{(l)}}$$

$$= \sum_{h=1}^H \delta_h^{(l+1)} \omega_{hi}^{(l+1)}$$

$$\delta_i^{(l)} = \begin{cases} \frac{\partial E}{\partial r_i^{(l)}} \cdot f'(l) & \rightarrow \text{if output neuron} \\ \sum_{h=1}^H \delta_h^{(l+1)} \cdot \omega_{hi}^{(l+1)} \cdot f'(l) & \rightarrow \text{For all other neurons} \end{cases}$$



BP Algo:

⇒ Apply input vector x^t to the network
and forward propagate

$$z_h^{(l)} = \sum_i w_{hi}^{(l)} r_i^{(l-1)}$$

$$r_h^{(l)} = f^{(l)}(z_h^{(l)})$$

⇒ Evaluate δ_n for all the output units

⇒ Backpropagate δ using

$$\delta_i^{(l)} = f'^{(l)} \sum_h \delta_h^{(l+1)} \cdot w_{hi}$$

⇒ Compute $\frac{\partial E}{\partial w_{ij}} = \delta_i^{(l)} \cdot r_j^{(l-1)}$

Improving Training Convergence

Gradient descent in generally has a slow convergence. we will talk about some simple tricks to make it faster.

Momentum:

Slow convergence may happen due to large oscillations in the successive gradient value for a weight.

The update in the current iteration is made dependent on the updates in the previous iterations. (idea of momentum)

Let t be the time index

\Rightarrow epoch no for batch learning

\Rightarrow iteration no in online learning

or minibatch learning

For each weight w_i ,

$$s_i^{(0)} = 0$$

$$s_i^{(t)} = \alpha s_i^{(t-1)} + (1-\alpha) \frac{\partial E^{(t)}}{\partial w_i} \quad [\text{running average}]$$

$$\Delta w_i^{(t)} = -\eta s_i^{(t)}$$

$$0 < \alpha < 1 \rightarrow \text{decay parameter of running average}$$

Adaptive Learning Rate :

Earlier, the η value was kept constant

- * What if we make the learning rate adaptive to the situation
- * In deep network, the contributions of weights to the errors are not the same,
 - ⇒ One can adapt the learning factor separately for each weight, depending on the convergence along that direction
 - ⇒ Accumulate the past error gradient magnitudes for each weight and make learning factor inversely proportional to that
 - ⇒ Learning takes place
 - ⇒ In the direction, where the gradient is small, the learning factor is kept larger to move faster and can be kept smaller along the direction where the gradient is large.

RMS Prop

Keep a running average of in past gradients giving more importance to recent values.

Accumulated past gradients:

$$r_i^{(t)} = \rho r_i^{(t-1)} + (1-\rho) \left| \frac{\partial E^t}{\partial w_i} \right|^2$$

$$\Delta w_i^{(t)} = - \frac{n}{\sqrt{r_i^{(t)}}} \frac{\partial E^t}{\partial w_i}$$

Adam (Adaptive Momentum)

It includes in momentum factor along with the accumulated past error gradients

$$s_i^{(t)} = \alpha s_i^{(t-1)} + (1-\alpha) \frac{\partial E^t}{\partial w_i}$$

$$r_i^{(t)} = \rho r_i^{(t-1)} + (1-\rho) \left| \frac{\partial E^t}{\partial w_i} \right|^2$$

$$\tilde{s}_i^{(t)} = \frac{s_i^{(t)}}{1-\alpha^t} \quad \tilde{r}_i^{(t)} = \frac{r_i^{(t)}}{1-\rho^t} \quad [\alpha^t \rightarrow \alpha \text{ raised to power } t]$$

$$\Delta w_i^t = -n \frac{\tilde{s}_i^{(t)}}{\sqrt{\tilde{r}_i^{(t)}}}$$

Batch Normalization:

The values of the input variables may belong to multiple scales. Thus converting them to a standardized scale is a good idea. In this respect all the input values are z-normalized to have zero mean and unit variance.

Batch normalization extends this idea further by normalizing the hidden unit values before applying the activation function.

Let us consider the weighted sum for a hidden unit j is a_j . For each batch or minibatch for j^{th} unit, we calculate the mean m_j and standard deviation s_j (considering the a_j values)

The z-normalized value :

$$\tilde{a}_j = \frac{a_j - m_j}{s_j}$$

Now it can be mapped to have arbitrary mean γ_j and scale β_j

$$\tilde{a}_j = \gamma_j \tilde{a}_j + \beta_j$$

is then fed to the activation fn.

During testing, we use m_j and s_j is calculated over the whole training data.

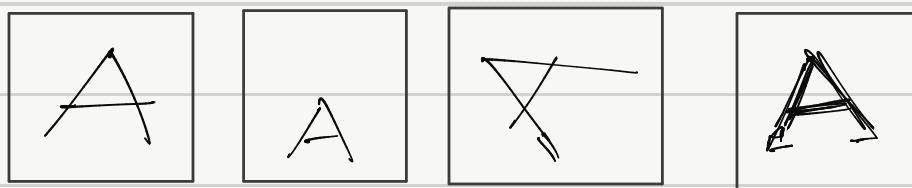
Regularization:

The complexity of an MLP needs to be adjusted to the complexity of the problem reflected by the underlying data.

⇒ Small networks have high bias and large or deep networks have high variances.

We want to deploy large networks, keeping the variance in check at the same time.

Data Augmentation:



⇒ Virtual example

⇒ We can generate multiple copies of the same training examples at different scale and add them to the training set with same label as the original example.
⇒ increase the training set.

⇒ Augmented example may be included in the training example. Let x' is an augmented example for x . Let $g(x|\theta)$ represent one neural network. we can define a penalty

$$E_h = \sum_{x, x'} [g(x|\theta) - g(x'|\theta)]^2$$

$$E' = E + \lambda_h \cdot E_h$$

Weight Decay:

- ⇒ If a weight is zero, the input or hidden unit before that weight does not propagate values.
- ⇒ Traditionally, we initialize the weights to be closer to zero.
 - ⇒ As the learning progresses, more and more weight moves away from zero
 - ⇒ As if complexity of the model increases with training.
 - ⇒ Though the training error may decrease, the validation error may increase

In weight decay, we add some constant force that always pulls a weight towards zero.

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} \rightarrow w_i$$

It is equivalent to perform gradient descent on the error function

$$E' = E + \left[\frac{\lambda}{2} \sum_i w_i^2 \right]$$

L2- Regularization factor for speed of decay

Sum over all the weights in all layers

** If we have two networks that have same training error, the simpler one (with fewer weights) has a higher probability of better generalization.

Dropout:

Adding noise to inputs helps in making the network robust as it does not rely too much on particular input.

⇒ Consider an image data set where, by luck a group of pixels always takes similar values for examples of a class.

⇒ Network will learn to base its decision for most classes on them.

⇒ If we have additional examples where these pixels are missing, the dependency will be reduced

⇒ We are putting off some inputs

⇒ Similar argument may be extended to the hidden layer inputs.

⇒ Don't train to make the decision dependent on a subset of hidden units

⇒ Distribute the decision.

Corrupting or missing out is modelled by dropping out hidden or input units with a probability p .
(set its output to be zero)

For each batch or minibatch, a smaller network is trained

⇒ Dropout is effectively sampling from a pool of networks.

