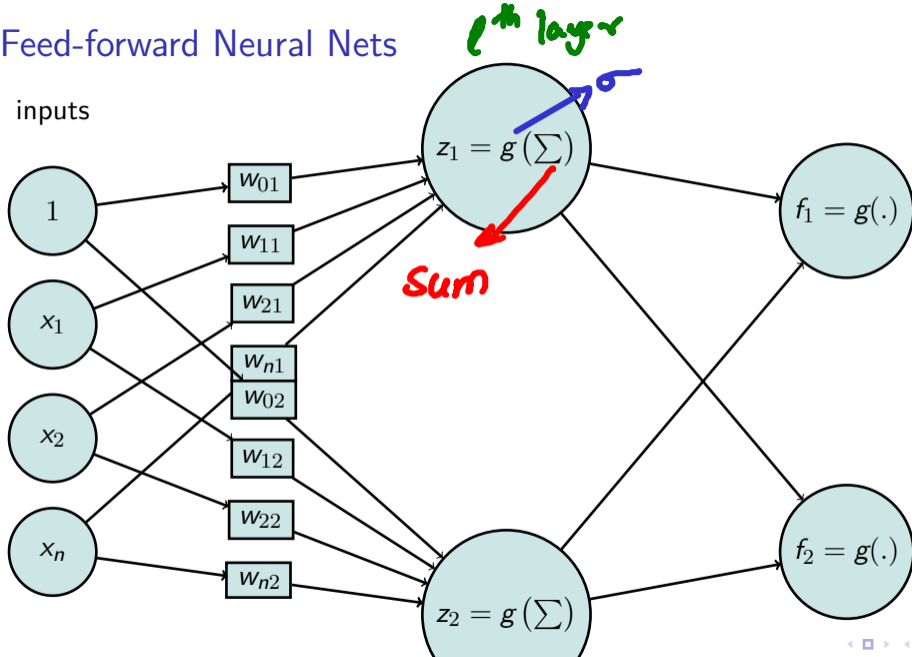


# Lecture 20 contd: Neural Network Training using Backpropagation, Convolutional And Recurrent Neural Networks

Instructor: Prof. Ganesh Ramakrishnan

# Feed-forward Neural Nets

inputs



# Training a Neural Network

STEP 0: Pick a network architecture

- Number of input units: Dimension of features  $\phi(\mathbf{x}^{(i)})$ .
- Number of output units: Number of classes.
- Reasonable default: 1 hidden layer, or if  $>1$  hidden layer, have same number of hidden units in every layer.
- Number of hidden units in each layer a constant factor (3 or 4) of dimension of  $x$ .
- We will use
  - ▶ the smooth sigmoidal function  $g(s) = \frac{1}{1+e^{-s}}$ : **We have now learnt how to train a single node sigmoidal (LR) neural network**
  - ▶ instead of the non-smooth step function  $g(s) = 1$  if  $s \in [\theta, \infty)$  and  $g(s) = 0$  otherwise.

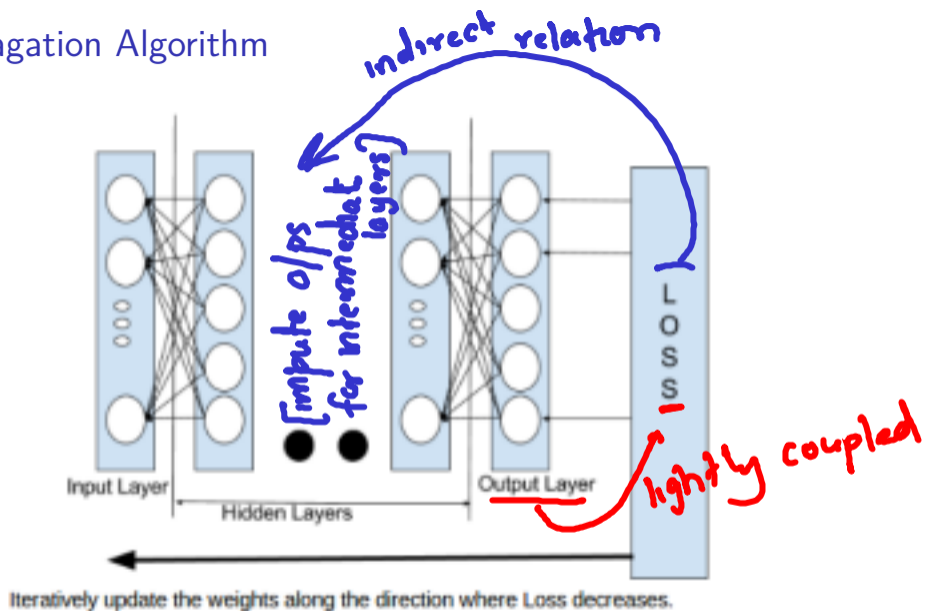
# High Level Overview of Backpropagation Algorithm for Training NN

- 1 Randomly initialize weights  $w_{ij}^l$  for  $l = 1, \dots, L$ ,  $i = 1, \dots, s_l$ ,  $j = 1, \dots, s_{l+1}$ .
- 2 Implement **forward propagation** to get  $f_w(\mathbf{x})$  for any  $x \in \mathcal{D}$ .
- 3 Execute **backpropagation on an misclassified example  $x \in \mathcal{D}$** 
  - 1 by computing partial derivatives  $\frac{\partial}{\partial w_{ij}^{(l)}} E(w)$  for  $l = 1, \dots, L$ ,  $i = 1, \dots, s_l$ ,  $j = 1, \dots, s_{l+1}$ .
  - 2 and using gradient descent to try to minimize (non-convex)  $E(w)$  as a function of parameters  $w$ .

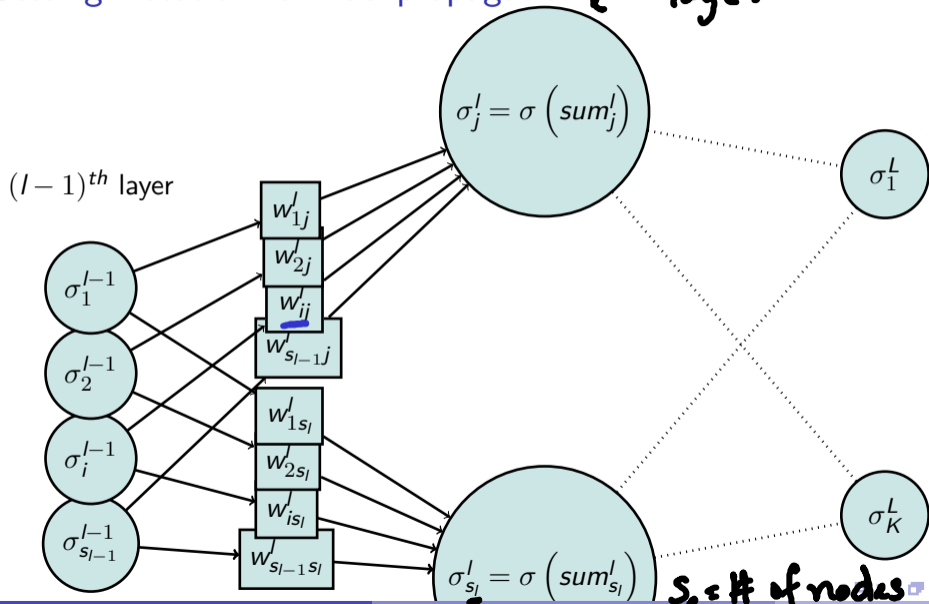
$$w_{ij}^l = w_{ij}^l - \eta \frac{\partial}{\partial w_{ij}^{(l)}} E(w)$$

- 4 Verify that the cost function  $E(w)$  has indeed reduced, else resort to some random perturbation of weights  $w$ .

# The Backpropagation Algorithm



# Setting Notation for Backpropagation $l^{th}$ layer



$s_l = \#$  of nodes in  $l^{th}$  layer

## Gradient Computation

- The Neural Network objective to be minimized:

$$E(\mathbf{w}) = -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log \left( \sigma_k^L(\mathbf{x}^{(i)}) \right) + (1 - y_k^{(i)}) \log \left( 1 - \sigma_k^L(\mathbf{x}^{(i)}) \right) \right] \quad (1)$$

Cross Ent loss  
computed  
only on  $L^{\text{th}}$  layer

$$+ \frac{\lambda}{2m} \sum_{l=1}^L \sum_{i=1}^{s_{l-1}} \sum_{j=1}^{s_l} (w'_{ij})^2$$

Regularization  
is computed  
for  $l=2$  to  $L$

# Gradient Computation

- The Neural Network objective to be minimized:

$$E(\mathbf{w}) = -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log \left( \sigma_k^L(\mathbf{x}^{(i)}) \right) + (1 - y_k^{(i)}) \log \left( 1 - \sigma_k^L(\mathbf{x}^{(i)}) \right) \right] + \frac{\lambda}{2m} \sum_{l=1}^L \sum_{i=1}^{s_{l-1}} \sum_{j=1}^{s_l} (w_{ij}^l)^2 \quad (1)$$

- $\underline{\text{sum}}_j^l = \sum_{k=1}^{s_{l-1}} \underline{w}_{kj}^l \sigma_k^{l-1}$  and  $\sigma_i^l = \frac{1}{1 + e^{-\text{sum}_i^l}}$
- $\frac{\partial E}{\partial w_{ij}^l} = \frac{\partial E}{\partial \sigma_j^l} \frac{\partial \sigma_j^l}{\partial \text{sum}_j^l} \frac{\partial \text{sum}_j^l}{\partial w_{ij}^l} + \frac{\lambda}{2m} w_{ij}^l \longrightarrow$
- $\frac{\partial \sigma_j^l}{\partial \text{sum}_j^l} = \left( \frac{1}{1 + e^{-\text{sum}_j^l}} \right) \left( 1 - \frac{1}{1 + e^{-\text{sum}_j^l}} \right) = \sigma_j^l (1 - \sigma_j^l)$
- $\frac{\partial \text{sum}_j^l}{\partial w_{ij}^l} = \frac{\partial}{\partial w_{ij}^l} \left( \sum_{k=1}^{s_{l-1}} w_{kj}^l \sigma_k^{l-1} \right) = \sigma_i^{l-1}$

$E$  depends on  $\sigma_j^l$  thru  $(\sigma^{l+1})$

How to compute

$$\frac{\partial E}{\partial \sigma_j^l}$$

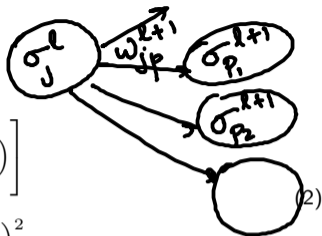
Ans: Back propagation!



- For a single example  $(\mathbf{x}, y)$ :

*Recall simple logistic case*  
 $\frac{\partial E}{\partial \sigma_j^L} = -\frac{y_j}{\sigma_j^L} - \frac{(1-y_j)}{(1-\sigma_j^L)}$

$$- \left[ \sum_{k=1}^K y_k \log(\sigma_k^L(\mathbf{x})) + (1 - y_k) \log(1 - \sigma_k^L(\mathbf{x})) \right] + \frac{\lambda}{2m} \sum_{l=1}^L \sum_{i=1}^{s_{l-1}} \sum_{j=1}^{s_l} (w_{ij}^l)^2$$



- $\frac{\partial E}{\partial \sigma_j^l} = \sum_{p=1}^{s_{l+1}} \frac{\partial E}{\partial \text{sum}_p^{l+1}} \frac{\partial \text{sum}_p^{l+1}}{\partial \sigma_j^l} = \sum_{p=1}^{s_{l+1}} \frac{\partial E}{\partial \sigma_p^{l+1}} \frac{\partial \sigma_p^{l+1}}{\partial \text{sum}_p^{l+1}} w_{jp}^{l+1}$  since  $\frac{\partial \text{sum}_p^{l+1}}{\partial \sigma_j^l} = w_{jp}^{l+1}$

- $\frac{\partial E}{\partial \sigma_j^L} = -\frac{y_j}{\sigma_j^L} + \frac{1-y_j}{1-\sigma_j^L}$

$\rightarrow \sigma_p^{l+1} (1 - \sigma_p^{l+1})$

$E = f_n(\sigma_{\{k\}}^L (\sigma_{\{k'\}}^{L-1} (\sigma_{\{k''\}}^{L-2} \dots (\sigma_{\{p\}}^{l+1} (\sigma_j^l \dots))))))$

$\sigma_p^{l+1} = g(\text{sum}_p^{l+1})$

$\text{sum}_p^{l+1} = \sum_j \sigma_j^l w_{jp}^{l+1}$

# Backpropagation in Action

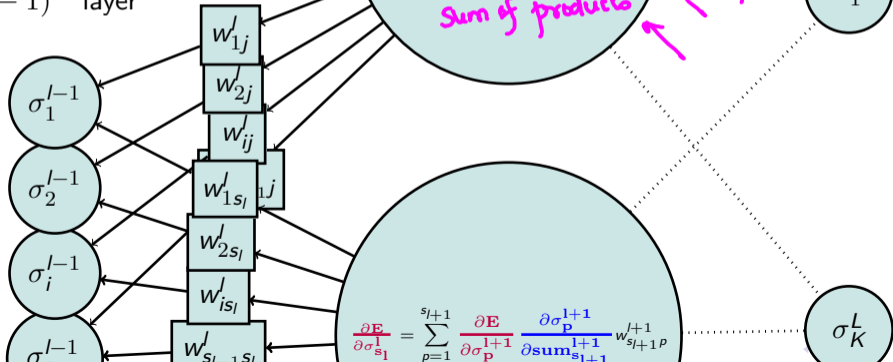
$p$ th layer

Note:  $\frac{\partial E}{\partial \sigma_p^{l+1}}$  is independent of  $j$  and can be memorized/cached

$$\frac{\partial E}{\partial \sigma_j^{l+1}} = \sum_{p=1}^{s_{l+1}} \frac{\partial E}{\partial \sigma_p^{l+1}} \frac{\partial \sigma_p^{l+1}}{\partial \text{sum}_p^{l+1}} w_{jp}^{l+1}$$

Sum of products

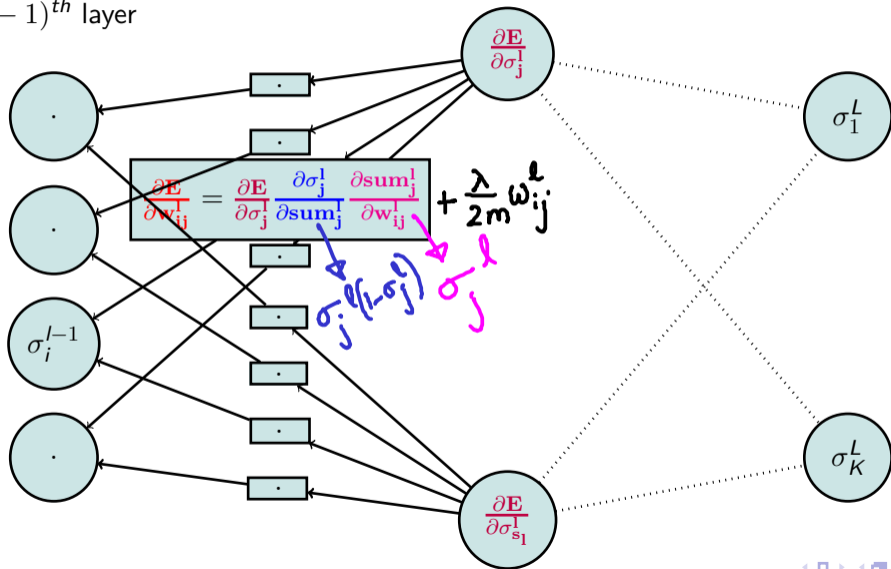
$(l-1)$ th layer



$$\frac{\partial E}{\partial \sigma_{s_1}^{l+1}} = \sum_{p=1}^{s_{l+1}} \frac{\partial E}{\partial \sigma_p^{l+1}} \frac{\partial \sigma_p^{l+1}}{\partial \text{sum}_{s_1+1}^{l+1}} w_{s_1+1 p}^{l+1}$$

# Backpropagation in Action

$(l-1)^{th}$  layer



## Recall and Substitute

- $sum_j^l = \sum_{k=1}^{s_{l-1}} w_{kj}^l \sigma_k^{l-1}$  and  $\sigma_i^l = \frac{1}{1+e^{-sum_i^l}}$

- $\frac{\partial E}{\partial w_{ij}^l} = \frac{\partial E}{\partial \sigma_j^l} \frac{\partial \sigma_j^l}{\partial sum_j^l} \frac{\partial sum_j^l}{\partial w_{ij}^l} + \frac{\lambda}{2m} w_{ij}^l$

- $\frac{\partial \sigma_j^l}{\partial sum_j^l} = \sigma_j^l (1 - \sigma_j^l)$

- $\frac{\partial sum_j^l}{\partial w_{ij}^l} = \sigma_i^{l-1}$

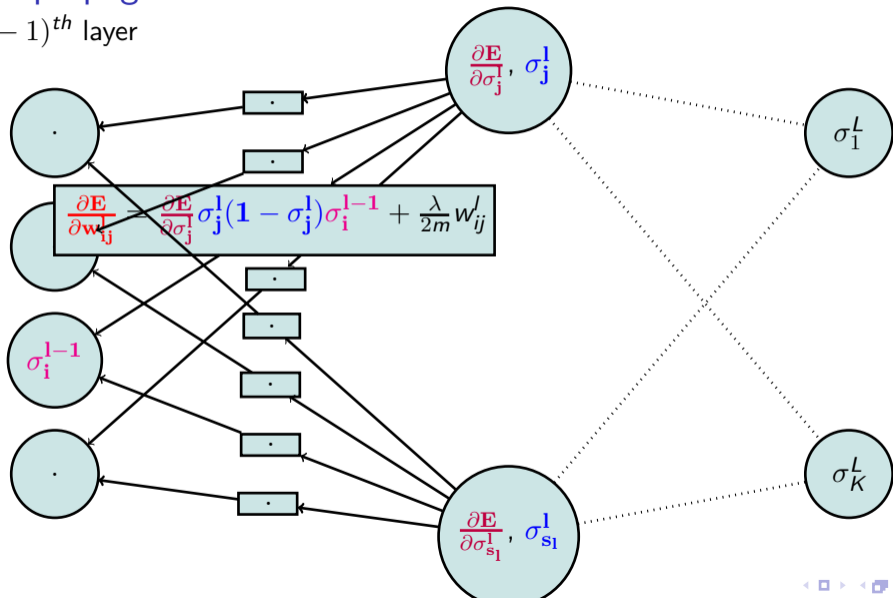
- $\frac{\partial E}{\partial \sigma_j^l} = \sum_{p=1}^{s_{l+1}} \frac{\partial E}{\partial \sigma_j^{l+1}} \frac{\partial \sigma_j^{l+1}}{\partial sum_p^{l+1}} w_{jp}^{l+1}$

- $\frac{\partial E}{\partial \sigma_j^L} = -\frac{y_j}{\sigma_j^L} - \frac{1-y_j}{1-\sigma_j^L}$

} Already done

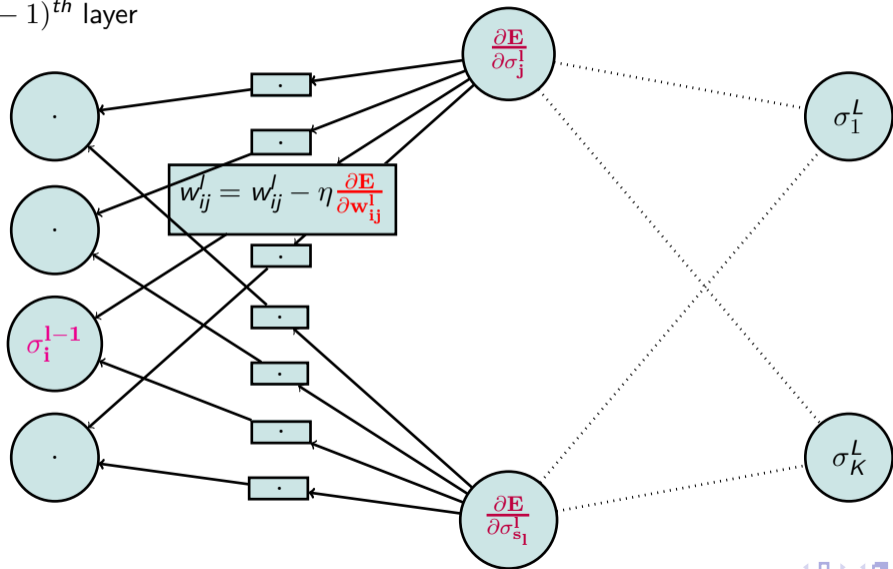
# Backpropagation in Action

$(l-1)^{th}$  layer



# Backpropagation in Action

$(l-1)^{th}$  layer



# The Backpropagation Algorithm for Training NN

- 1 Randomly initialize weights  $w_{ij}^l$  for  $l = 1, \dots, L$ ,  $i = 1, \dots, s_l$ ,  $j = 1, \dots, s_{l+1}$ .
- 2 Implement **forward propagation** to get  $f_w(\mathbf{x})$  for every  $\mathbf{x} \in \mathcal{D}$ .
- 3 Execute **backpropagation** on any misclassified  $\mathbf{x} \in \mathcal{D}$  by performing gradient descent to minimize (non-convex)  $E(\mathbf{w})$  as a function of parameters  $\mathbf{w}$ .

4  $\frac{\partial E}{\partial \sigma_j^L} = -\frac{y_j}{\sigma_j^L} - \frac{1-y_j}{1-\sigma_j^L}$  for  $j = 1$  to  $s_L$ . *→ For  $L^{\text{th}}$  layer*

- 5 For  $l = L - 1$  down to 2:

1  $\frac{\partial E}{\partial \sigma_j^l} = \sum_{p=1}^{s_{l+1}} \frac{\partial E}{\partial \sigma_j^{l+1}} \sigma_j^{l+1} (1 - \sigma_j^{l+1}) w_{jp}^{l+1}$

2  $\frac{\partial E}{\partial w_{ij}^l} = \frac{\partial E}{\partial \sigma_j^l} \sigma_j^l (1 - \sigma_j^l) \sigma_i^{l-1} + \frac{\lambda}{2m} w_{ij}^l$

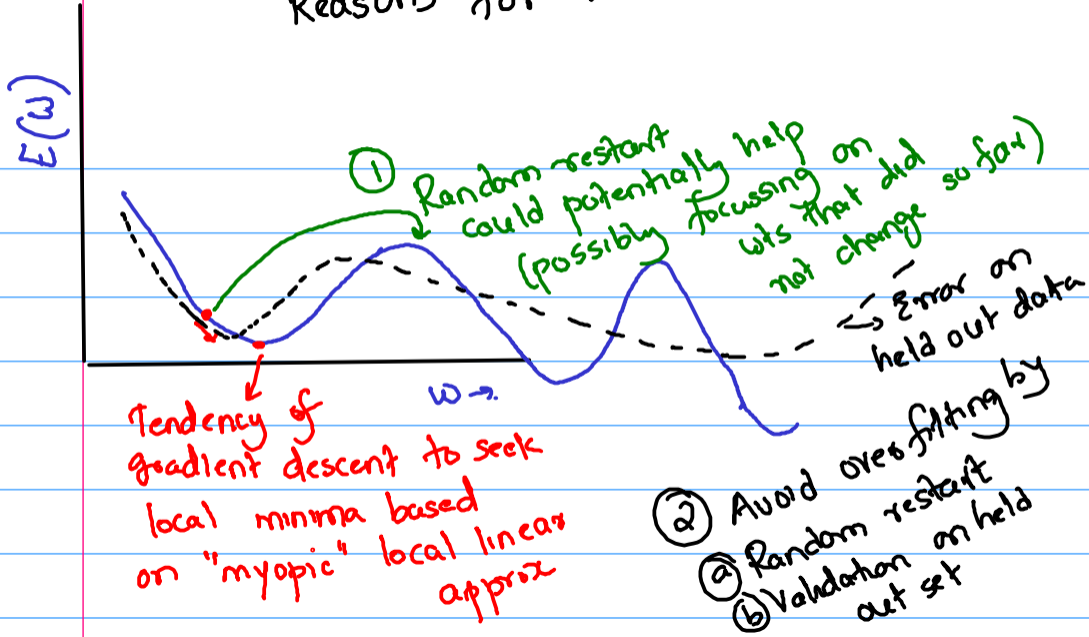
3  $w_{ij}^l = w_{ij}^l - \eta \frac{\partial E}{\partial w_{ij}^l}$

*Tut 8 problem: Compute min numbers of multiplications and additions for a single backprop assuming infinite memory*

- 6 Keep picking misclassified examples until the cost function  $E(\mathbf{w})$  shows significant reduction; else resort to some random perturbation of weights  $\mathbf{w}$  and restart a couple of times.

*Additionally estimate max. memory reqd!*

# Reasons for restart





## Challenges with Neural Network Training

- ① Local optima... only approx correct solution
- ② Slower, more memory than other methods
- ③ Architecture design & associating semantics
- ④ Numerical precision errors with layers

⑤ Representability vs learnability tradeoff

More the # of layers & nodes, more is ability to model

complex fns

# of examples required to learn increases exponentially (Curse of dimensionality)

[Articulate this claim in Tut 8 for all possible activation units]

## What Changed with Neural Networks?

Reason 1: Better computation infrastructure (GPUS, dist processing, MPI) & lots of data

- Origin: Computational Model of Threshold Logic from Warren McCulloch and Walter Pitts (1943)
- Big Leap: For ImageNet Challenge, AlexNet achieved 85 % accuracy (NIPS 2012). Prev best was 75 % (CVPR 2011).
- Present best is 96.5 % MSRA (arXiv 2015). Comparable to human level accuracy.
- Challenges involved were varied background, same object with different colors (e.g. cats), varied sizes and postures of same objects, varied illuminated conditions.
- Tasks like OCR, Speech recognition are now possible without segmenting the word image/signal into character images/signals.

Pretraining = Unsupervised learning of w's in first (few) layer(s)

Reason 2: ML based principles for sparser/specially constructed n/ws, param sharing, pretraining, convolution

## LeNet(1989 and 1998) v/s AlexNet(2012)

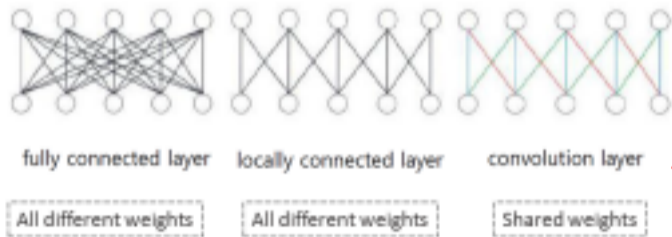
	LeeNet 1989	LeeNet 1998	AlexNet 2012
Task	Digit Recognition	Digit Recognition	Object Recognition
# Classes	10	10	1k
Image Size	16 X 16	28 X 28	256 X 256 X 3
# Examples	7k	60k	1.2 M
units	1256	8084	658 k
parameters	9760	60k	60 M
connections	65k	344k	652M
Summation	Sigmoid	Sigmoid	ReLU
GPU/ Non-GPU	Non-GPU based.	Non-GPU based.	GPU based.
Operations	11 billion	412 billions	200 quadrillions

## Reasons for Big Leap

- Why LeNet was not as successful as AlexNet, though the algorithm was same?
- Right algorithm at wrong time.
- Modern features.
- Advancement in Machine learning.
- Realistic data collection in huge amount due to: regular competitions, evaluation metrics or challenging problem statements.
- Advances in Computational Resources: GPUs, industrial scale clusters.
- Evolution of tasks: Classification of 10 objects to 100 objects to “structure of classes”.

# Convolutional Neural Network

- Variation of multi layer feedforward neural network designed to use minimal preprocessing with wide application in image recognition and natural language processing
- Traditional multilayer perceptron(MLP) models do not take into account spatial structure of data and suffer from curse of dimensionality
- Convolution Neural network has smaller number of parameters due to **local connections** and **weight sharing**



*H/W: How to compute shared wts?*