Open Elective Course [OE]
**Course Code: CSO507**
**Winter 2023-24**

Lecture#

# Deep Learning

## Unit-3: Artificial Neural Network (Part-III)
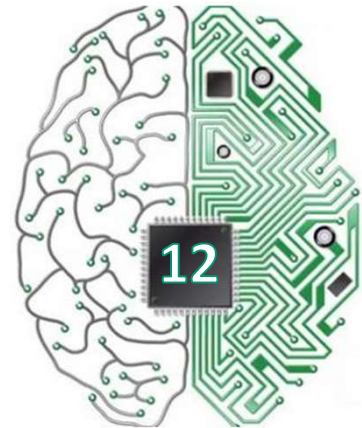### Unit-4: Convolutional Neural Network (Part-I: Motivation)

12

Course Instructor:

**Dr. Monidipa Das**

**Assistant Professor**

**Department of Computer Science and Engineering**

**Indian Institute of Technology (Indian School of Mines) Dhanbad, Jharkhand 826004, India**
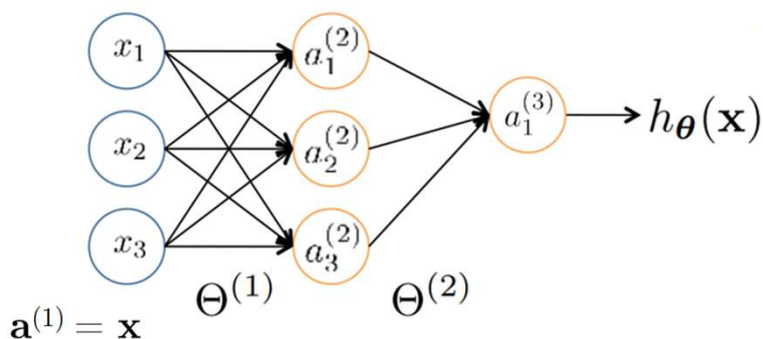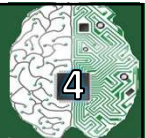
---

2

# Neural Network Learning

# Learning in NN: Backpropagation

- Similar to the perceptron learning algorithm, we cycle through our examples
  - If the output of the network is correct, no changes are made
  - If there is an error, weights are adjusted to reduce the error

- The trick is to assess the blame for the error and divide it among the contributing weights

# Forward Propagation

Feed-Forward Steps:

$$\mathbf{z}^{(2)} = \Theta^{(1)}\mathbf{x} = \Theta^{(1)}\mathbf{a}^{(1)}$$

$$\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$$

Add $a_0^{(2)} = 1$

$$\mathbf{z}^{(3)} = \Theta^{(2)}\mathbf{a}^{(2)}$$

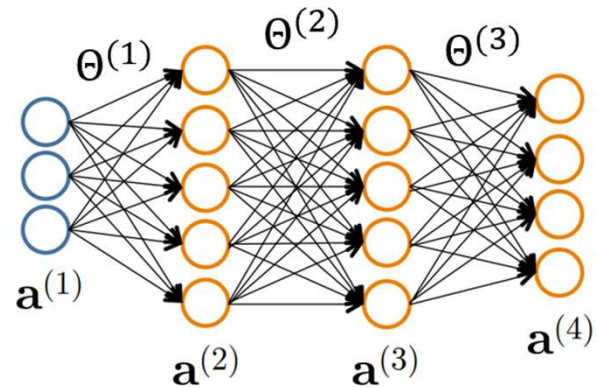$$h_\Theta(\mathbf{x}) = \mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$$

# Forward Propagation

- Given one labeled training instance $(\mathbf{x},\, y)$:

**Forward Propagation**

- $\mathbf{a}^{(1)} = \mathbf{x}$
- $\mathbf{z}^{(2)} = \Theta^{(1)}\mathbf{a}^{(1)}$
- $\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$     [add $\mathrm{a}_0{}^{(2)}$]
- $\mathbf{z}^{(3)} = \Theta^{(2)}\mathbf{a}^{(2)}$
- $\mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$     [add $\mathrm{a}_0{}^{(3)}$]
- $\mathbf{z}^{(4)} = \Theta^{(3)}\mathbf{a}^{(3)}$
- $\mathbf{a}^{(4)} = \mathrm{h}_\Theta(\mathbf{x}) = g(\mathbf{z}^{(4)})$



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Cost Function

**Logistic Regression:**

$$J(\theta) = -\frac{1}{n}\sum_{i=1}^{n}[y_i \log h_{\boldsymbol{\theta}}(\mathbf{x}_i) + (1-y_i)\log(1-h_{\boldsymbol{\theta}}(\mathbf{x}_i))] + \frac{\lambda}{2n}\sum_{j=1}^{d}\theta_j^2$$

**Neural Network:**

$$h_\Theta \in \mathbb{R}^K \qquad (h_\Theta(\mathbf{x}))_i = i^{th}\text{output}$$

$$J(\Theta) = -\frac{1}{n}\left[\sum_{i=1}^{n}\sum_{k=1}^{K} y_{ik}\log(h_\Theta(\mathbf{x}_i))_k + (1-y_{ik})\log\left(1-(h_\Theta(\mathbf{x}_i))_k\right)\right]$$

$$+ \frac{\lambda}{2n}\sum_{l=1}^{L-1}\sum_{i=1}^{s_{l-1}}\sum_{j=1}^{s_l}\left(\Theta_{ji}^{(l)}\right)^2$$

$k^{th}$ class:     true, predicted
not $k^{th}$ class: true, predicted

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Optimizing the Neural Network

$$J(\Theta) = -\frac{1}{n} \left[ \sum_{i=1}^{n} \sum_{k=1}^{K} y_{ik} \log(h_\Theta(\mathbf{x}_i))_k + (1 - y_{ik}) \log\left(1 - (h_\Theta(\mathbf{x}_i))_k\right) \right]$$

$$+ \frac{\lambda}{2n} \sum_{l=1}^{L-1} \sum_{i=1}^{s_{l-1}} \sum_{j=1}^{s_l} \left(\Theta_{ji}^{(l)}\right)^2$$

Solve via: $\min_{\Theta} J(\Theta)$

> $J(\Theta)$ is not convex, so GD on a neural net yields a local optimum
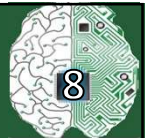> - But, tends to work well in practice

Need code to compute:
- $J(\Theta)$
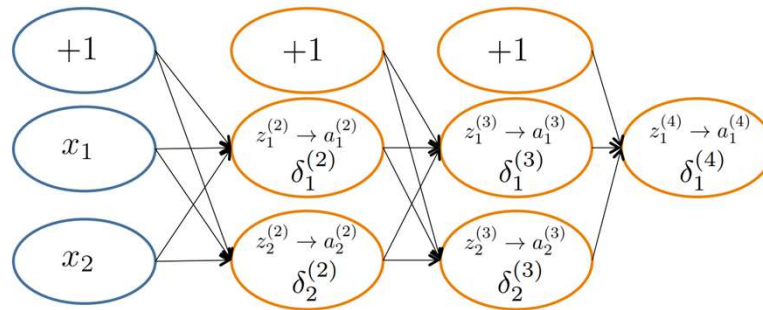- $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$

# Backpropagation Intuition

- Each hidden node $j$ is "responsible" for some fraction of the error $\delta_j^{(l)}$ in each of the output nodes to which it connects

- $\delta_j^{(l)}$ is divided according to the strength of the connection between hidden node and the output node

- Then, the "blame" is propagated back to provide the error values for the hidden layer
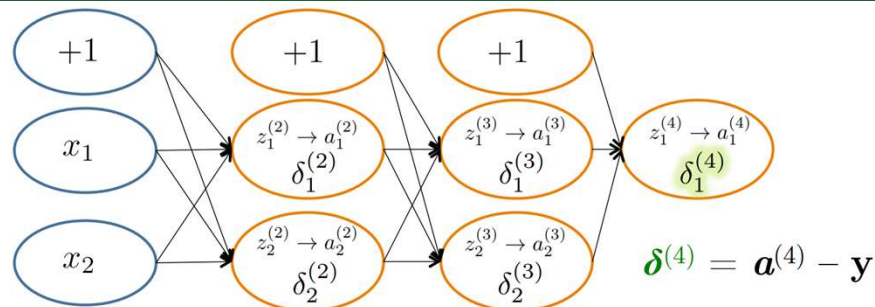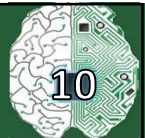
# Backpropagation Intuition



$\delta_j^{(l)}$ = "error" of node $j$ in layer $l$

Formally, $\delta_j^{(l)} = \dfrac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

where $\text{cost}(\mathbf{x}_i) = [y_i \log h_\Theta(\mathbf{x}_i) + (1 - y_i) \log(1 - h_\Theta(\mathbf{x}_i))]$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Backpropagation Intuition



$$\boldsymbol{\delta}^{(4)} = \boldsymbol{a}^{(4)} - \mathbf{y}$$

$\delta_j^{(l)}$ = "error" of node $j$ in layer $l$

Formally, $\delta_j^{(l)} = \dfrac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

where $\text{cost}(\mathbf{x}_i) = [y_i \log h_\Theta(\mathbf{x}_i) + (1 - y_i) \log(1 - h_\Theta(\mathbf{x}_i))]$

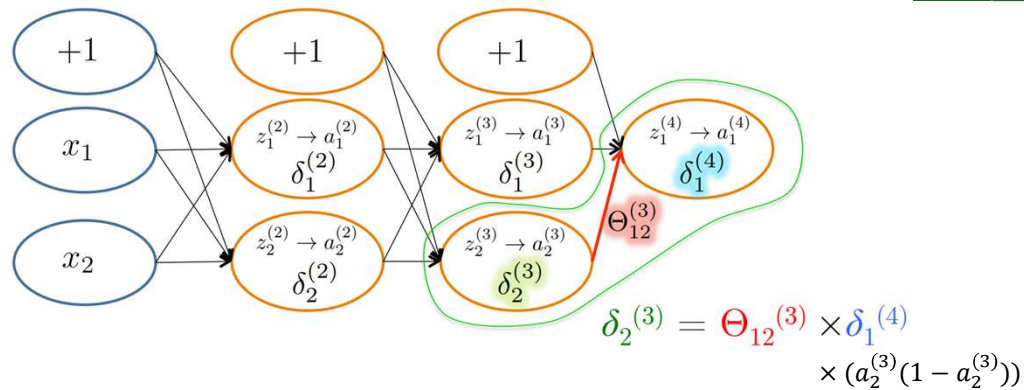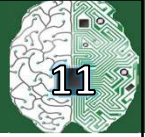Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Backpropagation Intuition

**11**

$$+1 \qquad +1 \qquad +1$$

$$x_1 \qquad \begin{array}{c} z_1^{(2)} \to a_1^{(2)} \\ \delta_1^{(2)} \end{array} \qquad \begin{array}{c} z_1^{(3)} \to a_1^{(3)} \\ \delta_1^{(3)} \end{array} \qquad \begin{array}{c} z_1^{(4)} \to a_1^{(4)} \\ \delta_1^{(4)} \end{array}$$

$$x_2 \qquad \begin{array}{c} z_2^{(2)} \to a_2^{(2)} \\ \delta_2^{(2)} \end{array} \qquad \begin{array}{c} z_2^{(3)} \to a_2^{(3)} \\ \delta_2^{(3)} \end{array} \qquad \Theta_{12}^{(3)}$$

$$\delta_2^{(3)} = \Theta_{12}^{(3)} \times \delta_1^{(4)}$$
$$\times (a_2^{(3)}(1 - a_2^{(3)}))$$

$\delta_j^{(l)} = $ "error" of node $j$ in layer $l$

Formally, $\delta_j^{(l)} = \dfrac{\partial}{\partial z_j^{(l)}} \mathrm{cost}(\mathbf{x}_i)$

where $\mathrm{cost}(\mathbf{x}_i) = [y_i \log h_\Theta(\mathbf{x}_i) + (1 - y_i) \log(1 - h_\Theta(\mathbf{x}_i))]$

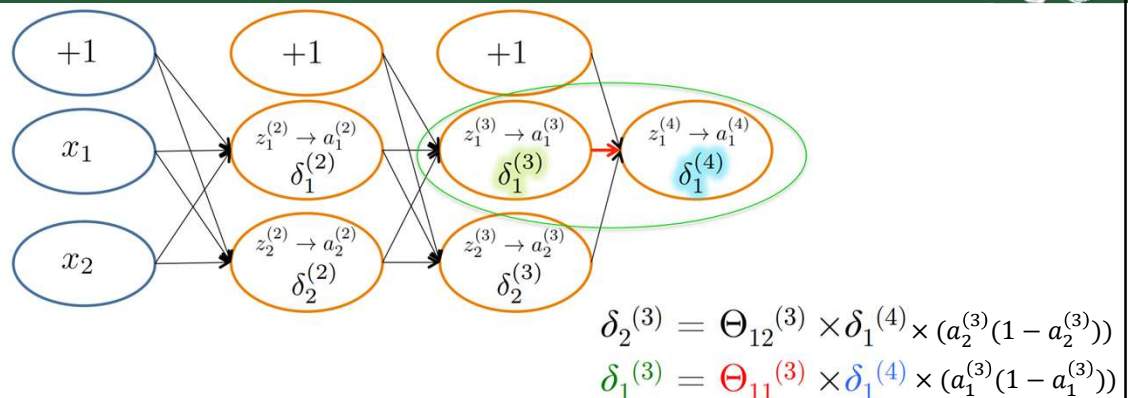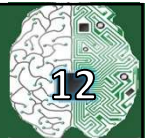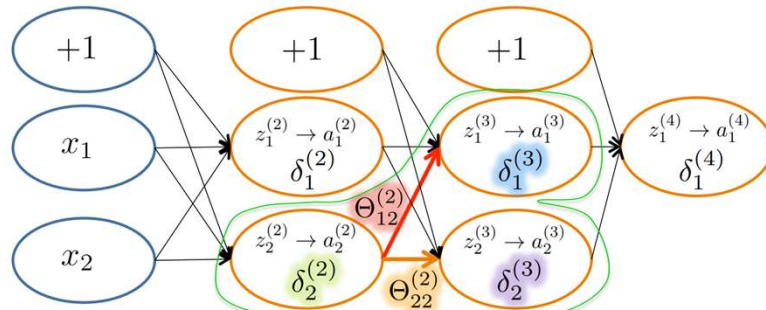Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Backpropagation Intuition

**12**

$$+1 \qquad +1 \qquad +1$$

$$x_1 \qquad \begin{array}{c} z_1^{(2)} \to a_1^{(2)} \\ \delta_1^{(2)} \end{array} \qquad \begin{array}{c} z_1^{(3)} \to a_1^{(3)} \\ \delta_1^{(3)} \end{array} \qquad \begin{array}{c} z_1^{(4)} \to a_1^{(4)} \\ \delta_1^{(4)} \end{array}$$

$$x_2 \qquad \begin{array}{c} z_2^{(2)} \to a_2^{(2)} \\ \delta_2^{(2)} \end{array} \qquad \begin{array}{c} z_2^{(3)} \to a_2^{(3)} \\ \delta_2^{(3)} \end{array}$$

$$\delta_2^{(3)} = \Theta_{12}^{(3)} \times \delta_1^{(4)} \times (a_2^{(3)}(1 - a_2^{(3)}))$$
$$\delta_1^{(3)} = \Theta_{11}^{(3)} \times \delta_1^{(4)} \times (a_1^{(3)}(1 - a_1^{(3)}))$$

$\delta_j^{(l)} = $ "error" of node $j$ in layer $l$

Formally, $\delta_j^{(l)} = \dfrac{\partial}{\partial z_j^{(l)}} \mathrm{cost}(\mathbf{x}_i)$

where $\mathrm{cost}(\mathbf{x}_i) = [y_i \log h_\Theta(\mathbf{x}_i) + (1 - y_i) \log(1 - h_\Theta(\mathbf{x}_i))]$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Backpropagation Intuition

$$\delta_2^{(2)} = (\Theta_{12}^{(2)} \times \delta_1^{(3)} + \Theta_{22}^{(2)} \times \delta_2^{(3)}) \times (a_2^{(2)}(1 - a_2^{(2)}))$$
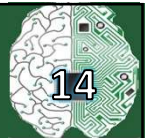
$\delta_j^{(l)}$ = "error" of node $j$ in layer $l$

Formally, $\delta_j^{(l)} = \dfrac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

where $\text{cost}(\mathbf{x}_i) = [y_i \log h_\Theta(\mathbf{x}_i) + (1 - y_i) \log(1 - h_\Theta(\mathbf{x}_i))]$

# Backpropagation Intuition: Gradient Computation

Let $\delta_j^{(l)}$ = "error" of node $j$ in layer $l$

(#layers $L$ = 4)

Element-wise product .*

### Backpropagation

- $\boldsymbol{\delta}^{(4)} = \boldsymbol{a}^{(4)} - \mathbf{y}$
- $\boldsymbol{\delta}^{(3)} = (\Theta^{(3)})^\mathsf{T} \boldsymbol{\delta}^{(4)} .* g'(\mathbf{z}^{(3)})$
- $\boldsymbol{\delta}^{(2)} = (\Theta^{(2)})^\mathsf{T} \boldsymbol{\delta}^{(3)} .* g'(\mathbf{z}^{(2)})$
- (No $\boldsymbol{\delta}^{(1)}$)
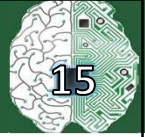
$g'(\mathbf{z}^{(3)}) = \mathbf{a}^{(3)} .* (1 - \mathbf{a}^{(3)})$

$g'(\mathbf{z}^{(2)}) = \mathbf{a}^{(2)} .* (1 - \mathbf{a}^{(2)})$

$\boldsymbol{\delta}^{(2)}$ $\boldsymbol{\delta}^{(3)}$ $\boldsymbol{\delta}^{(4)}$

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)}$$

(ignoring $\lambda$; if $\lambda = 0$)

# Backpropagation

Set $\Delta_{ij}^{(l)} = 0 \quad \forall l, i, j$        (Used to accumulate gradient)

For each training instance $(\mathbf{x}^{(s)}, y^{(s)})$:

     Set $\mathbf{a}^{(1)} = \mathbf{x}^{(s)}$

     Compute $\{\mathbf{a}^{(2)}, \ldots, \mathbf{a}^{(L)}\}$ via forward propagation

     Compute $\boldsymbol{\delta}^{(L)} = \mathbf{a}^{(L)} - y^{(s)}$

     Compute errors $\{\boldsymbol{\delta}^{(L-1)}, \ldots, \boldsymbol{\delta}^{(2)}\}$

     Compute gradients $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

Compute avg regularized gradient $D_{ij}^{(l)} = \begin{cases} \frac{1}{n}\Delta_{ij}^{(l)} + \lambda\Theta_{ij}^{(l)} & \text{if } j \neq 0 \\ \frac{1}{n}\Delta_{ij}^{(l)} & \text{otherwise} \end{cases}$

$\boldsymbol{D}^{(l)}$ is the matrix of partial derivatives of $J(\Theta)$

Note: Can vectorize $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$ as $\boldsymbol{\Delta}^{(l)} = \boldsymbol{\Delta}^{(l)} + \boldsymbol{\delta}^{(l+1)} \mathbf{a}^{(l)\mathsf{T}}$

---

# Training a Neural Network via Gradient Descent with Backprop

Given: training set $\{(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_n, y_n)\}$

Initialize all $\Theta^{(l)}$ randomly (NOT to 0!)

Loop // each iteration is called an epoch

     Set $\Delta_{ij}^{(l)} = 0 \quad \forall l, i, j$      (Used to accumulate gradient)

     For each training instance $(\mathbf{x}^{(s)}, y^{(s)})$:

         Set $\mathbf{a}^{(1)} = \mathbf{x}^{(s)}$

         Compute $\{\mathbf{a}^{(2)}, \ldots, \mathbf{a}^{(L)}\}$ via forward propagation

         Compute $\boldsymbol{\delta}^{(L)} = \mathbf{a}^{(L)} - y^{(s)}$

         Compute errors $\{\boldsymbol{\delta}^{(L-1)}, \ldots, \boldsymbol{\delta}^{(2)}\}$

         Compute gradients $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

     Compute avg regularized gradient $D_{ij}^{(l)} = \begin{cases} \frac{1}{n}\Delta_{ij}^{(l)} + \lambda\Theta_{ij}^{(l)} & \text{if } j \neq 0 \\ \frac{1}{n}\Delta_{ij}^{(l)} & \text{otherwise} \end{cases}$

     Update weights via gradient step $\Theta_{ij}^{(l)} = \Theta_{ij}^{(l)} - \alpha D_{ij}^{(l)}$

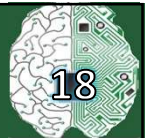Until weights converge or max #epochs is reached

*Backpropagation*

# More on Optimization Algorithm
## [with adaptive learning rate]

---

# AdaGrad

- Individually adapts learning rates of all parameters
  - Scale them inversely proportional to the sum of the historical squared values of the gradient

    Performs well for some but not all deep learning

**The AdaGrad Algorithm:**

**Require:** Global learning rate $\alpha$

**Require:** Initial parameter $\boldsymbol{\theta}$

**Require:** Small constant $\delta$, perhaps $10^{-7}$, for numerical stability

Initialize gradient accumulation variable $\boldsymbol{r} = \boldsymbol{0}$

**while** stopping criterion not met **do**

Sample a minibatch of $m$ examples from the training set $\{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(m)}\}$ with corresponding targets $\boldsymbol{y}^{(i)}$.

Compute gradient: $\boldsymbol{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i L(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), \boldsymbol{y}^{(i)})$

Accumulate squared gradient: $\boldsymbol{r} \leftarrow \boldsymbol{r} + \boldsymbol{g} \odot \boldsymbol{g}$

Compute update: $\Delta \boldsymbol{\theta} \leftarrow -\frac{\alpha}{\delta + \sqrt{\boldsymbol{r}}} \odot \boldsymbol{g}$.    (Division and square root applied element-wise)

Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta \boldsymbol{\theta}$

**end while**

# RMSProp

- RMSProp is an effective practical optimization algorithm

- Common optimization method for deep learning practitioners

- Modifies AdaGrad for a nonconvex setting

  – Changes gradient accumulation into a exponentially weighted moving average

  – Converges rapidly when applied to a convex function

**The RMSProp Algorithm**

**Require:** Global learning rate $\alpha$, decay rate $\rho$.
**Require:** Initial parameter $\theta$
**Require:** Small constant $\delta$, usually $10^{-6}$, used to stabilize division by small numbers.

Initialize accumulation variables $r = 0$
**while** stopping criterion not met **do**
    Sample a minibatch of $m$ examples from the training set $\{x^{(1)},\ldots,x^{(m)}\}$ with corresponding targets $y^{(i)}$.
    Compute gradient: $g \leftarrow \frac{1}{m}\nabla_\theta \sum_i L(f(x^{(i)};\theta),y^{(i)})$
    Accumulate squared gradient: $r \leftarrow \rho r + (1-\rho)g \odot g$
    Compute parameter update: $\Delta\theta = -\frac{\alpha}{\sqrt{\delta+r}} \odot g$.   ($\frac{1}{\sqrt{\delta+r}}$ applied element-wise)
    Apply update: $\theta \leftarrow \theta + \Delta\theta$
**end while**

# RMSProp Combined with Nesterov

**Require:** Global learning rate $\alpha$, decay rate $\rho$, momentum coefficient $\delta$.
**Require:** Initial parameter $\theta$, initial velocity $v$.
Initialize accumulation variable $r = 0$
**while** stopping criterion not met **do**
    Sample a minibatch of $m$ examples from the training set $\{x^{(1)},\ldots,x^{(m)}\}$ with corresponding targets $y^{(i)}$.
    Compute interim update: $\tilde{\theta} \leftarrow \theta + \delta v$
    Compute gradient: $g \leftarrow \frac{1}{m}\nabla_{\tilde{\theta}} \sum_i L(f(x^{(i)};\tilde{\theta}),y^{(i)})$
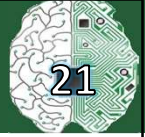    Accumulate gradient: $r \leftarrow \rho r + (1-\rho)g \odot g$
    Compute velocity update: $v \leftarrow \delta v - \frac{\alpha}{\sqrt{r}} \odot g$.   ($\frac{1}{\sqrt{r}}$ applied element-wise)
    Apply update: $\theta \leftarrow \theta + v$
**end while**

# Adam: Adaptive Moments

**Adam Algorithm**

**Require:** Step size $\alpha$ (Suggested default: 0.001)
**Require:** Exponential decay rates for moment estimates, $\rho_1$ and $\rho_2$ in $[0, 1)$. (Suggested defaults: 0.9 and 0.999 respectively)
**Require:** Small constant $\delta$ used for numerical stabilization. (Suggested default: $10^{-8}$)
**Require:** Initial parameters $\boldsymbol{\theta}$

    Initialize 1st and 2nd moment variables $\boldsymbol{s} = \boldsymbol{0}$, $\boldsymbol{r} = \boldsymbol{0}$
    Initialize time step $t = 0$
    **while** stopping criterion not met **do**
        Sample a minibatch of $m$ examples from the training set $\{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(m)}\}$ with corresponding targets $\boldsymbol{y}^{(i)}$.
        Compute gradient: $\boldsymbol{g} \leftarrow \frac{1}{m}\nabla_{\boldsymbol{\theta}} \sum_i L(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), \boldsymbol{y}^{(i)})$
        $t \leftarrow t + 1$
        Update biased first moment estimate: $\boldsymbol{s} \leftarrow \rho_1\boldsymbol{s} + (1 - \rho_1)\boldsymbol{g}$
        Update biased second moment estimate: $\boldsymbol{r} \leftarrow \rho_2\boldsymbol{r} + (1 - \rho_2)\boldsymbol{g} \odot \boldsymbol{g}$
        Correct bias in first moment: $\hat{\boldsymbol{s}} \leftarrow \frac{\boldsymbol{s}}{1-\rho_1^t}$
        Correct bias in second moment: $\hat{\boldsymbol{r}} \leftarrow \frac{\boldsymbol{r}}{1-\rho_2^t}$
        Compute update: $\Delta\boldsymbol{\theta} = -\alpha\frac{\hat{\boldsymbol{s}}}{\sqrt{\hat{\boldsymbol{r}}}+\delta}$   (operations applied element-wise)
        Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta\boldsymbol{\theta}$
    **end while**

# Choosing the Right Optimizer

- We have discussed several methods of optimizing deep models by adapting the learning rate for each model parameter

- Which algorithm to choose?
  - No consensus

- Most popular algorithms actively in use:
  - SGD, SGD with momentum, RMSProp, RMSProp with momentum, AdaDelta and Adam
  - Choice depends on user's familiarity with algorithm
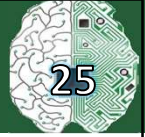
# Improved Learning

23

# Gradient Checking

24

- A method used during the implementation of the backward pass of a neural network. It compares the value of the analytical gradient to the numerical gradient at given points and plays the role of a sanity-check for correctness.

| Numerical gradient | Analytical gradient |
|---|---|
| $$\frac{df}{dx}(x) \approx \frac{f(x+h) - f(x-h)}{2h}$$ | $$\frac{df}{dx}(x) = f'(x)$$ |
| • Expensive; loss has to be computed two times per dimension<br>• Used to verify correctness of analytical implementation<br>• Trade-off in choosing $h$ not too small (numerical instability) nor too large (poor gradient approximation) | • 'Exact' result<br>• Direct computation<br>• Used in the final implementation |

# Training a Neural Network

1. Randomly initialize weights
2. Implement forward propagation to get $h_\Theta(\mathbf{x}_i)$ for any instance $\mathbf{x}_i$
3. Implement code to compute cost function $J(\Theta)$
4. Implement backprop to compute partial derivatives $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$
5. Use gradient checking to compare $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$ computed using backpropagation vs. the numerical gradient estimate.
    - Then, disable gradient checking code
6. Use gradient descent with backprop to fit the network

# Handling Overfitting: Early Stopping

- Most commonly used in deep learning
- Popularity is due to its effectiveness and its simplicity

In the case of early stopping we are controlling the effective capacity of the model by determining how many steps it can take to fit the training set

- Can think of early stopping as a very efficient hyperparameter selection algorithm
    - In this view number of training steps is just a hyperparameter
    - This hyperparameter has a U-shaped validation set performance curve
    - Most hyperparameters have such a U-shaped validation set performance curve

# Costs of Early Stopping

- Cost of this hyperparameter is doing the validation evaluation periodically during training
  - Ideally done in parallel to the training process on a separate machine
    - Separate CPU or GPU from main training process
    - Or use a smaller validation set or validating less frequently

- Need to maintain a copy of the best parameters
  - This cost is negligible because they can be stored on a slower, larger memory
    - E.g., training on GPU, but storing the optimal parameters in host memory or elsewhere

# Early Stopping as Regularization

- Early stopping is an unobtrusive form of regularization
- It requires almost no change to the underlying training procedure, the objective function, or the set of allowable parameter values



- As such it is easy to use early stopping without damaging the learning dynamics
  - In contrast to weight decay, where we must be careful not to use too much weight decay
    - Otherwise we trap the network in a bad local minimum corresponding to pathologically small weights

# Dropout Neural Net

- A simple way to prevent neural net overfitting



(a) Standard Neural Net

(b) After applying dropout.

Drop hidden and visible units from net, i.e., temporarily remove from the network with all input/output connections.
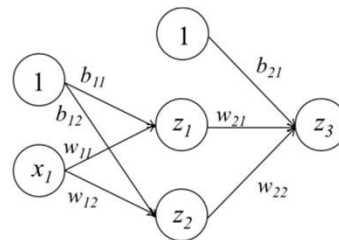
Choice of units to drop is random, determined by a probability $p$, chosen by a validation set, or equal to $0.5$

# Exercise

**Q1.** Consider the neural network architecture shown below for a binary classification problem.



We define:

$$a_1 = w_{11}x_1 + b_{11}$$
$$a_2 = w_{12}x_1 + b_{12}$$
$$a_3 = w_{21}z_1 + w_{22}z_2 + b_{21}$$
$$z_1 = ReLU(a_1)$$
$$z_2 = ReLU(a_2)$$
$$z_3 = g(a_3)$$
$$g(x) = \frac{1}{1 + e^{-x}}$$

Using the above neural network, write the partial derivative of the loss function $L(y, z_3)$ with respect to the bias term $b_{12}$, $\frac{\partial L}{\partial b_{12}}$, in terms of the partial derivatives $\frac{\partial \alpha}{\partial \beta}$, where $\alpha$ and $\beta$ can be any of $L, z_i, a_i, b_{ij}, w_{ij}, x_1$ for all valid values of i,j. Your answer should be as explicit as possible—that is, make sure each partial derivative $\frac{\partial \alpha}{\partial \beta}$ cannot be decomposed further into simpler partial derivatives. Do NOT evaluate the partial derivatives.

# Exercise

**Q2.** These plots were generated with gradient descent; with gradient descent with momentum ($\delta$ = 0.5) and gradient descent with momentum ($\delta$ = 0.9). Which curve corresponds to which algorithm?
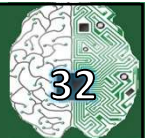


(1) is generated by?
(2) is generated by?
(3) is generated by?

**Q3.** Among the following associated with a neural network, which one(s) is/are the hyperparameter?

    A. No. of iteration
    B. Bias vectors
    C. Weight matrices
    D. Number of layers in the network
    E. Activation values

# Convolutional Neural Network:
## Motivation

# Importance of Local Features

- After training, the network parameters Z and V will become as follows :

# Importance of Local Features

This template will only be useful for this image.
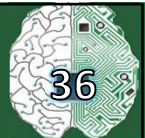
# Importance of Local Features

- **Local features** model the biological local reception fields :



- Sliding templates :
  - Define small templates and scan them across the image domain.
  - This operation is known as convolution.

---

# Questions?