# Lecture 19 contd: Neural Network Training using Backpropagation, Convolutional And Recurrent Neural Networks
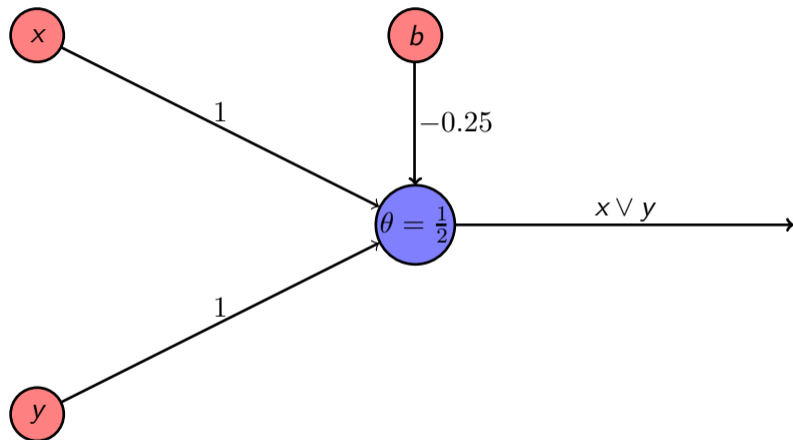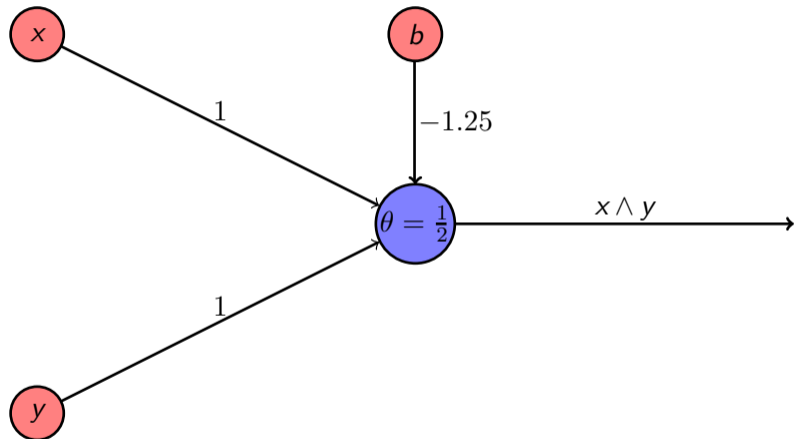
Instructor: Prof. Ganesh Ramakrishnan
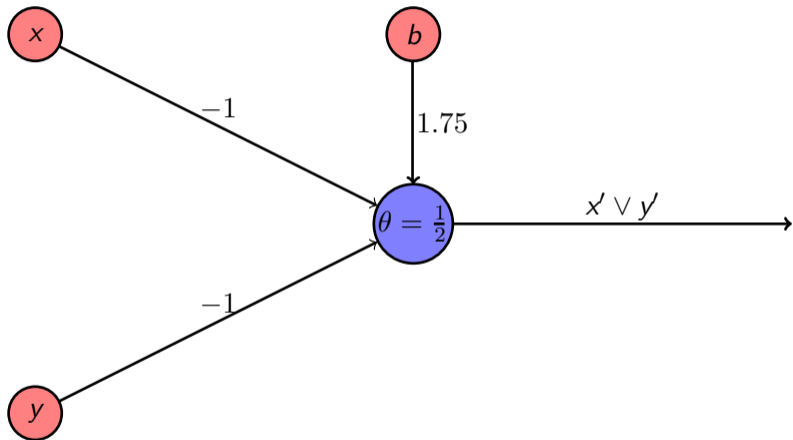
# OR using perceptron

# AND using perceptron

# x'Vy' using perceptron

How about XOR $((x' \lor y') \land (x \lor y))$?

$\underline{\neg x}$ $\underline{\neg y}$



functional composition
of linear seperators
$\equiv$ Cascade of perceptrons

How about (composing)
linear seperators?

$P\left(P(\bar{x} \lor \bar{y}) \land P(x \lor y)\right)$

# Cascade of Perceptrons for XOR($(x' \vee y') \wedge (x \vee y)$ with $\theta = 0$ )



inputs

Multiple outputs
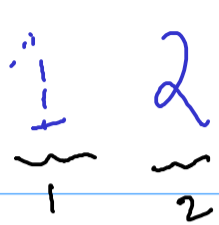
# Feed-forward Neural Nets

inputs

# Training a Neural Network

STEP 0: Pick a network architecture *[handwritten: should be → correctly suggested either by data or by expert]*

- Number of input units: Dimension of features $\phi\left(\mathbf{x}^{(i)}\right)$.
- Number of output units: Number of classes. → *[handwritten: $y_1, y_2 \cdots y_k \in \{0,1\}$ .. Designed for multi labeling]*
- 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 $x$. *[handwritten: Rule of thumb]*
- 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.

$1 \quad 2$  } Multiple class, Single label (per example)

$\underbrace{\phantom{1}}_{1} \quad \underbrace{\phantom{2}}_{2}$

[20 NewsGroups Dataset]



} Multiple class, Multilabel

ML  ↙ ↘ Text mining

M2  ↙ ↘ Car driving

} Multi instance Multi label (MIML)

labels 1,2 ← { _ _ _ } → set of instances ← { _ _ _ } → labels 3,4

Input layer

Hidden layers

output layer

$\phi_1(x)$

$\phi_2(x)$

$\phi_3(x)$

$\phi_p(x)$

# **High Level Overview** of Backpropagation Algorithm for Training NN



1. Randomly initialize weights $w_{ij}^l$ for $l = 1, \ldots, L$, $i = 1, \ldots, s_l$, $j = 1, \ldots, s_{l+1}$.
2. Implement **forward propagation** to get $f_w(\mathbf{x})$ for any $x \in \mathcal{D}$. [eg: For XOR]
3. Execute **backpropagation**
   1. by computing partial derivatives $\frac{\partial}{\partial w_{ij}^{(l)}} E(w)$ for $l = 1, \ldots, L$, $i = 1, \ldots, s_l$, $j = 1, \ldots, s_{l+1}$.
   2. and using gradient descent to try to minimize (non-convex) $E(w)$ as a function of parameters $\mathbf{w}$.
4. Verify that the cost function $E(w)$ has indeed reduced, else resort to some random perturbation of weights $\mathbf{w}$.

later

# Intuition for Backpropagation

$f(x, y, z) = (w_x x + w_y y)w_q + w_z z, \quad q = w_x x + w_y y$

At input (-3,7,-5):-  → *from $(\ell-1)$ layers*

(gradient of f w.r.t. v shown in red where $v \in \{w_x, w_y, w_z, w_q\}$)



*$(\ell-1)^{th}$ layers*

*$\ell^{th}$ layers*

*$(\ell+1)$ layer*

x    -3
     -5

y    7
     -5

     q 4
     -5

z    -5
     4

     f -20
     1

*from current (1) to desired (-20)*

Should we decrease $w_y$ or should we increase $w_x$ to decrease f?

# Intuition for Backpropagation



Iteratively update the weights along the direction where Loss decreases.

# Setting Notation for Backpropagation



$S_L = $ # of perceptrons in $l^{th}$ layer

$\sigma = $ activation fn (sigmoid)

Circular nodes are perceptrons

$\sigma_1^l = \sigma\left(sum_1^l\right)$

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

$\sigma_1^{l-1}$

$\sigma_2^{l-1}$

$\sigma_i^{l-1}$

$\sigma_{s_{l-1}}^{l-1}$

$w_{11}^l$

$w_{21}^l$

$w_{i1}^l$

$w_{s_{l-1}1}^l$

$w_{1j}^l$

$w_{2j}^l$

$w_{ij}^l$

$w_{s_{l-1}j}^l$

$\sigma_j^l = \sigma\left(sum_j^l\right)$

$\sigma_1^L$ → $L = $ final or output layer.

$\sigma_K^L$

$K = $ # of output nodes.

## Gradient Computation

- The Neural Network objective to be minimized:

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

$$+ \frac{\lambda}{2m}\sum_{l=1}^{L}\sum_{i=1}^{s_{l-1}}\sum_{j=1}^{s_l}\left(w_{ij}^l\right)^2$$

*output at final ($L^{th}$) layer*

*Summation over all examples*

*summation over all outputs*

*Regularizer computed on weights across all layers.*

Cross entropy loss component computed only in terms of final output layer & input nodes

$$\left(x^{(i)}, y^{(i)}\right)$$

Expand on: $w_{ij}^{l,(k+1)} = w_{ij}^{l,(k)} - \left(\frac{\partial E_g}{\partial w_{ij}}\right)_{\omega^{(k)}}$
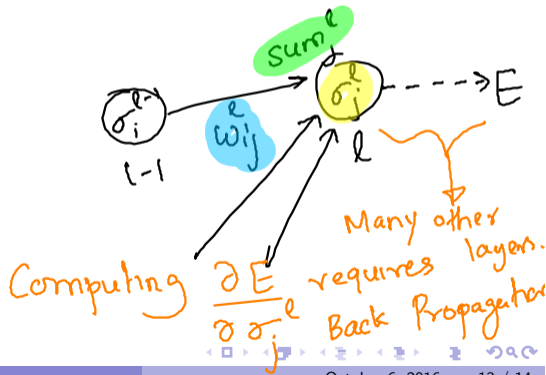
# Gradient Computation

- The Neural Network objective to be minimized:

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

$$+ \frac{\lambda}{2m}\sum_{l=1}^{L}\sum_{i=1}^{s_{l-1}}\sum_{j=1}^{s_l}\left(w_{ij}^l\right)^2 \tag{1}$$
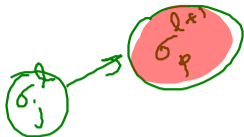
In $\frac{\partial E}{\partial \sigma_j^l}$ if $l = L$ then directly differentiate

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

- $\dfrac{\partial E}{\partial w_{ij}^l} = \dfrac{\partial E}{\partial \sigma_j^l}\dfrac{\partial \sigma_j^l}{\partial sum_j^l}\dfrac{\partial sum_j^l}{\partial w_{ij}^l}$ → chain rule

- $\dfrac{\partial sum_j^l}{\partial w_{ij}^l} = \dfrac{\partial}{\partial w_{ij}^l}\left(\sum\limits_{k=1}^{s_{l-1}} w_{kj}^l \sigma_k^{l-1}\right) = \sigma_i^{l-1}$

- $\dfrac{\partial \sigma_j^l}{\partial sum_j^l} = \left(\dfrac{1}{1+e^{-sum_i^l}}\right)\left(1 - \dfrac{1}{1+e^{-sum_i^l}}\right)$

Many other layers.

Computing $\frac{\partial E}{\partial \sigma_j^l}$ requires Back Propagation

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

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

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L} \sum_{i=1}^{s_{l-1}} \sum_{j=1}^{s_l} \left( w_{ij}^l \right)^2 \tag{2}$$

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

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

*Recall from logistic update*