

Lecture 21: Neural Networks, Back-propagation, etc

Instructor: Prof. Ganesh Ramakrishnan

Non-linear perceptron?

- Kernelized perceptron: $f(x) = \text{sign} \left(\sum_i \alpha_i y_i K(x, x_i) + b \right)$
 - INITIALIZE: $\alpha = \text{zeroes}()$
 - REPEAT: for $\langle x_i, y_i \rangle$
 - If $\text{sign} \left(\sum_j \alpha_j y_j K(x_j, x_i) + b \right) \neq y_i$
 - then, $\alpha_i = \alpha_i + 1$
 - endif
- Neural Networks: Cascade of layers of perceptrons giving you non-linearity
 - $\text{sign}((w^*)^T \phi(x))$ replaced by $g((w^*)^T \phi(x))$ where $g(s)$ is a
 - 1 step function: $g(s) = 1$ if $s \in [\theta, \infty)$ and $g(s) = 0$ otherwise OR
 - 2 sigmoid function: $g(s) = \frac{1}{1+e^{-s}}$ with possible thresholding using some θ (such as $\frac{1}{2}$).
 - 3 Rectified Linear Unit (ReLU): $g(s) = \max(0, s)$: A most popular activation function
 - 4 Softplus: $g(s) = \ln(1 + e^s)$

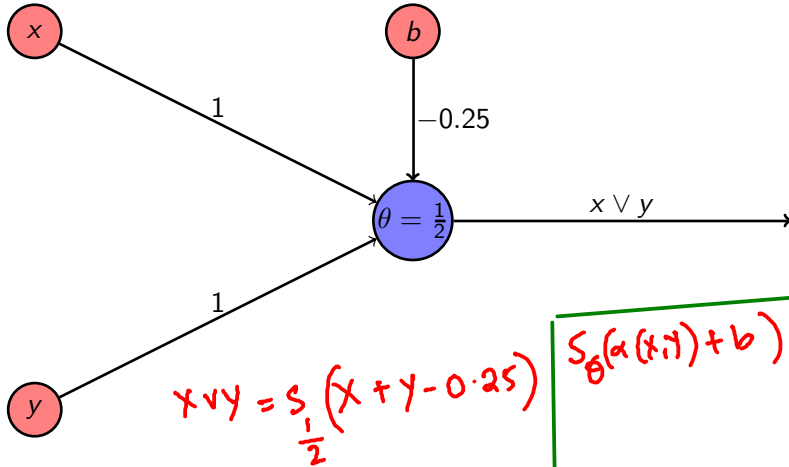
Options 2, 3 and 4 have the thresholding step deferred. Threshold changes as bias is changed.

Recall: Measure for Linear non-separability?

- **Aspect 1: Number of functions that can be represented** Recall from Tutorial 6, Problem 1: Given n boolean variables how many of 2^{2^n} boolean functions can be represented by a perceptron? Ans: For 2 it is 14, for 3 it is 104, for 4 it is 1882
- **Aspect 2: Cardinality of largest set of points that can be shattered**
VC (VapnikChervonenkis) dimension \Rightarrow A measure of the richness of a space of functions that can be learned by a statistical classification algorithm.
 - A classification function $f(\mathbf{w})$ is said to shatter a set of data points (x_1, x_2, \dots, x_n) if, for all assignments of labels to those points, there exists a \mathbf{w} such that $f(\mathbf{w})$ makes no errors when evaluating that set of data points.
 - Cardinality of the largest set of points that $f(\mathbf{w})$ can **shatter** is its VC-dimension.
 - Eg: For f as a threshold interval, VC dimension = 1
 - Eg: For f as an interval classifier, VC dimension = 2
 - Eg: For f as linear classifier (in 2 dimensions), VC dimension = 3
 - Eg: For f as linear classifier (in \mathbb{R}^n), VC dimension = $n + 1$
 - Eg: For f as a neural network with sigmoid function, V nodes and E edges, then the VC dimension is at least $\Omega(|E|^2)$ and at most $O(|E|^2 \cdot |V|^2)$

OR using (step) perceptron

(For AND, $b = -1.25$)



$$x \vee y = s_{\frac{1}{2}}(x + y - 0.25)$$

$$s_{\theta}(a(x,y) + b) = \text{step fn} = 1$$

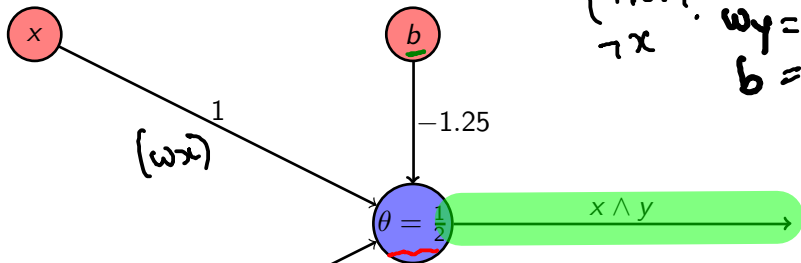
if $a + b > \theta$

$= 0$

o/w.

AND using (step) perceptron

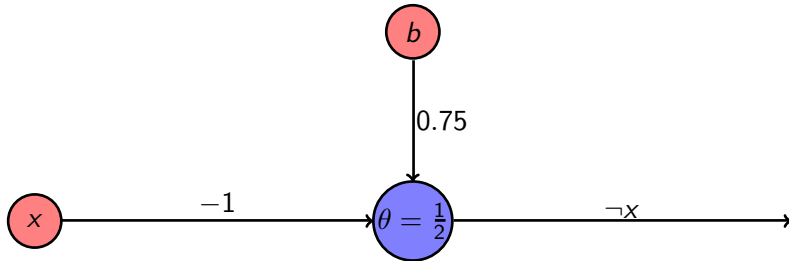
(NOT? $\neg x$ $w_x = -1$
 $w_y = 0$
 $b = 0.75$)



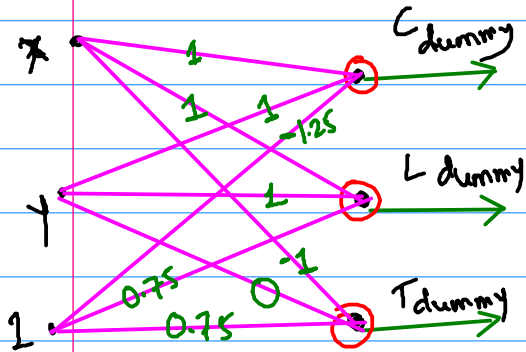
Note: θ restricted to $\frac{1}{2}$
to retain probabilistic
interpretation

1	1	1	$(2) + b > \frac{1}{2}$
1	0	0	$(1) + b < \frac{1}{2}$
0	1	0	$(1) + b < \frac{1}{2}$
0	0	0	$(0) + b < \frac{1}{2}$

NOT using perceptron

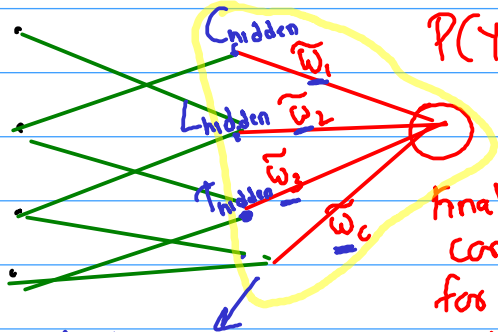


$(x \wedge y)$ $(x \vee y)$ $(\neg x)$
 C_{dummy} L_{dummy} T_{dummy}



•: Step fn with $\Theta = 1/2$
 (or sigmoid fn with $\Theta = 1/2$)

In practice when only one of k classes should be predicted, use a softmax generalization of sigmoid as last layer

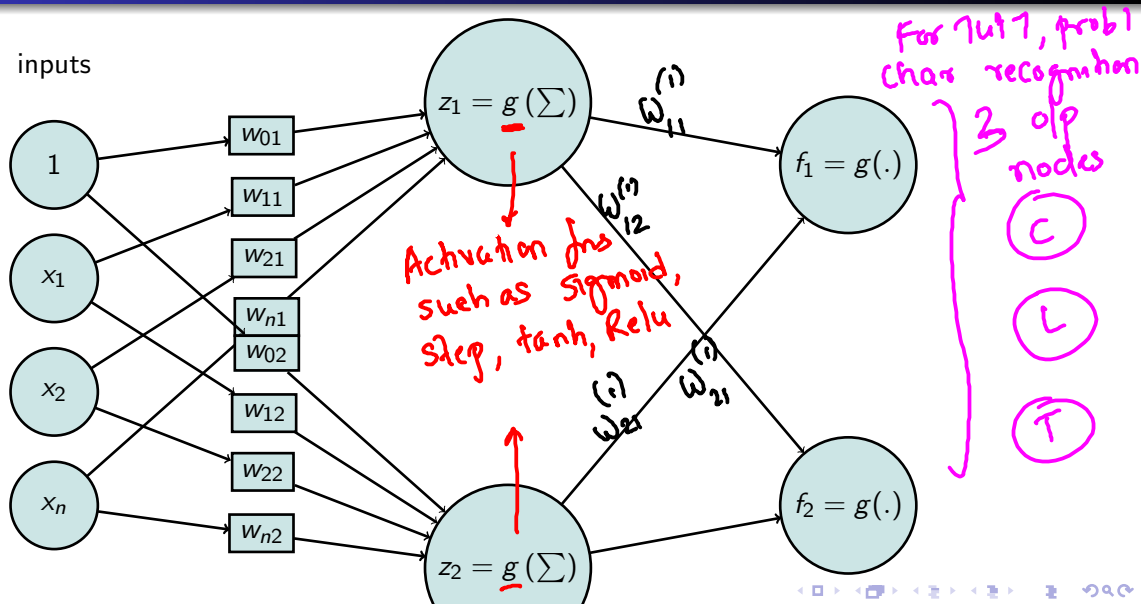


$$P(Y=c|x) = \frac{e^{-\tilde{w}_c \phi(x)}}{\sum_{c=1}^k e^{-\tilde{w}_c \phi(x)}}$$

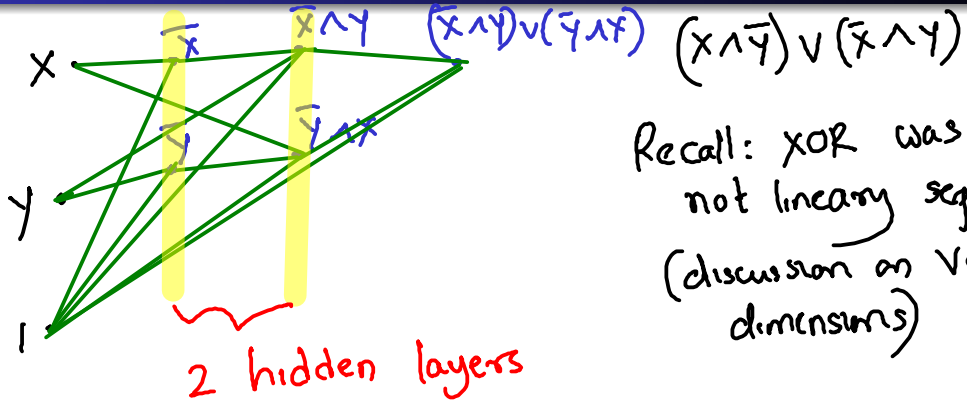
Final softmax
consolidation
for single label,
multiclass classification

At penultimate layer, you get
multilabel multiclass classification

Feed-forward Neural Nets



Eg: Feed-forward Neural Net for XOR ($\theta = 0$)

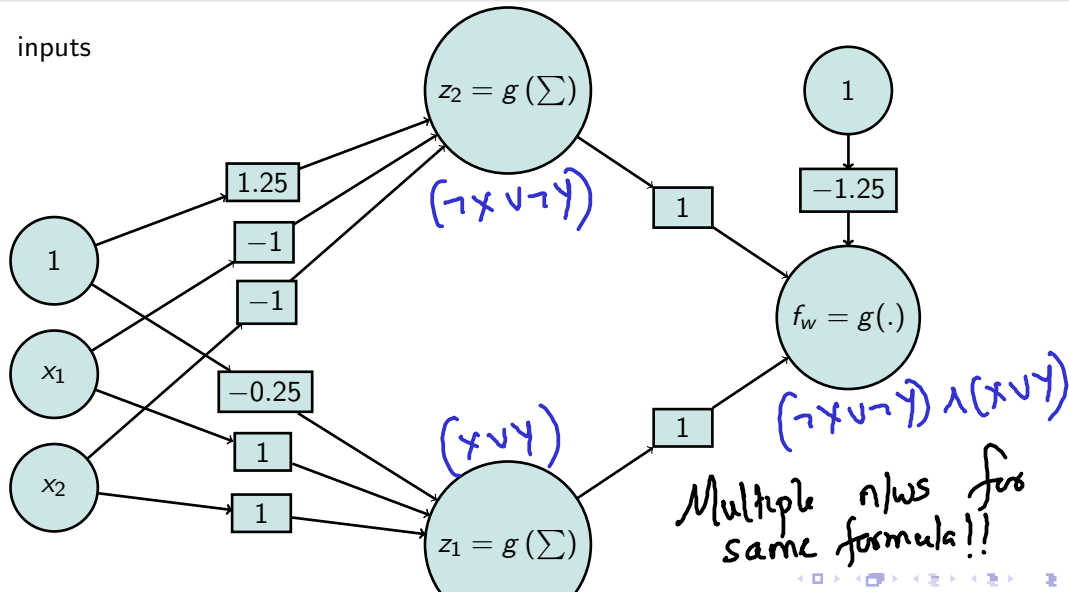


Recall: XOR was
not linearly separable
(discussion on VC
dimensions)

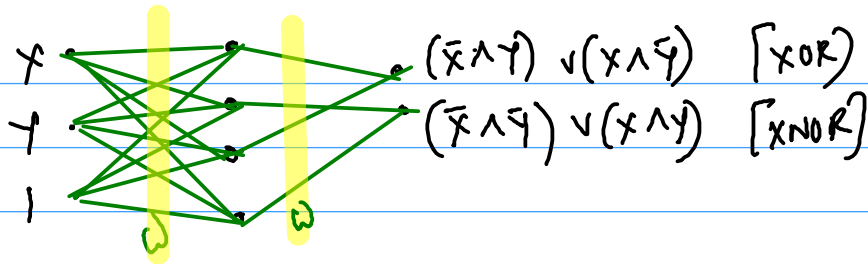
How about fewer hidden layers for XOR?
Ans: $x \wedge \bar{y}$ & $\bar{x} \wedge y$ are representable using linear separator!

Eg: Feed-forward Neural Net for XOR ($\theta = 0$)

inputs



For XOR & XNOR in same n/w



Q: How to "train" the network (i.e. learn the w 's)?

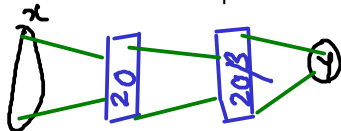
Training a Neural Network

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 \mathbf{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 (using [stochastic] gradient desc)**
- instead of the non-smooth step function $g(s) = 1$ if $s \in [\theta, \infty)$ and $g(s) = 0$ otherwise.



Tutorial 7, problem 5

ϕ is left
to be \mathbf{x}
(trivial ϕ)

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 $\mathbf{x} \in \mathcal{D}$.

③ Execute **backpropagation**

① by computing partial derivatives $\frac{\partial}{\partial w_{ij}^{(l)}} E(w)$ for $l = 1, \dots, L$, $i = 1, \dots, s_l$,

$j = 1, \dots, s_{l+1}$.

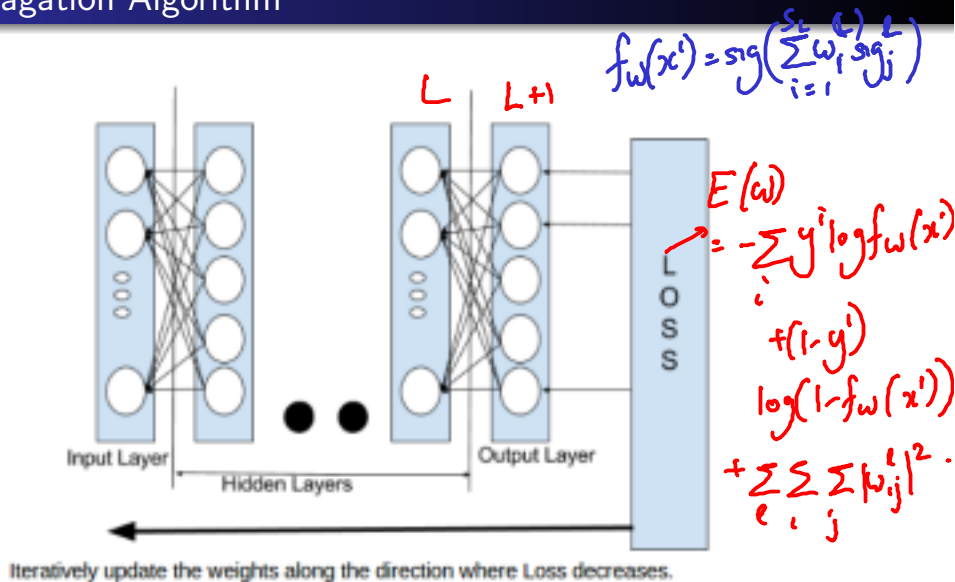
② and using gradient descent to try to minimize (non-convex) $E(w)$ as a function of parameters \mathbf{w} .

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

→ key step

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

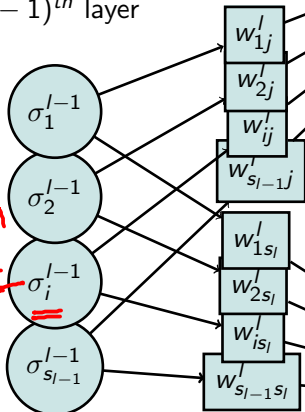
The Backpropagation Algorithm



Setting Notation for Backpropagation

$$sum_j^l = \sum_{i=1}^{s_{l-1}} \sigma_i^{(l-1)} w_{ij}$$

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



$$\sigma_j^l = \sigma(\underline{sum_j^l})$$

$$\sigma_{s_l}^l = \sigma(sum_{s_l}^l)$$

Step by step (backward) procedure to obtain $\frac{\partial E}{\partial w_{ij}^l}$ in terms of

$$\frac{\partial E}{\partial w_{ij}^{l+1}} \dots \frac{\partial E}{\partial w_{ij}^L}$$

(chain rule)

$$\sigma_1^l$$

$$\sigma_K^l$$

Gradient Computation

- The Neural Network objective to be minimized:

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

easy to handle

$$\frac{\partial E(\mathbf{w})}{\partial w_{ij}^{(l)}} = \sum_{i=1}^m \sum_{k=1}^K \frac{\partial E(\mathbf{w})}{\partial \sigma_k^L(\mathbf{x}^{(i)})} \sum_{j=1}^{s_{L-1}} \frac{\partial \sigma_k^L(\mathbf{x}^{(i)})}{\partial \text{sum}_j^{L-1}(\mathbf{x}^{(i)})} \sum_{p=1}^{s_{L-1}} \frac{\partial \text{sum}_j^{L-1}}{\partial \sigma_p^{L-1}} \dots \frac{\partial \text{sum}_j^{(l)}}{\partial w_{ij}^{(l)}}$$

Gradient Computation

- The Neural Network objective to be minimized:

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

- $sum_j^l = \sum_{k=1}^{s_{l-1}} w_{kj}^l \sigma_k^{l-1}$ and $\sigma_i^l = \frac{1}{1+e^{-sum_i^l}}$
- $\frac{\partial E}{\partial w_{ij}^l} = \frac{\partial E}{\partial \sigma_j^l} \frac{\partial \sigma_j^l}{\partial sum_j^l} \frac{\partial sum_j^l}{\partial w_{ij}^l} + \frac{\lambda}{2m} w_{ij}^l$
- $\frac{\partial \sigma_j^l}{\partial sum_j^l} = \left(\frac{1}{1+e^{-sum_j^l}} \right) \left(1 - \frac{1}{1+e^{-sum_j^l}} \right) = \sigma_j^l (1 - \sigma_j^l)$
- $\frac{\partial sum_j^l}{\partial w_{ij}^l} = \frac{\partial}{\partial w_{ij}^l} \left(\sum_{k=1}^{s_{l-1}} w_{kj}^l \sigma_k^{l-1} \right) = \sigma_i^{l-1}$