

4/11/15

# (FINISH) SVMs + NEURAL NETWORKS

## ① Multi-class SVM

→ So far: Binary SVM  $y \in \{-1, +1\}$

→ For multiclass  $y \in \{1, 2, \dots, K\}$ ; your options:

→ One-vs-Rest (All)  $K$ -binary SVMs

→ One-vs-One  $\binom{K}{2}$ -binary SVMs

→ [Crammer & Singer JMLR'01]  
1-natively multi-class SVM

$K$  weight vectors  $\{\vec{w}^{(1)}, \dots, \vec{w}^{(K)}\}$   $\vec{w}^{(i)} \in \mathbb{R}^d$

$K$  "scores" for each input  $\vec{x}$

$$S^{(1)}(\vec{x}) = \vec{w}^{(1)} \cdot \vec{x}$$

$$S^{(K)}(\vec{x}) = \vec{w}^{(K)} \cdot \vec{x}$$

} Prediction  $\Rightarrow \hat{y} = \underset{C \in \{1, \dots, K\}}{\text{argmax}} S^{(C)}(\vec{x})$

How to learn parameters?

[Similar structure as binary SVM OP]

$$\min \text{Norm}^2(\vec{w}) + C \sum_i \text{slack}_i$$

s.t. correct classification  $\geq$  Margin - slack<sub>i</sub>

$$\Leftrightarrow \min_{\substack{\vec{w}^{(1)} \\ \vec{w}^{(2)} \\ b^{(1)} \\ b^{(2)}}} \frac{1}{2} \sum_{c=1}^K \|\vec{w}^{(c)}\|^2 + C \sum_{i=1}^N \zeta_{ei}$$

$$\text{s.t.} \quad \underbrace{[\vec{w}^{(y_i)} \cdot \vec{x}_i + b^{(y_i)}]}_{\text{Score of correct class } y_i} - \underbrace{[\vec{w}^{(y)} \cdot \vec{x}_i + b^{(y)}]}_{\text{Score of incorrect class } y} \geq 1 - \zeta_{ei}$$

Score of correct class  $y_i$

Score of incorrect class  $y$

$\forall i, y \neq y_i$   
for all incorrect classes

$$\zeta_{ei} \geq 0$$

Still a QP!

→  $K(d+1)$  vars

→  $N \cdot (K-1) + N$  constraints

Exercise: For  $K=2$ , can we reduce this to a "regular" binary SVM?  
[Yes, with an assumption!]

$$\min_{\substack{\vec{w}^{(1)} \quad \vec{w}^{(2)} \\ b^{(1)} \quad b^{(2)}}} \frac{1}{2} [\vec{w}^{(1)} \cdot \vec{w}^{(1)} + \vec{w}^{(2)} \cdot \vec{w}^{(2)}] + C \sum_{i=1}^N \zeta_{ei}$$

$$\text{s.t.} \quad \begin{cases} \vec{w}^{(1)} \cdot \vec{x}_i + b^{(1)} - \vec{w}^{(2)} \cdot \vec{x}_i - b^{(2)} \geq 1 - \zeta_{ei} & \forall i, y_i = 1 \\ \vec{w}^{(2)} \cdot \vec{x}_i + b^{(2)} - \vec{w}^{(1)} \cdot \vec{x}_i - b^{(1)} \geq 1 - \zeta_{ei} & \forall i, y_i = 2 \end{cases}$$

assume/constrain  $\vec{w}^{(1)} = -\vec{w}^{(2)}$   
 $b^{(1)} = -b^{(2)}$

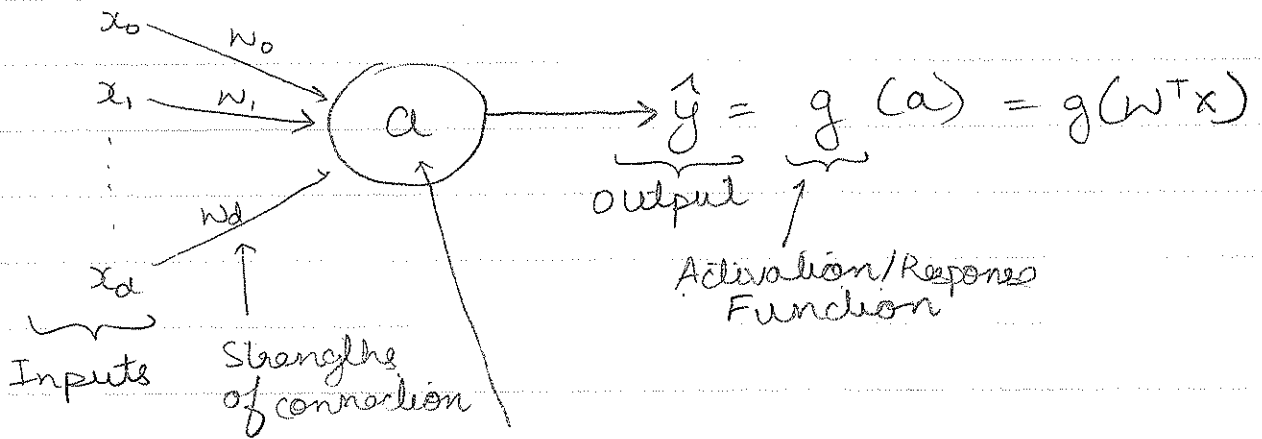
→ get back binary SVM QP (up to constants)

② Where we are going with these classifiers:

Purely Linear Classifiers	Kernel Classifiers	"Natively" Non-linear Classifiers
[Linear in $\vec{x}$ & $\vec{w}$ ]	[Linear in $\phi(\vec{x})$ non-linear in $\vec{x}$ , but linear in $\vec{w}$ ]	[Non-linear in $\vec{x}$ & $\vec{w}$ ]
GNB, (linear) LR SVM	Kernel SVMs	Neural Nets, Decision Trees

$\xrightarrow{\text{Flow of class}}$   
 $\xleftarrow{\hspace{10em}}$

③ Artificial "Neuron"

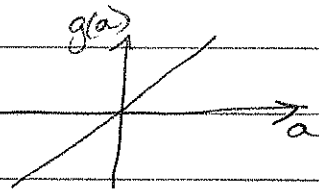


Activation

$$a = \sum_{j=0}^d w_j x_j = \vec{w}^T \vec{x}$$

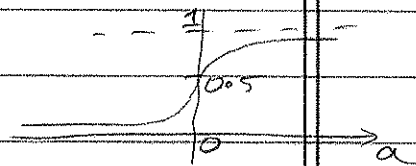
Many different activation functions possible

(a) Linear  $g(a) = a$



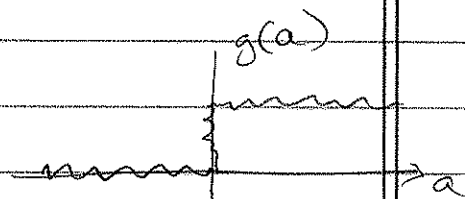
$$\Rightarrow \hat{y} = g(\vec{w}^T \vec{x}) = \vec{w}^T \vec{x} \Leftrightarrow \text{Linear Regression}$$

(b) Sigmoid / Logistic  $g(a) = \frac{1}{1 + e^{-a}}$



$$\Rightarrow \hat{y} = \frac{1}{1 + e^{-w^T x}} \quad \left\{ \begin{array}{l} \text{Logistic Regression} \\ P(Y=1 | \vec{x}, \vec{w}) \end{array} \right.$$

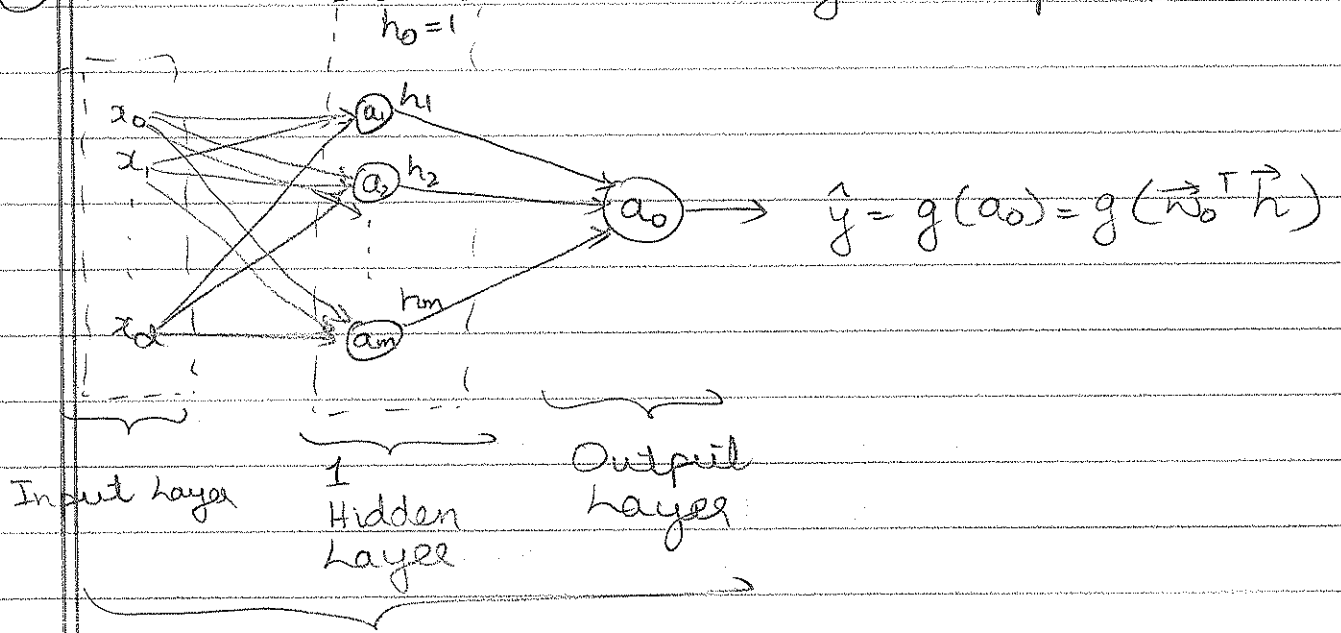
(c) Step  $g(a) = \begin{cases} 1 & a \geq 0 \\ 0 & \text{else} \end{cases}$



$$\hat{y} = \begin{cases} 1 & \text{if } w^T x \geq 0 \\ 0 & \text{if } w^T x < 0 \end{cases} \quad \text{Binary Classifier}$$

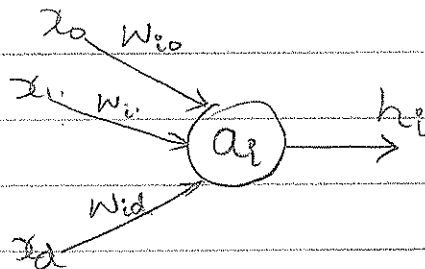
called "Perceptron"  
in this context.

### ④ Neural Network: Multi-layer Perceptrons



3-layer NN

"Zoom-in"



→ Hidden activation 
$$a_i = w_{i0} \cdot x_0 + w_{i1} x_1 + \dots + w_{id} x_d = \vec{w}_i^T \vec{x}$$

→ Hidden response 
$$h_i = g(a_i) = g(\vec{w}_i^T \vec{x})$$

→ Stack all hidden responses into a vector

$$\vec{h} = \begin{bmatrix} h_0 = 1 \\ h_1 \\ \vdots \\ h_m \end{bmatrix}$$

→ Output Activation  $a_0 = w_{00}h_0 + w_{01}h_1 + \dots + w_{0m}h_m$   
 $= \vec{w}_0^T \vec{h}$

→ Output  $\hat{y} = g(\vec{w}_0^T \vec{h})$



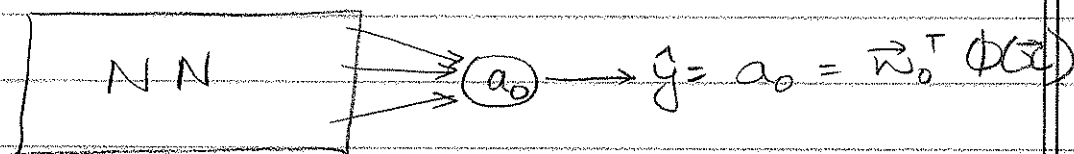
### ⑤ Output Layer & Response Functions

Think about a NN (till the penultimate layer) as a "feature extractor" converting  $\vec{x} \rightarrow \phi(\vec{x})$  but unlike kernel SVMs  $\phi(\vec{x})$  is really  $\phi(\vec{x}; \vec{w})$

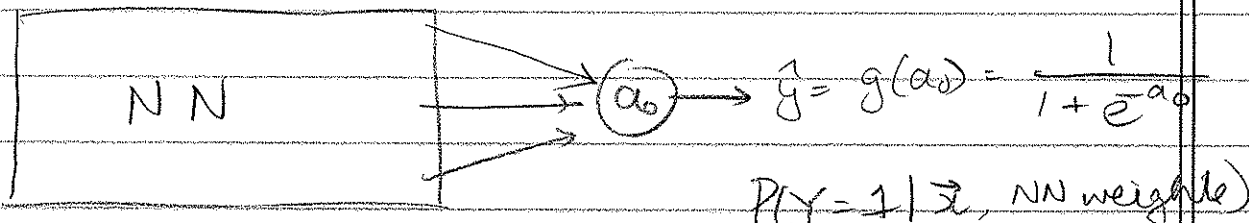
↖ we will learn  $\vec{w}$  to extract different features

The last layer in a NN will depend on the task

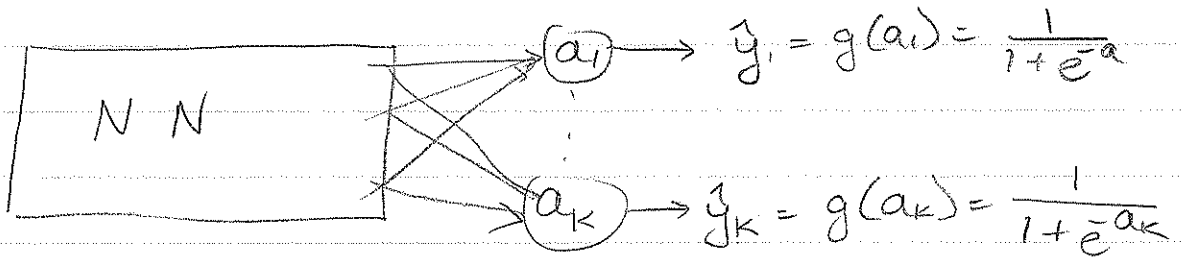
(a) Regression: 1 Output Unit; Linear Output Response



(b) Binary Classification: 1 Output Unit; Sigmoid Output Response



(c) Multi-label Classification [Note: not multi-class]  
K- Output Units ; each Sigmoid



→ An input can be assigned to multiple classes!

(d) Multi-class classification  
K- Output Units : "soft-max" response

