

Machine Learning I: Supervised Methods

B. Keith Jenkins

Announcements

- Slido event code: 1583423
- Graded midterms returned at end of lecture today
- Homework 6 is due Friday
- Project assignment is coming

Reading

- Bishop 5.0-5.3
 - Feedforward ANN

Today's lecture

- Artificial neural networks (part 1)
 - Single neuron unit
 - Activation functions
 - Learning algorithms
 - Single layer of neuron units
- Multilayer ANNs (feedforward)
 - Capabilities
 - Interpretations of multiple layers for pattern classification
 - Learning algorithms

- Annotations to figures on pp. 11-12 revised for clarity post lecture.

ARTIFICIAL NEURAL NETWORKS (ANN)

Ref.'s: 1. Bishop, Ch.5

2. I. Goodfellow, et al.

3. C.M. Bishop, H. Bishop, Deep Learning: Foundations and Concepts (Springer, 2024).

What is A.N.N.?

- Sequence of linear combinations of inputs, followed by thresholding

- Network of neuron units and interconnections

Why A.N.N.?

- Parallel, distributed representation of algorithms

- Mammalian brains are very good at pattern recognition

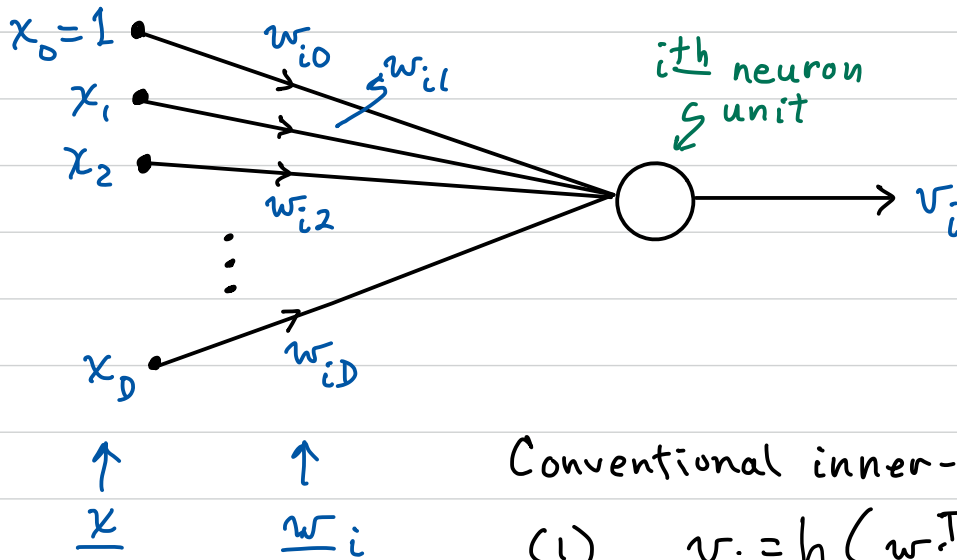
- Recent successes in applying A.N.N. to problems

- Can be proven to be universal function approximators (with a few caveats)

- Annotations to figures on pp. 11-12 revised for clarity post lecture.

Consider one neuron unit and its input connections

→ augmented notation



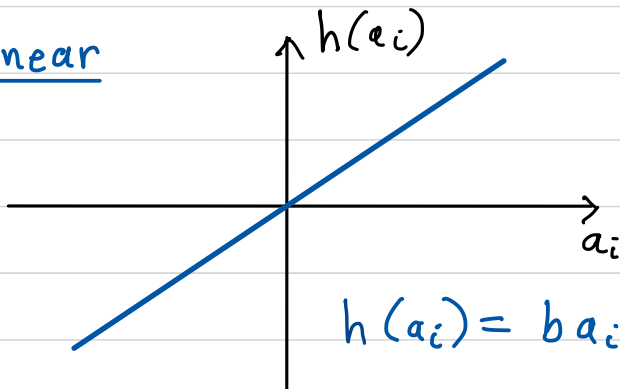
Conventional inner-product neuron unit:

$$(1) \quad v_i = h(\underline{w}_i^T \underline{x}) = h(a_i)$$

↑ activation function (membrane) potential or activation

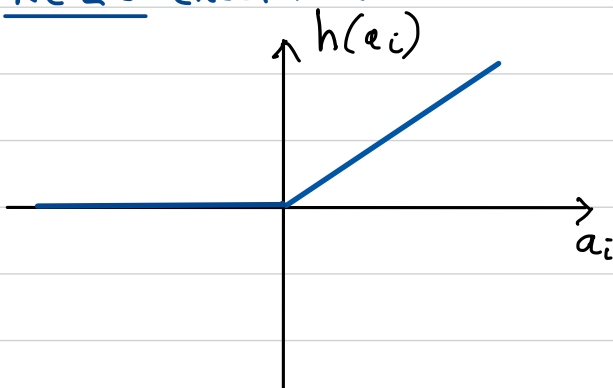
Common Activation Functions

Linear

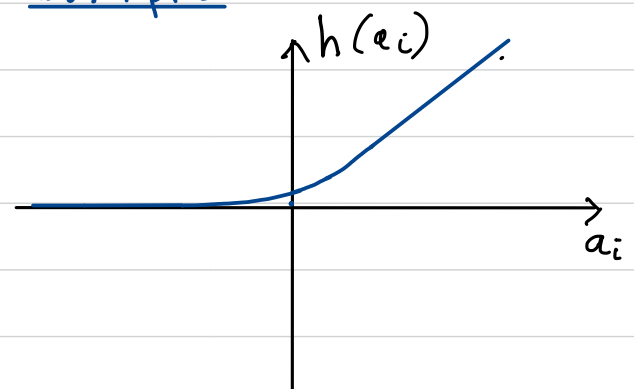


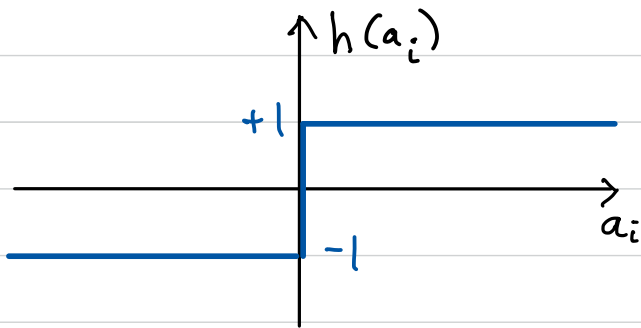
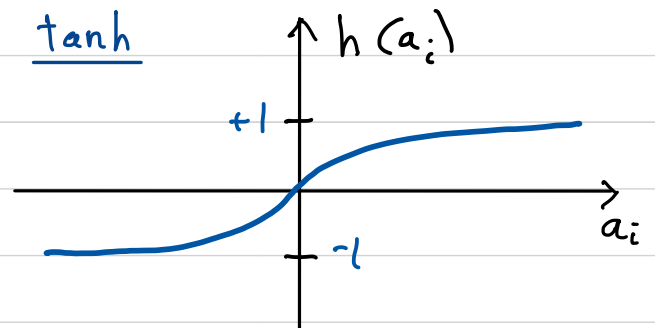
Note: $v_i = h(\underline{w}_i^{(0)T} \underline{x}^{(0)} - \theta_i)$
 $\theta_i \triangleq \text{threshold} = -w_{i0}$.

ReLU (Rectified linear unit)

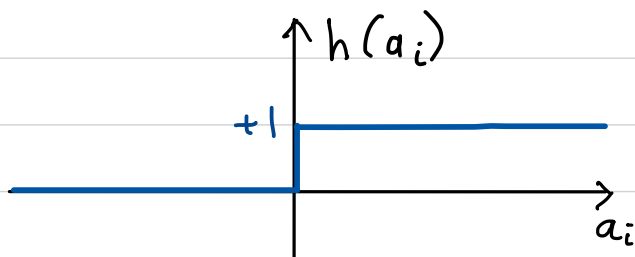
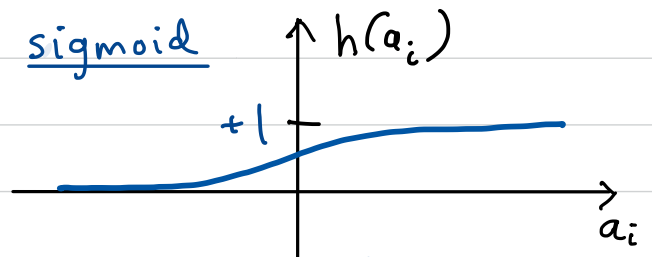


Softplus



Hard thresholdSoft thresholdtanh

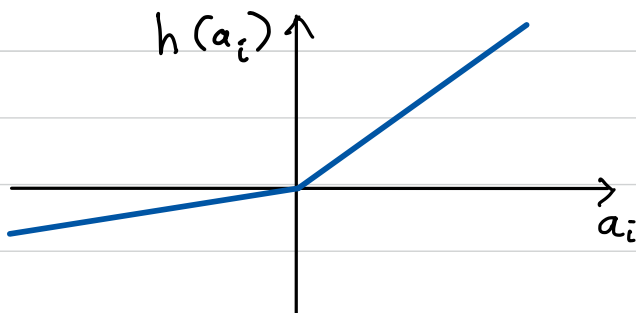
$$h(a_i) = \tanh(a_i)$$

sigmoid

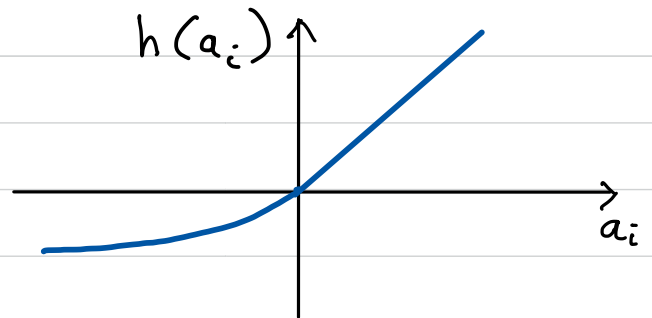
$$h(a_i) = \frac{1}{1 + e^{-a_i}}$$

Leaky ReLU

(Leaky rectified linear unit)

ELU

(Exponential LU)



⑤

$$h(a_i) = \begin{cases} a_i, & \text{if } a_i > 0 \\ b a_i, & \text{otherwise} \end{cases}$$

typically $0 < b < 1$.

$$h(a_i) = \begin{cases} a_i, & \text{if } a_i \geq 0 \\ b(e^{a_i} - 1), & \text{otherwise} \end{cases}$$

$$b > 0.$$

Consider 1 neuron unit, h = hard threshold

$$v = h(\underline{w}^T \underline{x}) = \begin{cases} 1, & \underline{w}^T \underline{x} > 0 \\ 0, & \underline{w}^T \underline{x} < 0 \end{cases}$$

What can this neuron unit implement in pattern classification?

→ a 2-class linear classifier: $g(\underline{x}) = \underline{w}^T \underline{x}$

output representation:

$$\begin{aligned} v = 1 &\Rightarrow \underline{x} \in S_1 \\ v = 0 &\Rightarrow \underline{x} \in S_2. \end{aligned}$$

ADALINE:
adaptive linear
element

w 's are the adapting mechanism, and provide storage for "long-term memory".

LEARNING ALGORITHM EXAMPLES (classification)

1 neuron unit

Let:

$$\begin{aligned} \hat{v} &= \text{actual output (classifier prediction)} \in \{0, 1\} \\ v_t &= \text{target output (true class label)} \in \{0, 1\} \end{aligned}$$

Perceptron Learning

Original perceptron:

$$n = (i \bmod N) + 1$$

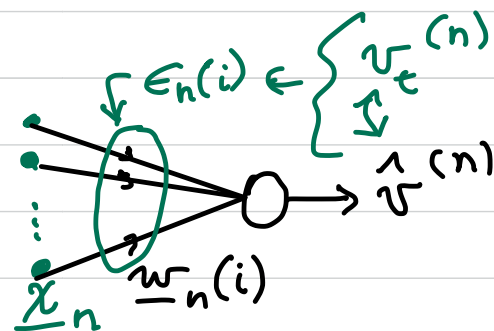
If \underline{x}_n is misclassified:

$$\underline{w}(i+1) = \underline{w}(i) + \eta(i) z_n \underline{x}_n$$

Otherwise

$$\underline{w}(i+1) = \underline{w}(i)$$

Next i



Now:

$$(2) \quad \underline{w}(i+1) = \underline{w}(i) + \eta(i) \underbrace{\left[v_t^{(n)} - \hat{v}^{(n)} \right]}_{E_n(i)} \underline{x}_n$$

$$\Delta \underline{w} = \underline{w}(i+1) - \underline{w}(i) = \eta \dots$$

Target v_t	NN o/p \hat{v}	z_n	$v_t^{(n)} - \hat{v}^{(n)}$
-----------------	---------------------	-------	-----------------------------

miscl. {
correct (o.w.) {

0 (s_2)	1 (s_1)	-1	-1
1 (s_1)	0 (s_2)	+1	+1
$v_t = \hat{v}$			0

Many other linear-classifier learning algorithms can also be posed in A.N.N. framework:

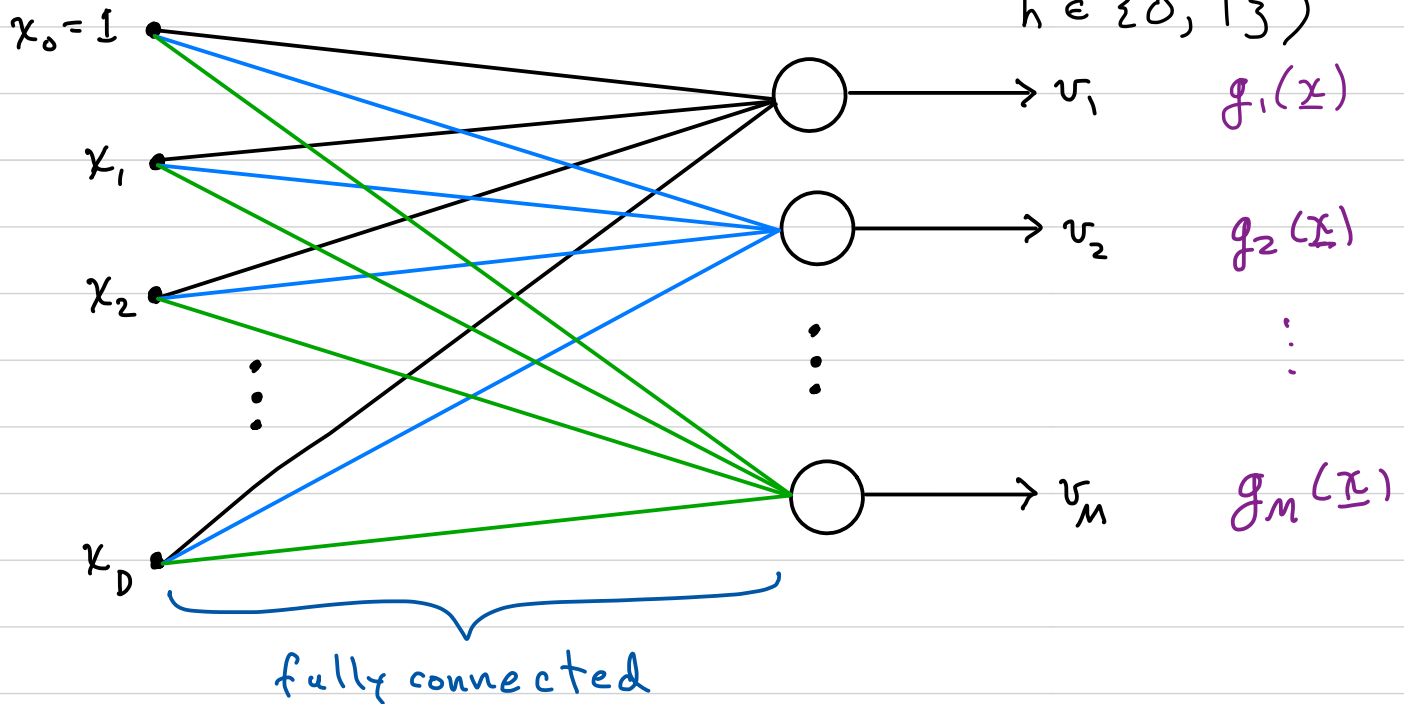
e.g., LMS for regression:

$$\underline{w}(i+1) = \underline{w}(i) - \eta(i) \left(\underbrace{y_n}_{v_t^{(n)}} - \underbrace{\underline{w}^T(i) \underline{x}_n}_{\hat{v}^{(n)}(i)} \right) \underline{x}_n, \quad n = (i \bmod N) + 1$$

Use $h(\cdot) = \text{linear} \rightarrow$

$$\underline{w}(i+1) = \underline{w}(i) - \eta(i) \underbrace{\left(v_t^{(n)} - \hat{v}^{(n)}(i) \right)}_{E_n(i)} \underline{x}_n$$

$M > 1$ neuron units in a single layer ($h = \text{hard threshold}$, $h \in \{0, 1\}$)



What can this implement in pattern classification?

$$\begin{aligned} v_1 &= h(a_1) = h(\underline{w}_1^T \underline{x}) \\ v_2 &= h(a_2) = h(\underline{w}_2^T \underline{x}) \\ &\vdots \\ v_M &= h(a_M) = h(\underline{w}_M^T \underline{x}) \end{aligned}$$

→ M 2-class classifiers

⑤

or → Single OvR M -class classifier, using "one-hot" representation:

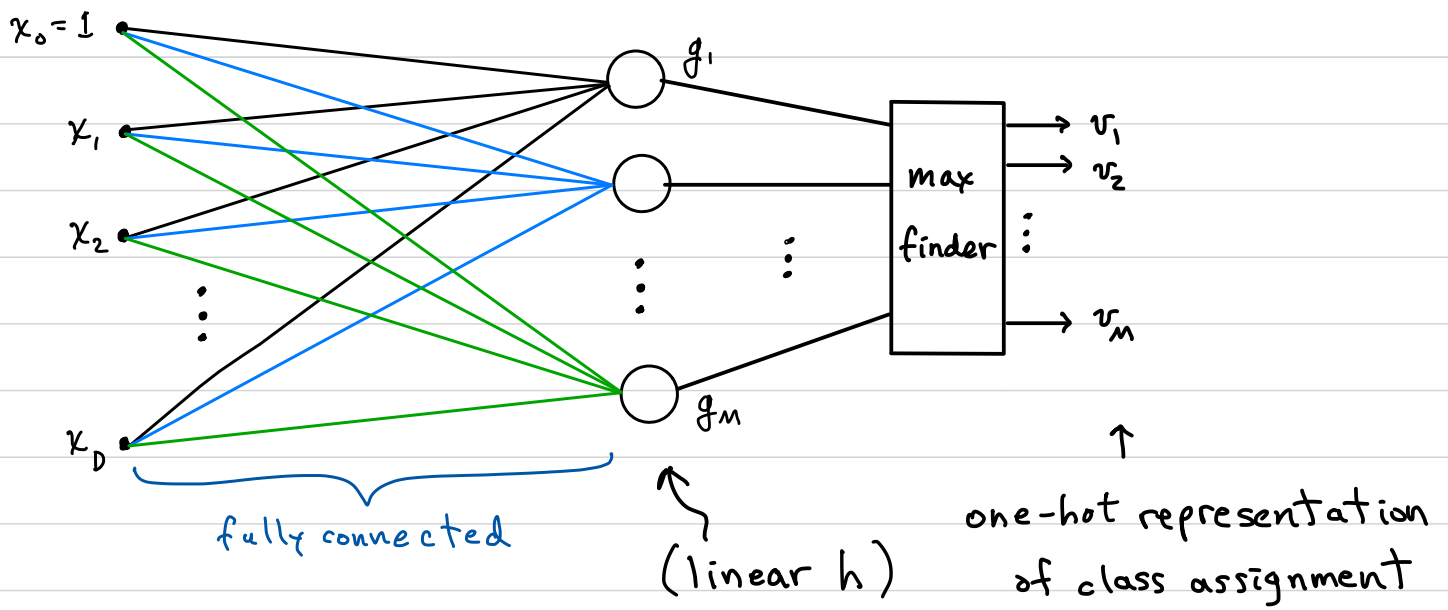
→ Use hard thresholding $h \in \{0, 1\}$.

$$\rightarrow v_i = 1, v_j = 0 \quad \forall j \neq i \Rightarrow \underline{x} \in S_i.$$

Decision rule

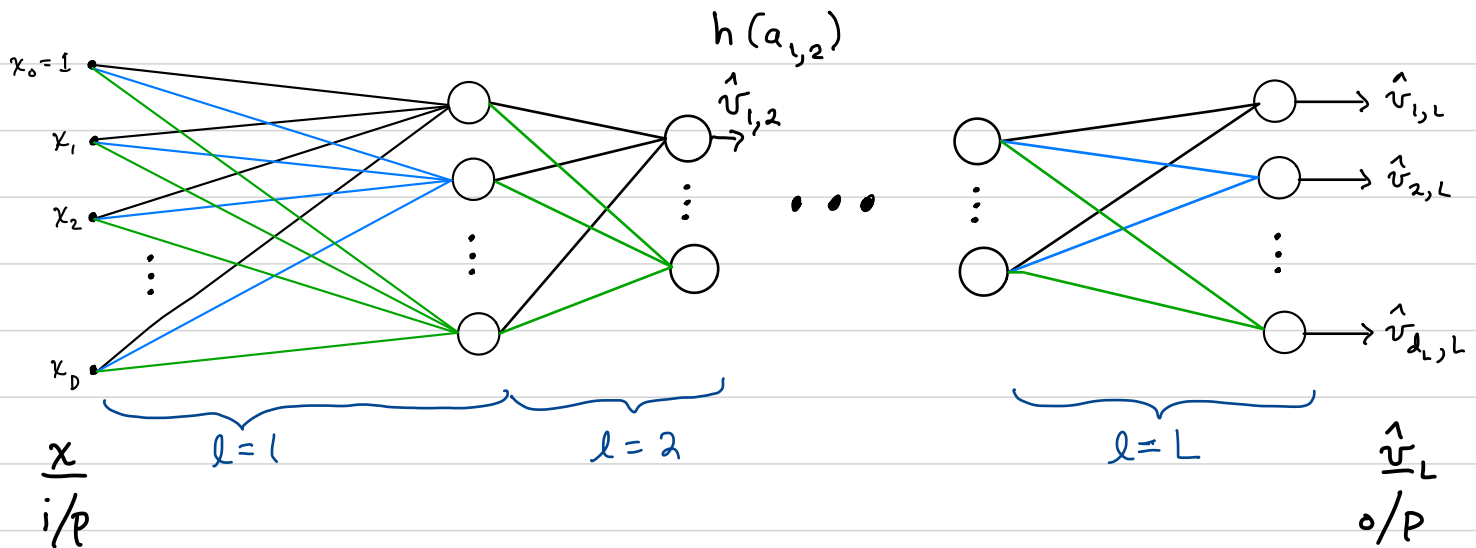
MVM?

→ Can do it with additional layer(s):



MULTILAYER A.N.N.'s [Bishop Ch.5]

Consider feedforward networks

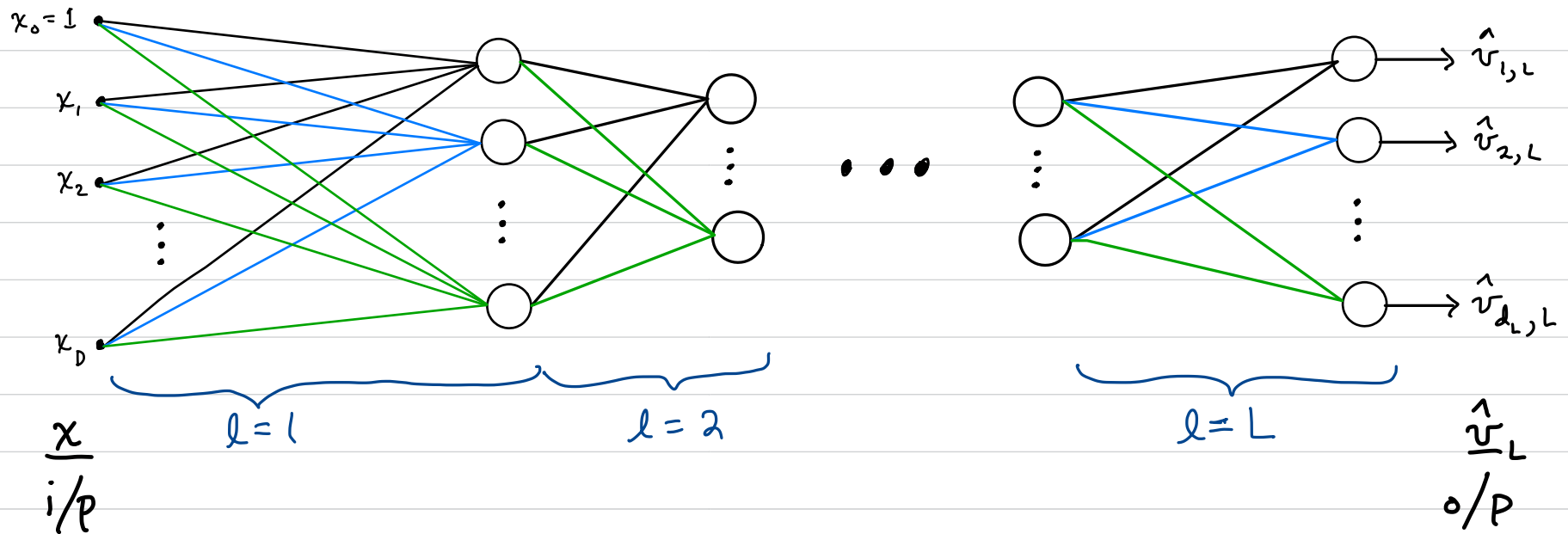


$h(a) = h(a_{i,l})$ can be a nonlinear fcn.

What can a multilayer feedforward A.N.N. compute in pattern classification? (for nonlinear h)

1. i/p is feature space; A.N.N. is a nonlinear classifier.
2. i/p is original feature space; layers $l=1$ to $l=L-1$ perform nonlinear mapping to new feature space. Layer L is a linear classifier.
3. Combination of 1 & 2. (n.l. mapping \rightarrow n.l. classifier).
4. Inputs are pattern space (inputs before feature extraction). First F layers perform feature extraction. Subsequent $(L-F)$ layers perform classification.

Multilayer Feedforward ANN's - Learning algorithms



The number of units in layer l is d_l .

Criterion function : MSE

$$J = \sum_{n=1}^N \mathcal{E}^{(n)} = \sum_{n=1}^N \frac{1}{2} \sum_{i=1}^{d_L} \left[v_{i,L}^{(n)} - \hat{v}_{i,L}^{(n)} \right]^2$$

For classification problems, typ. $d_L = C = \# \text{ classes}$.

For regression " , $d_L = \# \text{ output values to predict for each input } \underline{x}$
(often $d_L = 1$ for regression)

Minimize J by gradient descent.

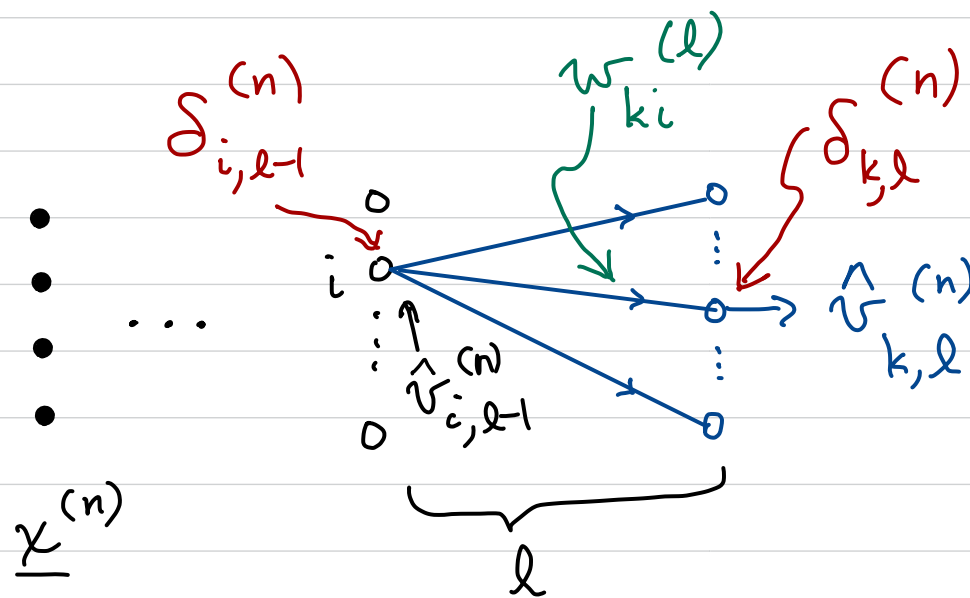
Choose activation functions ($h(a)$).

Want differentiable $h(a)$ (almost everywhere, at least). \leftarrow Assume we have this

Take $\nabla_{\underline{w}} J(\underline{w})$ [Done in Bishop]

\uparrow
all weights in the ANN.

Results:



$$\delta_{i,L}^{(n)} = h'(a_{i,L}^{(n)}) (\hat{v}_{i,L}^{(n)} - \hat{v}_{i,L}^{(n)})$$

The diagram shows the output of the final layer $l=L$. A vertical column of dots represents the output vector $\hat{\underline{v}}^{(n)}$. The output of a specific node is $\hat{v}_{i,L}^{(n)}$. The error term $\delta_{i,L}^{(n)}$ is calculated as the derivative of the activation function h' evaluated at $a_{i,L}^{(n)}$ multiplied by the difference between the predicted output $\hat{v}_{i,L}^{(n)}$ and the target output $\hat{v}_{i,L}^{(n)}$. The output of another node is $\hat{v}_{d_L,L}^{(n)}$.

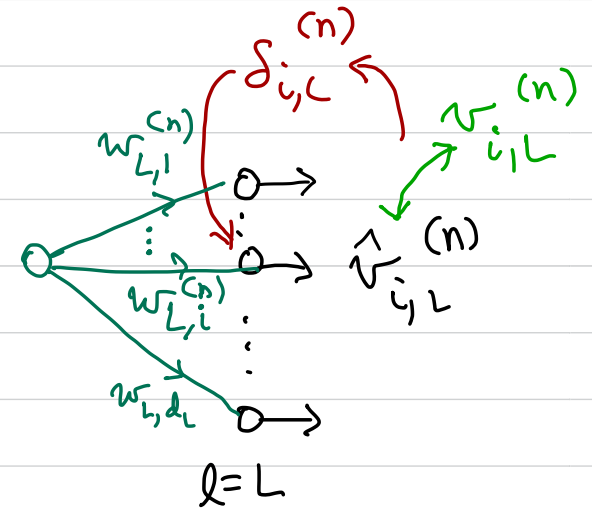
→ Backward error propagation

$$(1) \quad \Delta w_{ki}^{(l)} = \eta \delta_{k,l}^{(n)} \hat{v}_{i,l-1}^{(n)}$$

$$\begin{cases} \Delta w_{ki}^{(l)} = w_{ki}^{(l)}(j+1) - w_{ki}^{(l)}(j) \\ j = \text{iteration index}; \quad \eta = \eta(j) \end{cases}$$

$$(2) \quad \delta_{i,L}^{(n)} = [v_{i,L}^{(n)} - \hat{v}_{i,L}^{(n)}] h'(a_{i,L}^{(n)})$$

$$(3) \quad \delta_{i,l-1}^{(n)} = h'(a_{i,l-1}^{(n)}) \sum_k \delta_{k,l}^{(n)} w_{ki}^{(l)}$$



(i) Forward pass: calculate $\hat{v}_L^{(n)}$ for input \underline{x}_n .

(ii) Backward pass: calculate error terms $\delta_{k,l}^{(n)}$ and weight updates $\Delta w_{ki}^{(l)}$