

# Machine Learning I: Supervised Methods

B. Keith Jenkins

## Announcements

---

- Course Project Assignment has been posted on D2L
  - Includes descriptions of 3 datasets
  - Requirements
  - Tips to help you
  - Grading criteria
- There will be 2 more homeworks (HW7, HW8)
- No Slido poll questions today

\* - Added post-lecture:

- Equation numbers for all  $w_{pq}^{(l)}$ ;
- Last comment on p. 12.

## Reading

---

- Bishop 9.1
  - K-means clustering (we will use this with Radial Basis Function ANNs)

## Today's lecture

---

- Multiple layer feedforward ANNs (part 2)
  - Feedforward ANNs as universal function approximators
    - Proof by construction ( $D=1$ )
  - Summary and comments

## Feedforward ANN's: fundamental capabilities and limitations

- Classification
- Regression
- Function approximation.

ANN, for any input  $\underline{x}$ , will give outputs:  $\hat{y}_i(\underline{x})$ ,  $i=1, \dots, d_L$ .

> Regression: outputs  $\hat{y}_1(\underline{x}), \hat{y}_2(\underline{x}), \dots, \hat{y}_{d_L}(\underline{x})$ .

Like "approximating" some unknown function  $y_i(\underline{x})$

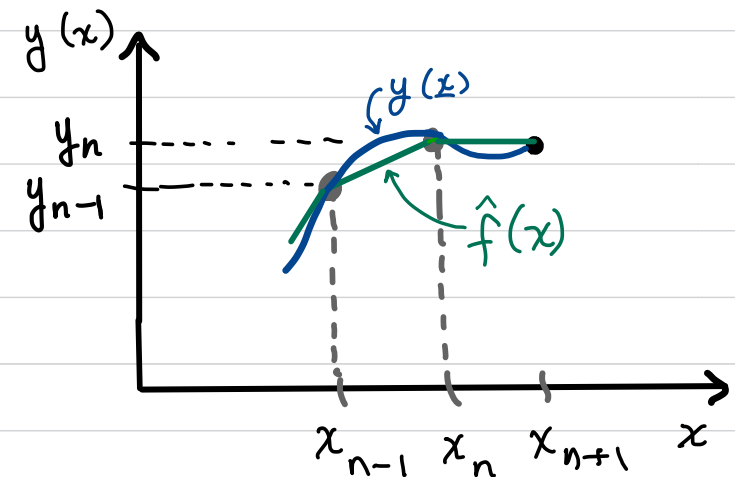
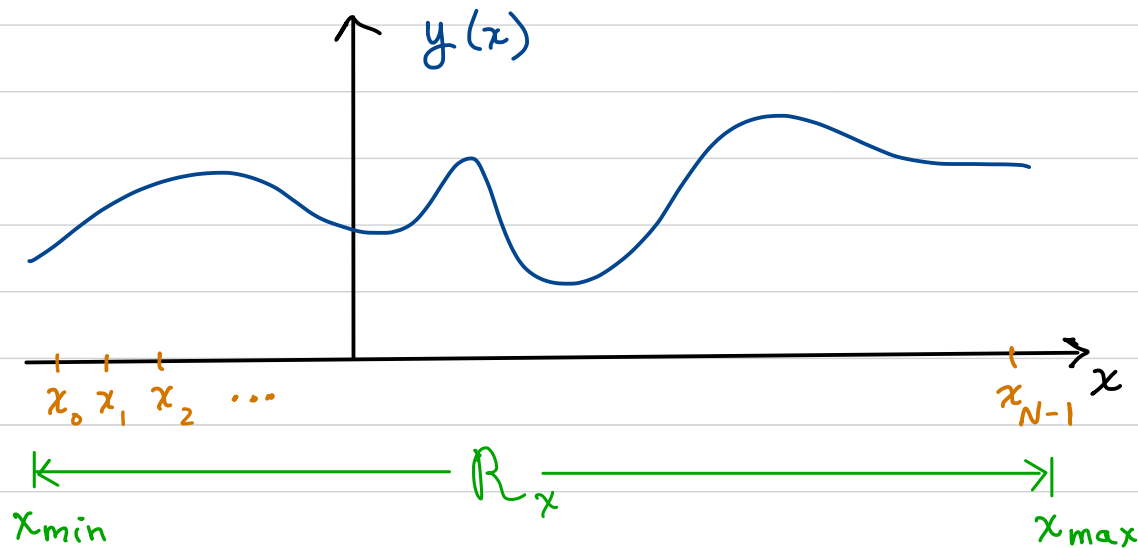
> Classification: (i) approx. binary-valued (unknown) fcn's.  $y_i(\underline{x})$   
 (ii) approx. discriminant fcn.  $g_k(\underline{x})$  or  $g_{kj}(\underline{x})$ .

For ML (reg. & class'n) the fcn. we are approximating is unknown.

Question: how general of a function can a multilayer feedforward ANN approximate, to (ideally) any degree of accuracy?

# Theoretical example of universal function approximation using ReLU

Problem: approximate any function  $y(x)$  using a 2-layer ANN, over the region  $R_x$ .  
Assume a 1D problem, ReLU for the hidden layer.



Use a piecewise linear approximation.

Choose data points (grid points)  $(x_i, y_i)$ ,  $y_i = y(x_i)$ ,  $i = 0, 1, 2, \dots, N-1$ .

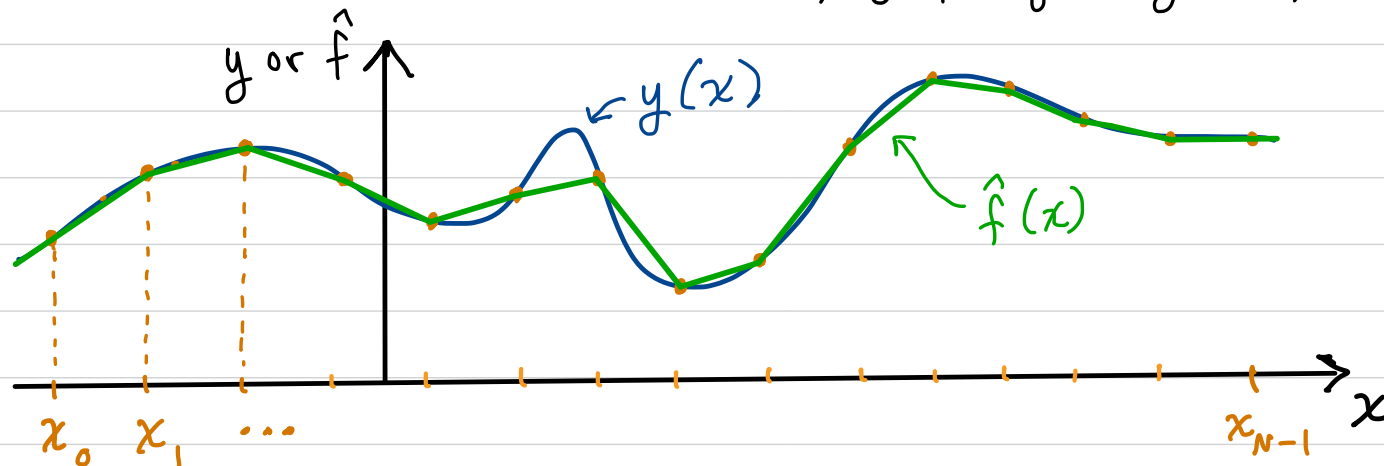
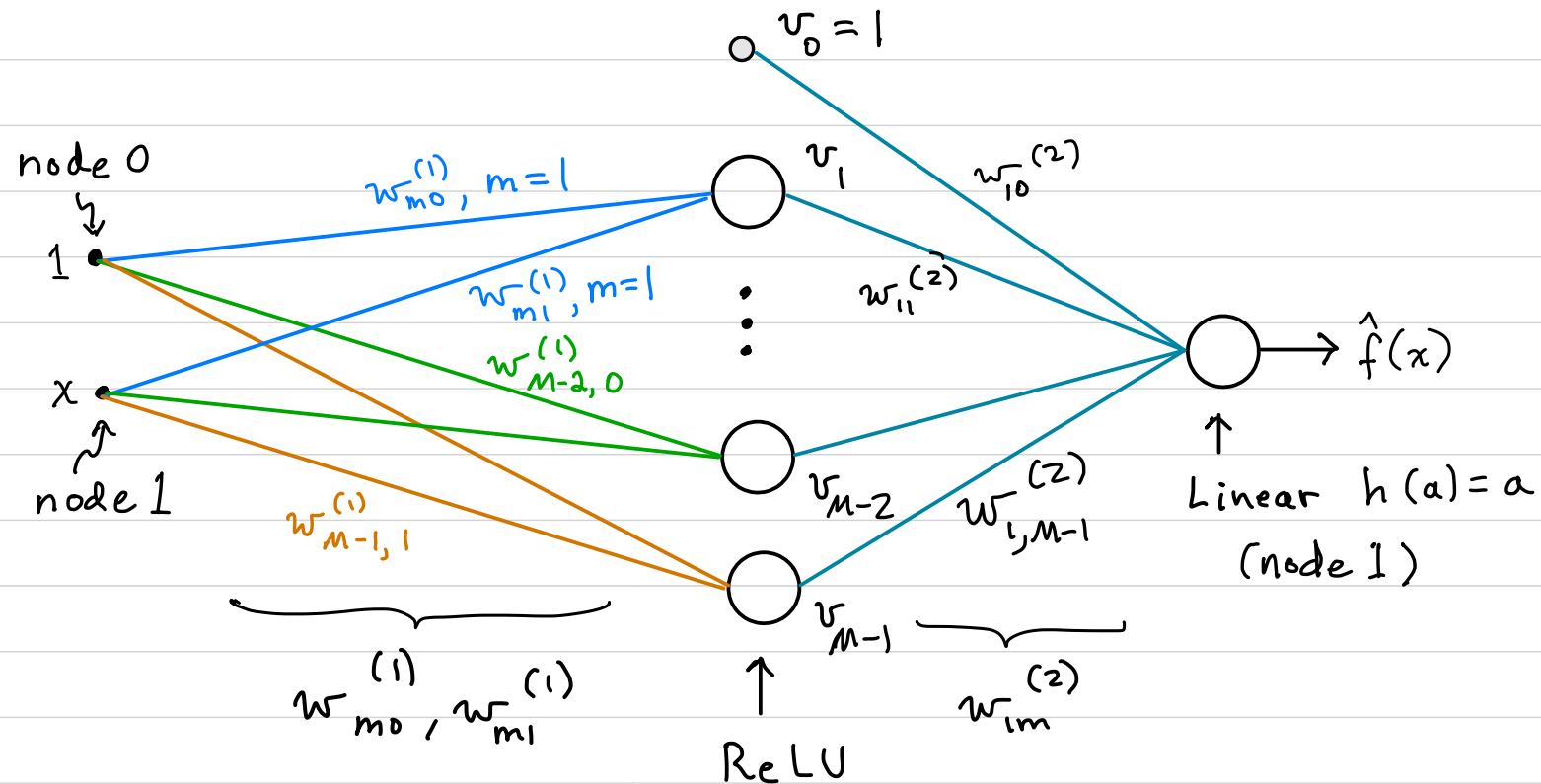


Fig. 1



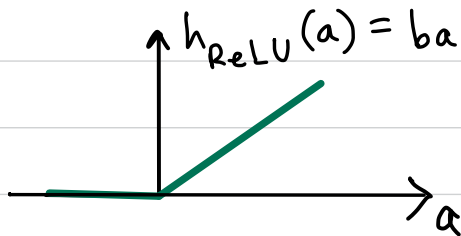
$\Rightarrow$  hidden unit outputs  $v_m(x_n)$ ,  $m = 1, \dots, M-1$ ;  $M-1 = \#$  hidden units.  
Hidden bias node output  $v_0 = 1$ .

Let  $M = N = \#$  "data points" (or "grid points")

Goal: find weights that will give  $\hat{f}(x)$  as a piecewise-linear approximation to  $y(x)$

2<sup>nd</sup> layer1<sup>st</sup> layer  
activation fcn. of  $v_m$ 

$$\hat{f}(x) = w_{i0}^{(2)} v_0(x) + \sum_{m=1}^{M-1} w_{im}^{(2)} v_m(x), \quad v_m(x) = h_{\text{ReLU}}(w_{m1}^{(1)} x + w_{m0}^{(1)})$$



$$v_m(x) = \max \{ 0, b(w_{m1}^{(1)} x + w_{m0}^{(1)}) \}$$

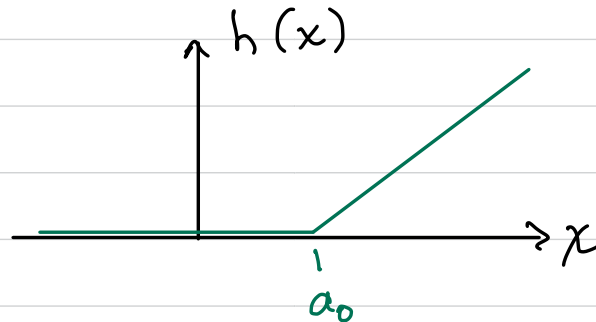
↑ ↑

Let ReLU parameter  $b=1$  ( $\Rightarrow$  slope 1).Let  $w_{m1}^{(1)} = 1 \quad \forall m$ 

$$\left. \begin{array}{l} \text{Let ReLU parameter } b=1 \text{ (}\Rightarrow\text{ slope 1).} \\ \text{Let } w_{m1}^{(1)} = 1 \quad \forall m \end{array} \right\} v_m(x) = \max \{ 0, x + w_{m0}^{(1)} \}$$

$$\hat{f}(x) = w_{i0}^{(2)} v_0(x) + \sum_{m=1}^{M-1} w_{im}^{(2)} \max \{ 0, x + w_{m0}^{(1)} \}$$

We will want:

How to express using  $h_{\text{ReLU}}$ ?

$$\begin{aligned} \rightarrow h(x) &= h_{\text{ReLU}}(x - a_0) \\ &= \max \{ 0, x - a_0 \} \end{aligned}$$

Order the  $N$  data points so that:  $x_0 < x_1 < x_2 < \dots < x_{N-1}$   
in which duplicate data points are omitted.

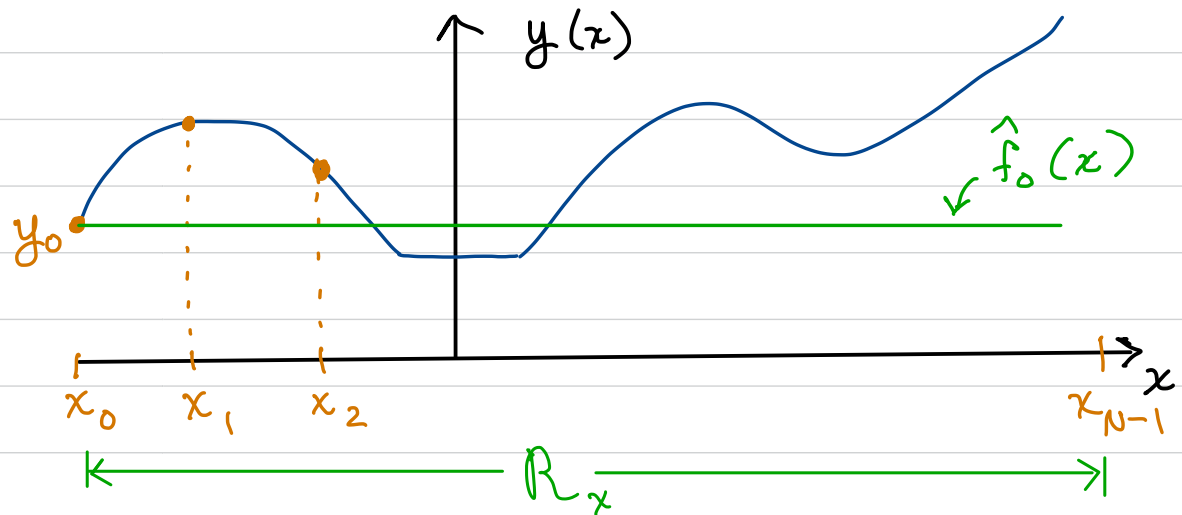
For simplicity, assume  $x_0 = x_{\min}$  and  $x_{N-1} = x_{\max}$  (we have data points at edges of region  $R_x$ ).

Comment: this proof also applies (with minor modifications) if  $x_{\min} < x_0$   
and  $x_{N-1} < x_{\max}$ .

Let  $\hat{f}_{n-1}(x)$  be the function approximation based on data points  $x_0, x_1, \dots, x_{n-1}$ .

$\Rightarrow \hat{f}_0(x)$  uses only data point  $(x_0, y_0)$ .

(2) Let:  $\hat{f}_0(x) = y_0$ .

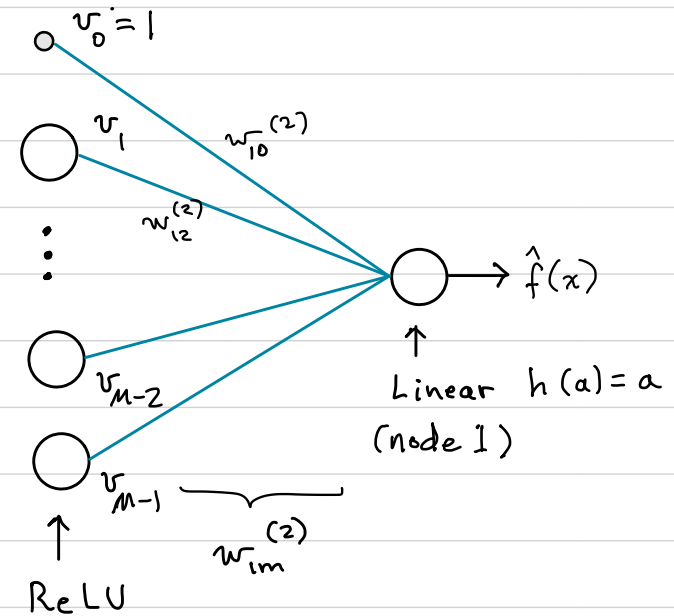


In the ANN (Fig. 1), this can be done by:

$$(3) \left[ w_{10}^{(2)} = y_0 \Rightarrow \hat{f}_0(x) = y_0 \right]$$

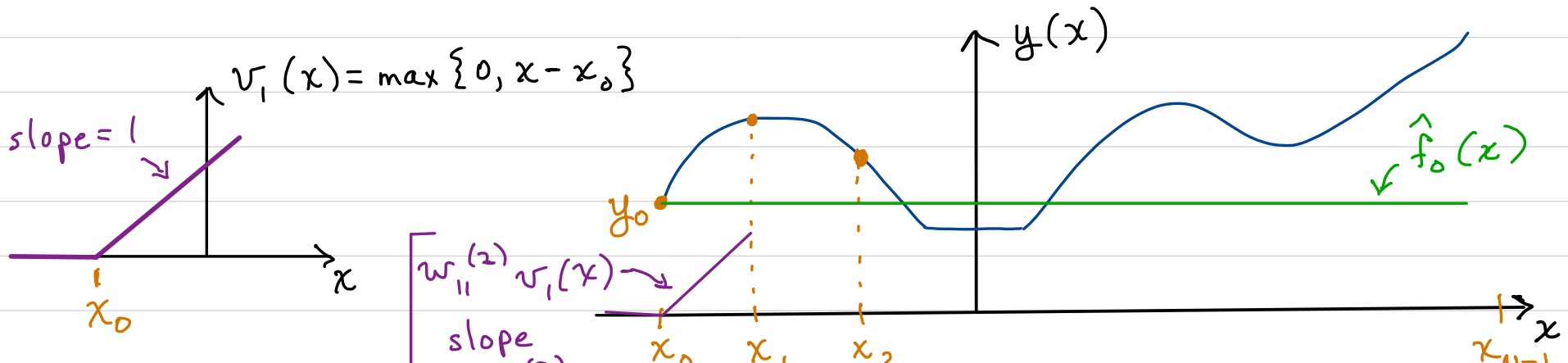
For  $\hat{f}_1(x)$  use data  $(x_0, y_0), (x_1, y_1)$ .  
 $\rightarrow$  Or, use  $\hat{f}_0(x)$  and  $(x_1, y_1)$ .

What is  $\hat{f}_0(x_1)$ ?  $y_0$   
 $\Rightarrow \hat{f}_0(x_1) \neq y_1$  in general



Use  $v_1(x)$  to correct  $\hat{f}_0(x_1)$ , so that  $\hat{f}_1(x_1) = y_1$ .

Place "hinge" of the ReLU for  $v_1$  at  $x_0$ :



$$(1) \Rightarrow \hat{f}_1(x) = \underbrace{w_{10}^{(2)} v_0(x)}_{y_0} + w_{11}^{(2)} \max\{0, x + \underbrace{w_{10}^{(1)}}_{\leftarrow}\}\}$$

$$(4) \quad \left[ \text{Choose } w_{10}^{(1)} = -x_0 \right]$$

$$\hat{f}_1(x) = \hat{f}_0(x) + w_{11}^{(2)} \max\{0, x - x_0\}$$

$$\left[ \hat{f}_1(x) = \hat{f}_0(x) + w_{11}^{(2)} v_1(x) = \hat{f}(x) \text{ if } v_m(x) = 0, m > 1. \right]$$

What slope  $w_{11}^{(2)}$  will make  $\hat{f}_1(x_1) = y_1$ ?

$$\hat{f}_1(x) - \hat{f}_0(x) = w_{11}^{(2)} v_1(x)$$

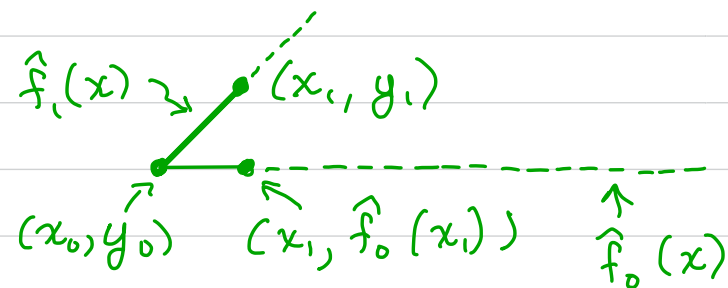
$$\text{slope} = w_{11}^{(2)} = \frac{y_1 - y_0}{x_1 - x_0}$$

$$\hat{f}_1(x_1) - \hat{f}_0(x_1) = w_{11}^{(2)} v_1(x_1)$$

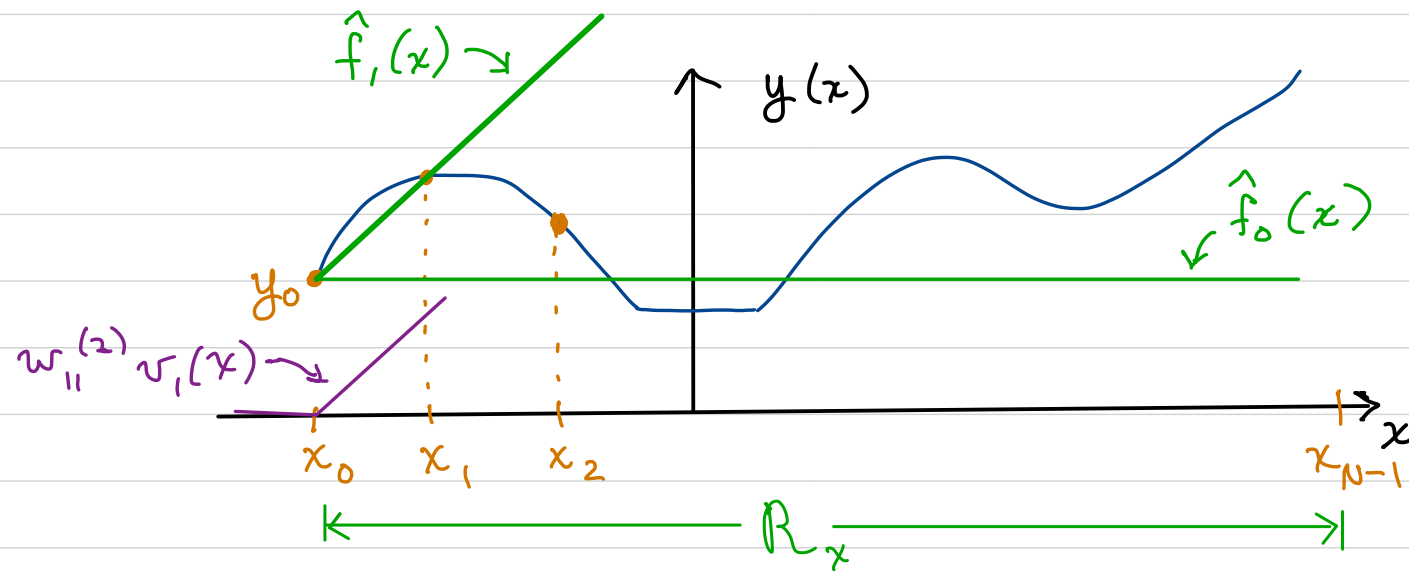
$$y_1 - y_0 = w_{11}^{(2)} (x_1 - x_0) \quad \text{because } x_1 > x_0.$$

$$(5) \quad \left[ w_{11}^{(2)} = \frac{y_1 - \hat{f}_0(x_1)}{x_1 - x_0} \right]$$

$$\Rightarrow \hat{f}_1(x_1) = y_1 \Rightarrow 0 \text{ error at } x_1.$$







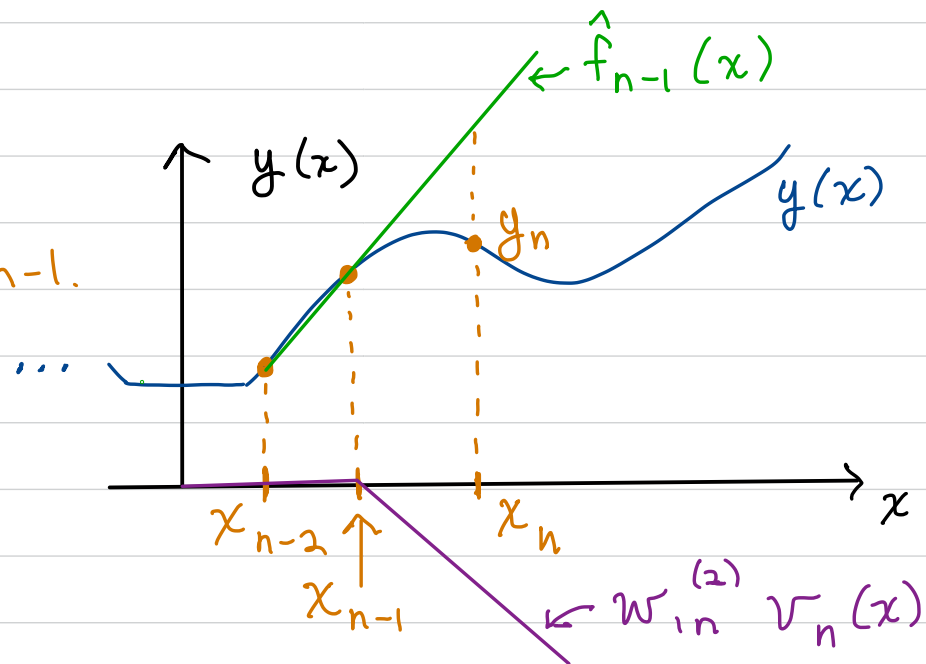
Now consider  $n^{\text{th}}$  segment (for  $n \geq 2$ ):

Given:

$$\hat{f}_{n-1}(x)$$

$$x_n, y_n$$

$$x_m, y_m, m \leq n-1.$$



$$(1) \Rightarrow \hat{f}(x) = \underbrace{w_{10}^{(2)} v_0(x)}_{y_0} + \sum_{m=1}^{M-1} w_{1m}^{(2)} \max \{0, x - \underbrace{w_{m0}^{(1)}}_{\leftarrow}\} \quad \hat{f}_0(x)$$

$$(6) \quad \left[ \text{Choose } w_{m0}^{(1)} = -x_{m-1}, \quad \forall m > 1. \right]$$

$$\begin{aligned} \hat{f}(x) &= \hat{f}_0(x) + \underbrace{\sum_{m=1}^{n-1} w_{1m}^{(2)} \max \{0, x - x_{m-1}\}}_{\triangleq \hat{f}_{n-1}(x)} + \underbrace{\sum_{m=n+1}^{M-1} w_{1m}^{(2)} \max \{0, x - x_{m-1}\}}_{=0 \text{ in region } x_{n-1} \leq x \leq x_n \text{ because } x \leq x_n} \\ &\quad + w_{1n}^{(2)} \max \{0, x - x_{n-1}\} \end{aligned}$$

This is the estimate based on  $(x_i, y_i)$ ,  $i=0, 1, \dots, n$ .

$$\left[ \hat{f}_n(x) = \hat{f}_{n-1}(x) + w_{1n}^{(2)} v_n(x), \quad v_n(x) = \max \{0, x - x_{n-1}\} \right]$$

Assume we have  $\hat{f}_{n-1}(x_{n-1}) = y_{n-1} \Rightarrow 0$  error at  $x_{n-1}$ .

Generally  $\hat{f}_{n-1}(x_n) \neq y_n$ .

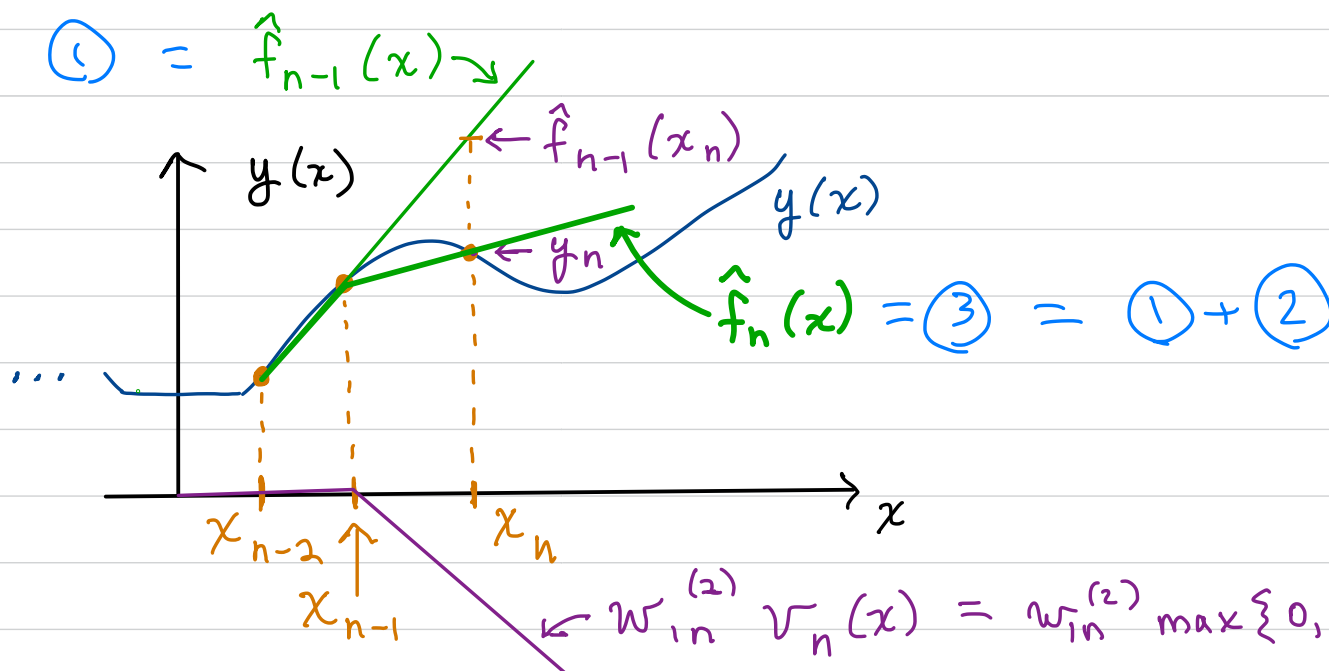
$$\hat{f}_n(x) = \hat{f}_{n-1}(x) + \underbrace{w_{1n}^{(2)} \max \{0, x - x_{n-1}\}}_{\text{Choose to make } \hat{f}_n(x_n) = y_n}$$

$$\hat{f}_n(x_n) = \hat{f}_{n-1}(x_n) + w_{1n}^{(2)}(x_n - x_{n-1}) \quad \text{because } x_n > x_{n-1}.$$

$$\hat{f}_n(x_n) - \hat{f}_{n-1}(x_n) = w_{1n}^{(2)}(x_n - x_{n-1})$$

$$y_n - \hat{f}_{n-1}(x_n) = w_{1n}^{(2)}(x_n - x_{n-1})$$

$$(7) \quad \left[ w_{1n}^{(2)} = \frac{y_n - \hat{f}_{n-1}(x_n)}{x_n - x_{n-1}}, \quad n \geq 2. \right.$$



[ By induction, error can be made 0 at every data point, and  $\hat{f}(x)$  gives linear approximation to  $y(x)$  between neighboring data points.

## Summary

- Feedforward ANN with 1 hidden layer (2 layers of weights).
- Weights  $w_{m1}^{(1)} = 1 \quad \forall m$  (all weights from input  $x$  in first layer)
- Hidden units use ReLU activation functions  $h(a) = \max\{0, a\}$  (unit slope)
- Weights  $w_{m0}^{(1)} = -x_{m-1}$  give offset of each hidden-unit ReLU

$$v_m(x) = \max\{0, x + w_{m0}^{(1)}\} = \max\{0, x - x_{m-1}\}$$

( $x_{m-1}$  is hinge point of  $v_m(x)$  ReLU function.)

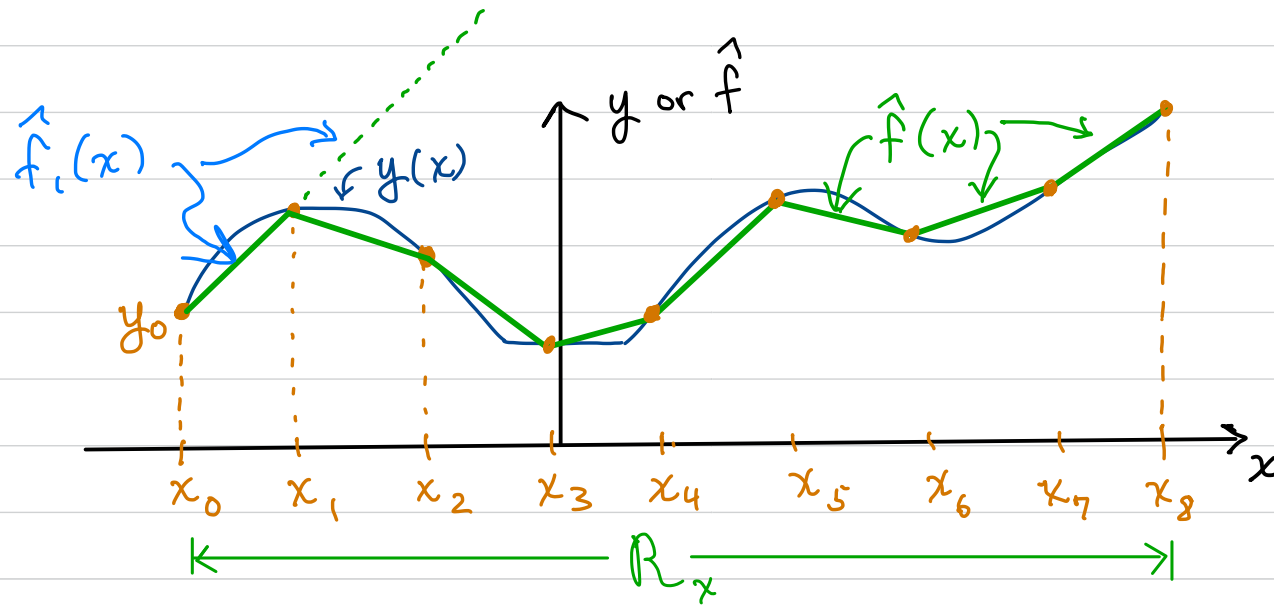
- Second layer weights  $w_{1n}^{(2)}$  give slope of ReLU of  $v_n(x)$  so that

$$\hat{f}_n(x) = \hat{f}_0(x) + \underbrace{\left( \sum_{m=1}^{n-1} w_{1m}^{(2)} v_m(x) \right)}_{\text{Sum of contributions from all } v_m(x), m < n.} + w_{1n}^{(2)} v_n(x)$$

gives error-free estimate at  $x_n$ :  $\hat{f}_n(x_n) = y_n$ .

- \* ◦ All weight values are given in Eqs. (1), (3), (4)-(7).

- Example result ( $N=9$  data points)



## Comments

1.  $\hat{f}(x)$  is a piecewise-linear approximation of  $y(x)$ .

2. By taking  $N (=M)$  sufficiently large, error of approximation

$$\mathcal{E} = \frac{1}{x_{\max} - x_{\min}} \int_{R_x} (\hat{f}(x) - y(x))^2 dx$$

can be made arbitrarily small, for essentially any realistic function  $y(x)$ .

3. Thus, this ANN shows that a (1D) universal function approximator is possible with a 2-layer feedforward ANN.

4. Can be extended to higher dimensional inputs  $D \geq 2$ .  
( $D > 2$  most easily proven using Fourier-series approach instead of piecewise-linear).