

Nonlinear Function Approximation: Neural Networks

Qiang Liu
UT Austin

Approaches to Nonlinear Regression:

- Fixed basis function:

$$f(\underline{x}; w) = \sum_i w_i \phi_i(x),$$

e.g., polynomial regression: $f(x; w) = w_0 + w_1 x + w_2 x^2 + \dots$

- Adaptive basis functions:

- Kernel method:

$$\underline{f(x; w)} = \sum_{i \in \text{data}} w_i k(x, x_i).$$

- Neural networks:

$$\underline{f(x; w, v)} = \sum_i w_i \phi(x; v_i),$$

$$\underline{\phi_i(x)} = \underline{\phi(x, v_i)}$$

$$\min_{w, v} E_D[(y - f(x; w, v))^2]$$

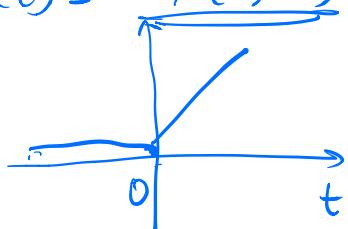
$$\boxed{\underline{\phi(x)} = \sigma\left(\sum_{e=1}^d v_e x_e + v_0\right)}$$

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$\sigma(\cdot)$: activation function

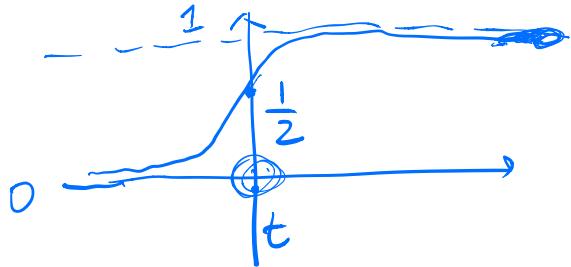
$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_d \end{bmatrix}$$

$\sigma(t) = \max(0, t)$: Rectified Linear Unit (ReLU)



$$\sigma(t) = \begin{cases} 0 & t \leq 0 \\ t & t > 0 \end{cases}$$

$$\sigma(t) = \frac{\exp(t)}{1 + \exp(t)} = \frac{1}{1 + \exp(-t)}$$



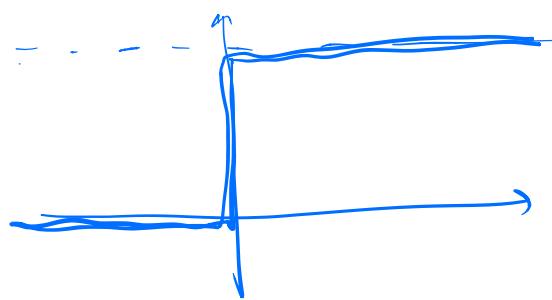
Sigmoid function.

$\phi(x)$: Neuron.

$\sigma(\cdot)$: activation

v : weights

$$\sigma(t) = \mathbb{I}(x > 0) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$$



$$f(x; w, v) = \sum_{i=1}^m w_i \sigma\left(\sum_{l=1}^d v_{il} x_l + v_{i0}\right)$$

$$\underline{v}_i = [v_{i0}, v_{i1}, \dots, v_{id}]$$

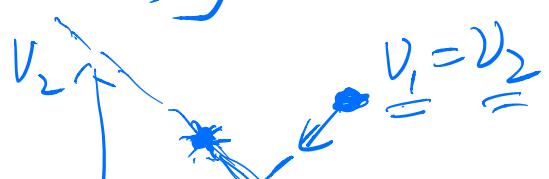
$$\underline{v} = [v_1, \dots, v_m]$$

* Learning Neural Nets. is Non-Convex Opt.

$$f(x, v) = \sigma(x - v_1) + \sigma(x - v_2)$$

$$L(v) = E_D[(y - \sigma(x - v_1) - \sigma(x - v_2))^2]$$

Assume $[v_1^*, v_2^*]$ is op.

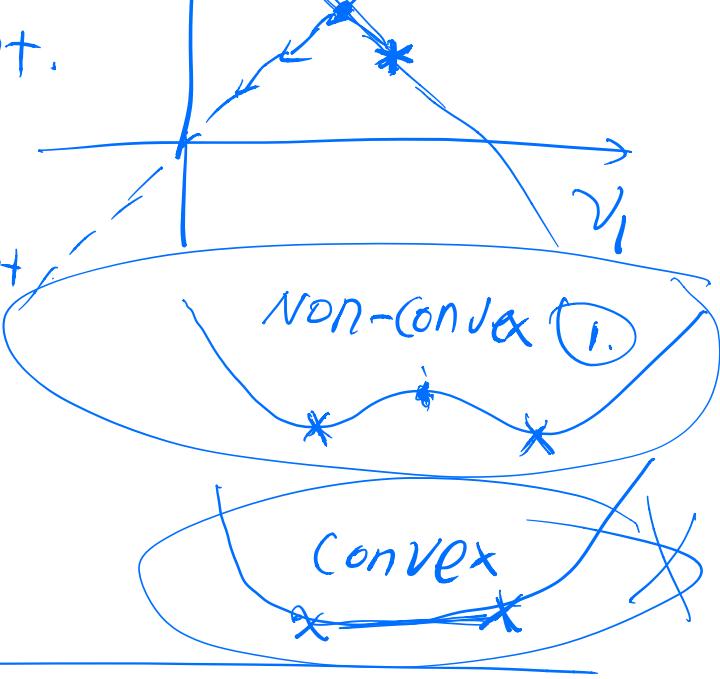


then $[v_2^*, v_1^*]$ is also opt.

If $L(v)$ is convex

then $\underline{[v^*, v^*]}$ is also opt

$$\bar{v}^* = \frac{v_1^* + v_2^*}{2}$$



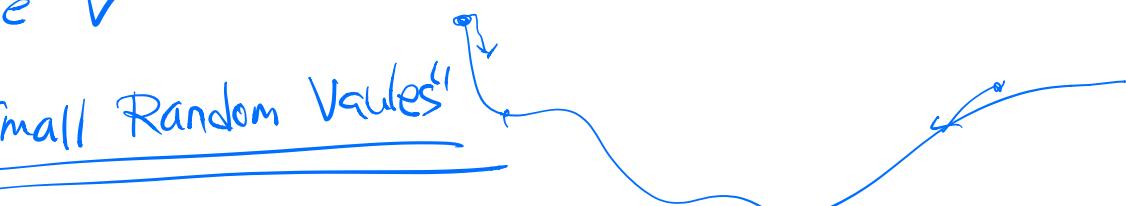
Gradient descent

$$\min_{w, v} L(w, v)$$

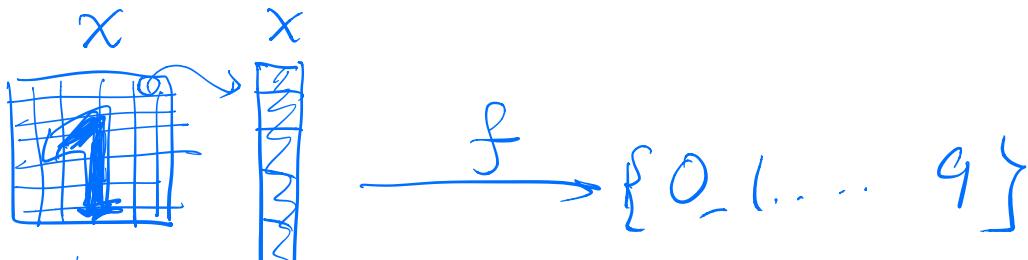
$$\begin{cases} w \leftarrow w - \epsilon \nabla_w L(w, v) \\ v \leftarrow v - \epsilon \nabla_v L(w, v) \end{cases}$$

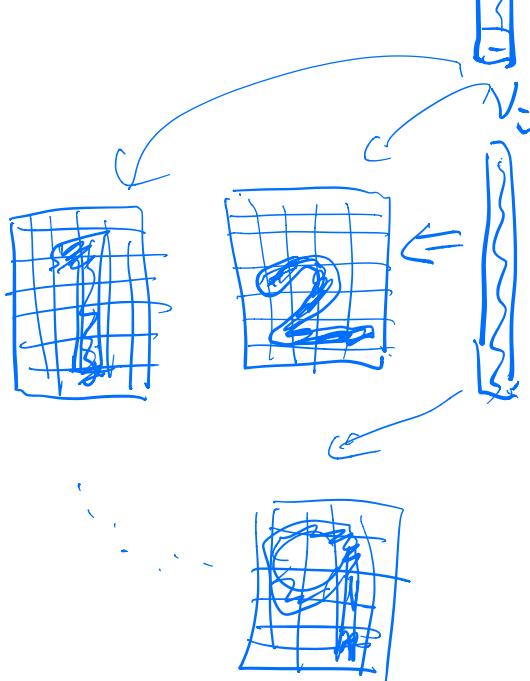
Initialize v

to be "small Random Values"



Digit classification

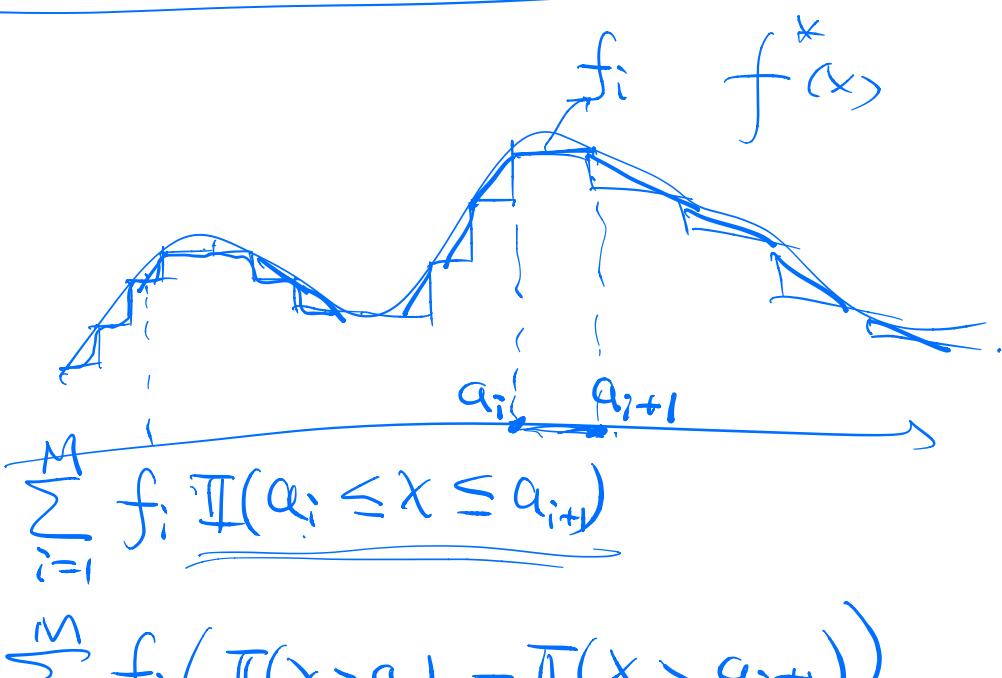




$$\sum_{i=1}^m w_i \mathbb{I}(v_i^T x - v_0)$$

$$\sum_{i=1}^m w_i \mathbb{I}(v_i^T x \geq v_0)$$

If m is large enough, $f(x, \underline{w}, \underline{v})$
can approximate any continuous function



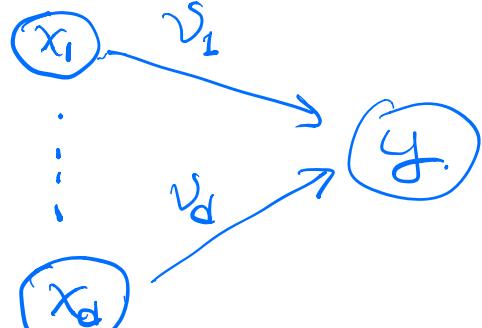
$$\sum_{i=1}^M f_i \mathbb{I}(a_i \leq x \leq a_{i+1})$$

$$\sum_{i=1}^M f_i (\mathbb{I}(x > a_i) - \mathbb{I}(x > a_{i+1}))$$

$$= \left[\sum_{i=1}^M (f_i - f_{i-1}) \mathbb{I}(x \geq a_i) \right]$$

Graphical Representation of NN

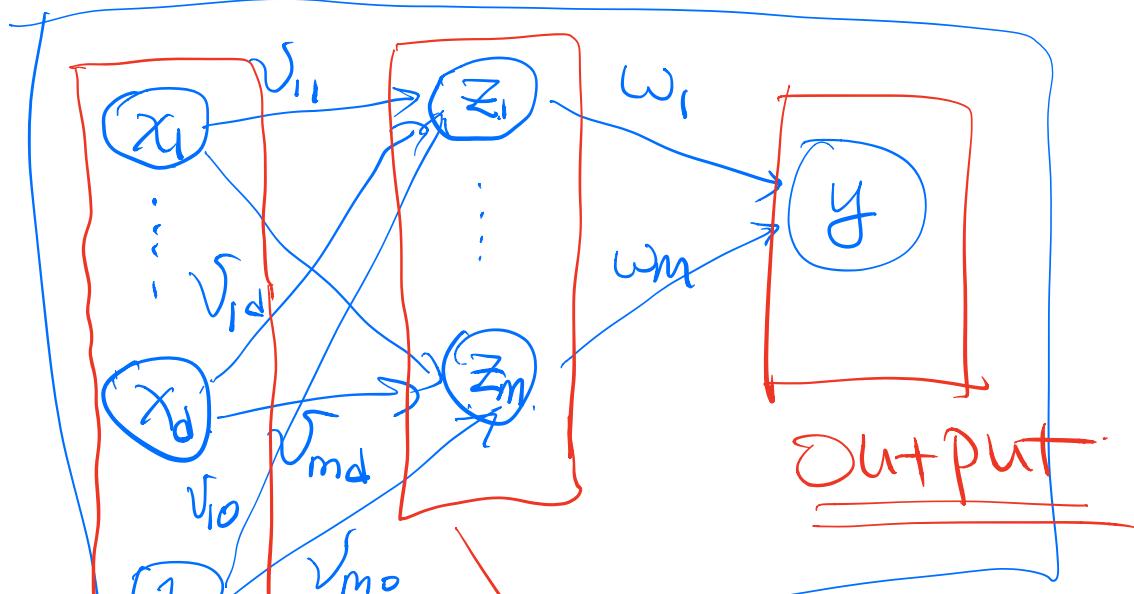
$$y = \sum_{\ell=1}^d v_\ell x_\ell$$



$$y = \sum_{i=1}^m w_i \sigma\left(\sum_{\ell=1}^d v_{i\ell} x_\ell + v_{i0}\right)$$

$\hat{=} z_i$

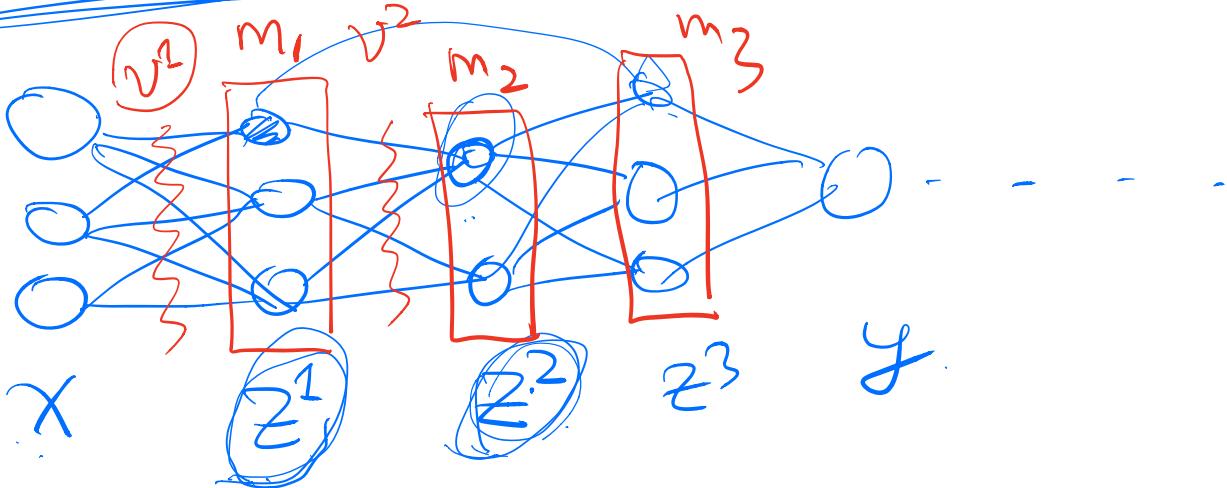
Neural Network



Input
Layer

Hidden
Layer

Deep Neural Networks

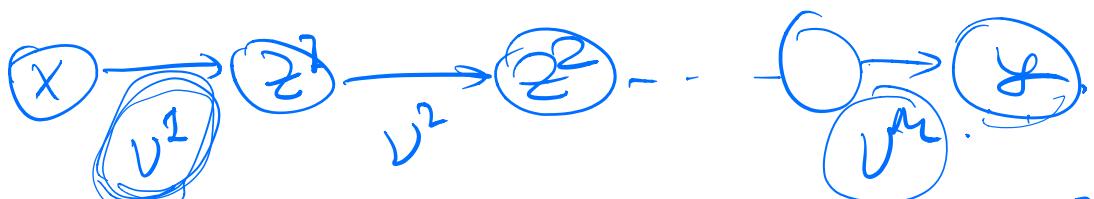


$$z^1 = \underline{\Phi}(x, \underline{V}) = \begin{bmatrix} \underline{\sigma}((\underline{v}_1^1)^T x) \\ \vdots \\ \underline{\sigma}((\underline{v}_{m_1}^1)^T x) \end{bmatrix}$$

$$z^2 = \underline{\Phi}(z^1, \underline{V}^2) = \begin{bmatrix} \underline{\sigma}((\underline{v}_1^2)^T \underline{z}^1) \\ \vdots \\ \underline{\sigma}((\underline{v}_{m_2}^2)^T \underline{z}^1) \end{bmatrix}$$

$$y = \sum_{i=1}^{m_3} w_i z_i^3$$

Back-propagation = Calculation
Gradient of
Deep Neural Nets



$$L(V) = E_D[(y - f(x, V^1 \dots V^m))^2]$$

$$\nabla_{V^1} L(V) = 2 E_D[(f(x, V) - y) \frac{\partial f(x, V)}{\partial V^1}]$$

$$\frac{\partial f(x, V)}{\partial V^1} = \frac{dy}{\partial V^1} = \left(\frac{dy}{\partial z^m} \right) \left(\frac{\partial z^m}{\partial z^{m-1}} \right) \dots \left(\frac{\partial z^3}{\partial z^2} \right) \left(\frac{\partial z^2}{\partial z^1} \right) \left(\frac{\partial z^1}{\partial V^1} \right)$$