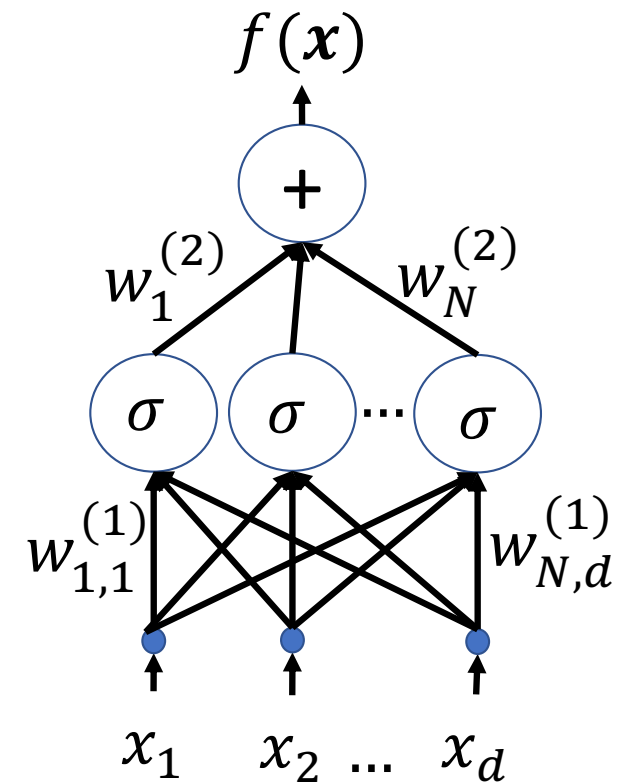


Lecture 10: Nonlinear Regression

Mark Hasegawa-Johnson

These slides are in the public domain



Outline

- From linear to nonlinear regression
- Rectified linear units (ReLU)
- Training a two-layer network: Back-propagation

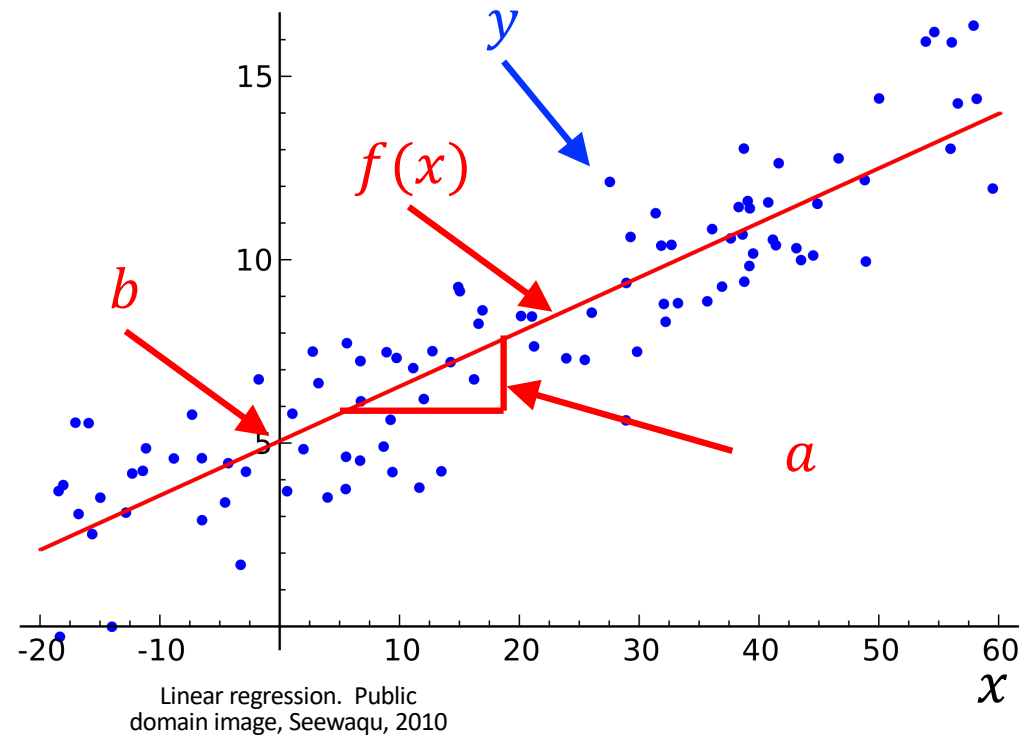
Linear regression

Linear regression is used to estimate a real-valued target variable, y , using a linear function of another variable, x :

$$f(x) = ax + b$$

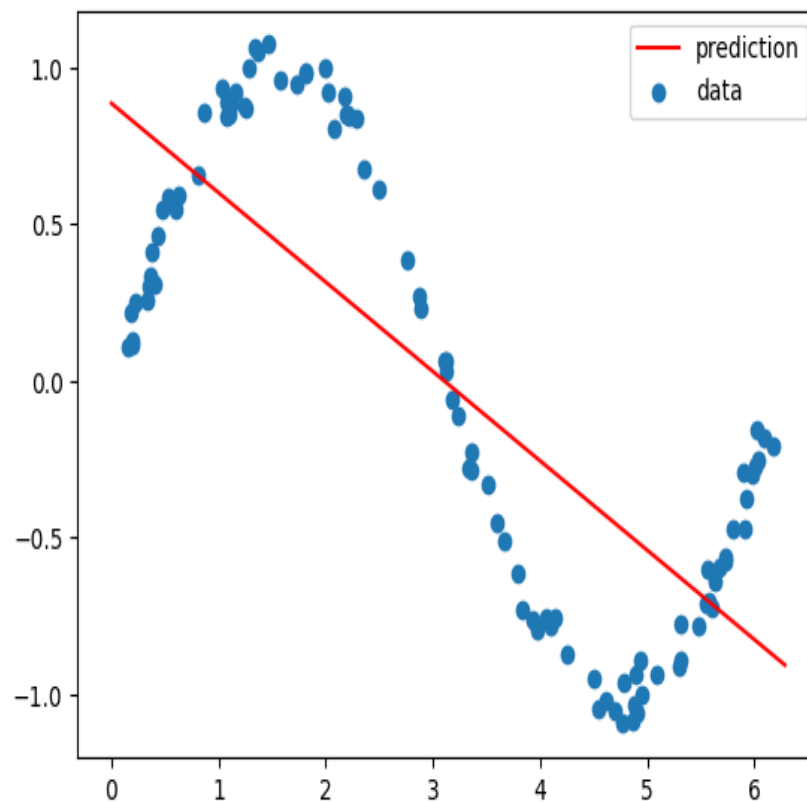
... or of a vector, \mathbf{x} , ...

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$



Linear regression can't fit nonlinear data without nonlinear features

Here's an example from MP1.
Linear regression can't fit nonlinear data unless it has nonlinear features.



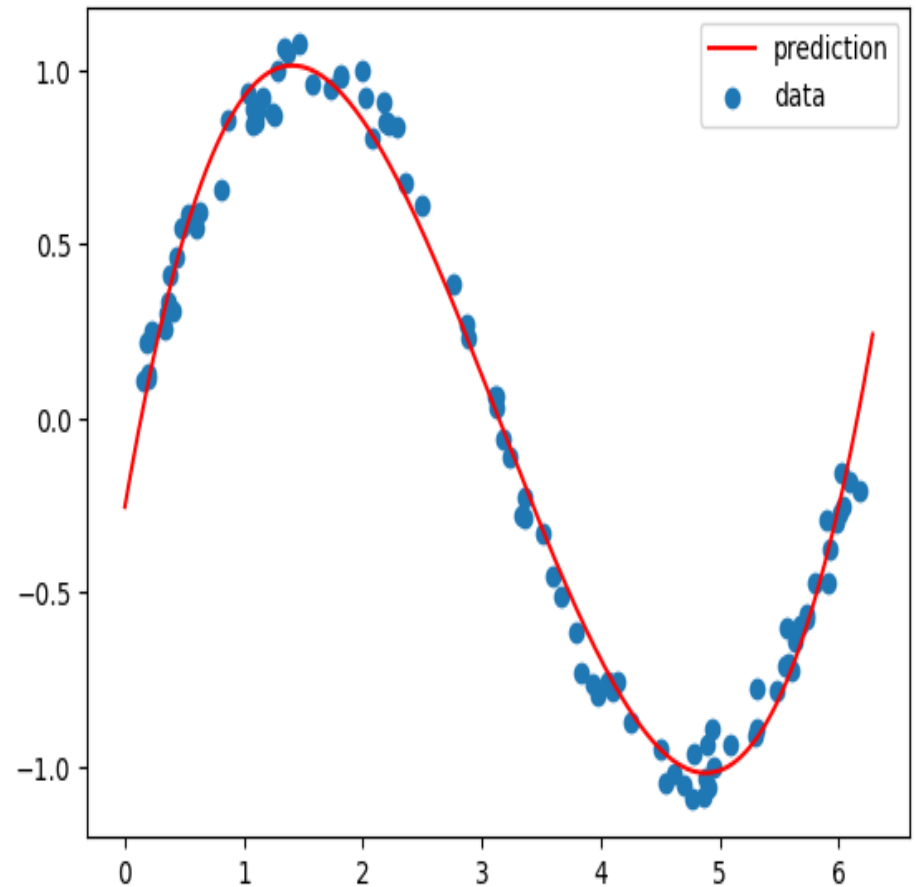
Linear regression can fit nonlinear data by using nonlinear features.

Here's the same example from MP1. Linear regression can fit nonlinear data if it uses nonlinear features! For example, if

$$\mathbf{x} = \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \end{bmatrix}$$

...then $f(x)$ can fit any cubic polynomial, which is sometimes a good enough approximation of a sine wave:

$$f(x) = \mathbf{w}^T \mathbf{x} = w_1 + w_2x + w_3x^2 + w_4x^3$$



Today: Two-layer neural nets

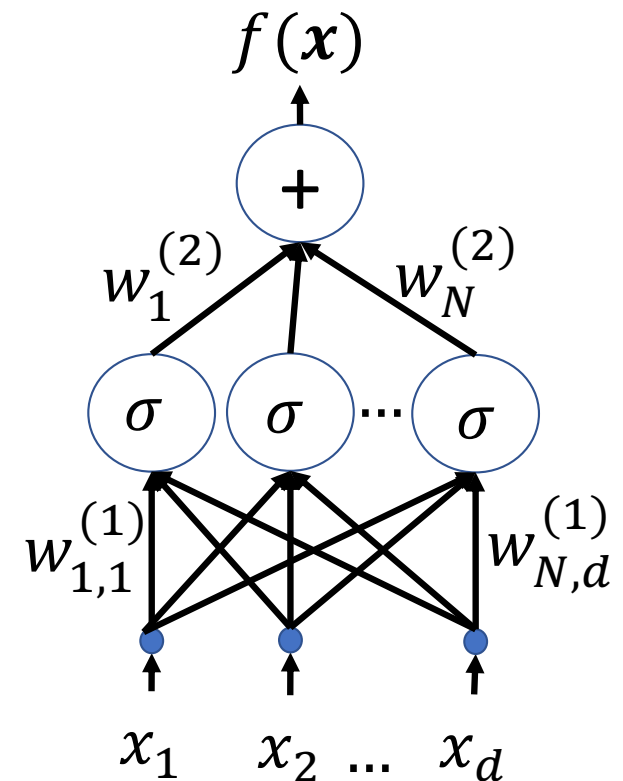
Today we want to consider a more general method for nonlinear regression, called a “two-layer neural net,” given by:

$$f(\mathbf{x}) = \mathbf{w}^{(2),T} \sigma(\mathbf{W}^{(1)} \mathbf{x})$$

...where the first and second-layer weights and biases are

$$\mathbf{W}^{(1)} = \begin{bmatrix} w_{1,1}^{(1)} & \cdots & w_{1,d}^{(1)} \\ \vdots & \ddots & \vdots \\ w_{N,1}^{(1)} & \cdots & w_{N,d}^{(1)} \end{bmatrix}, \quad \mathbf{w}^{(2)} = \begin{bmatrix} w_1^{(2)} \\ \vdots \\ w_N^{(2)} \end{bmatrix}$$

... and $\sigma(\mathbf{z})$ applies the logistic sigmoid to every element of the vector \mathbf{z} .



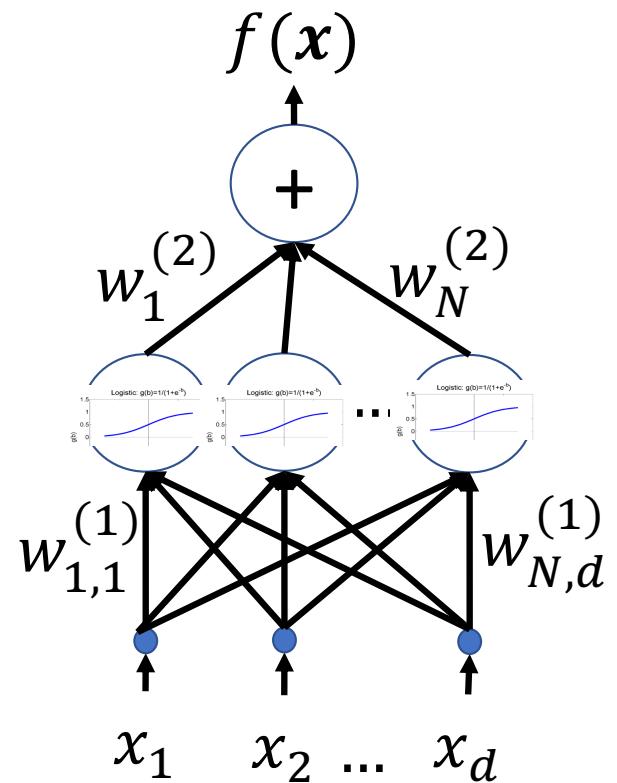
Why does it work?

Suppose we write the first-layer weight matrix as a stack of row vectors, like this: $\mathbf{W}^{(1)} = [\mathbf{w}_1^{(1)} \quad \dots \quad \mathbf{w}_N^{(1)}]^T$. Then the hidden node vector is

$$\sigma(\mathbf{W}^{(1)}\mathbf{x}) = \begin{bmatrix} \sigma(\mathbf{w}_1^{(1),T}\mathbf{x}) \\ \vdots \\ \sigma(\mathbf{w}_N^{(1),T}\mathbf{x}) \end{bmatrix}$$

Each of its elements is

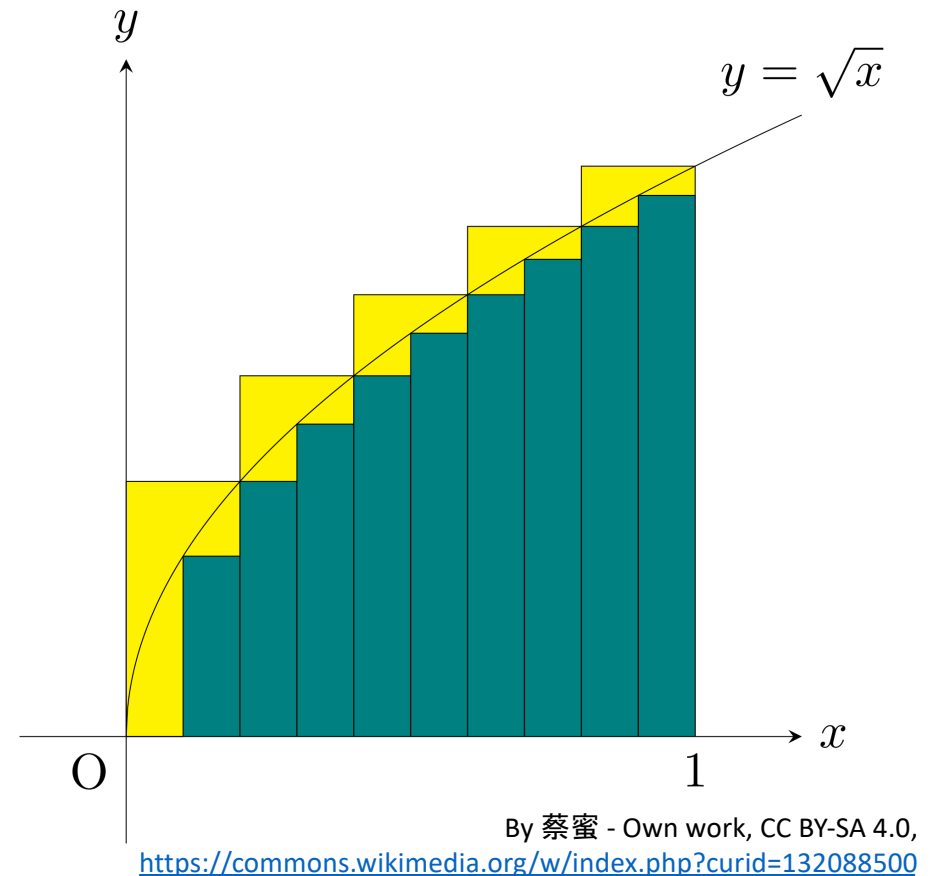
$$\sigma(\mathbf{w}_j^{(1),T}\mathbf{x}) \approx \begin{cases} 0 & \mathbf{w}_j^{(1),T}\mathbf{x} < 0 \\ 1 & \mathbf{w}_j^{(1),T}\mathbf{x} > 0 \end{cases}$$



Why does it work?

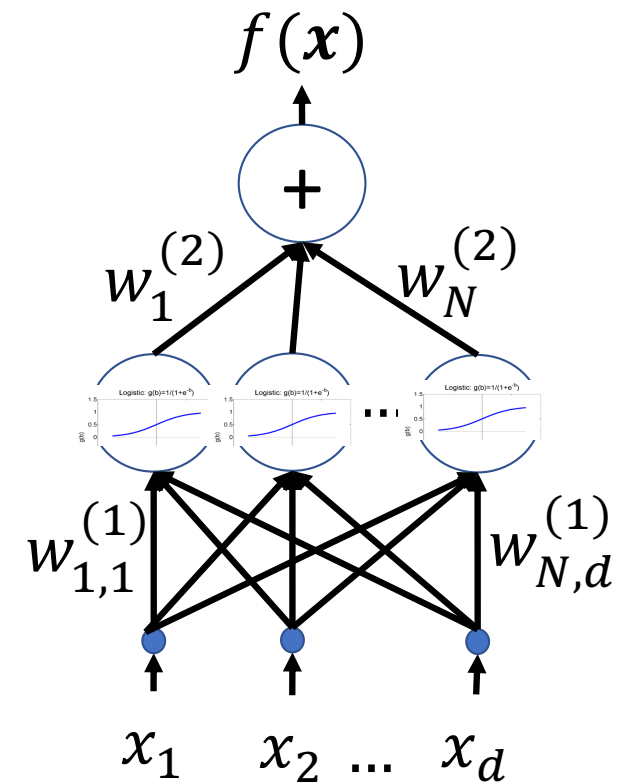
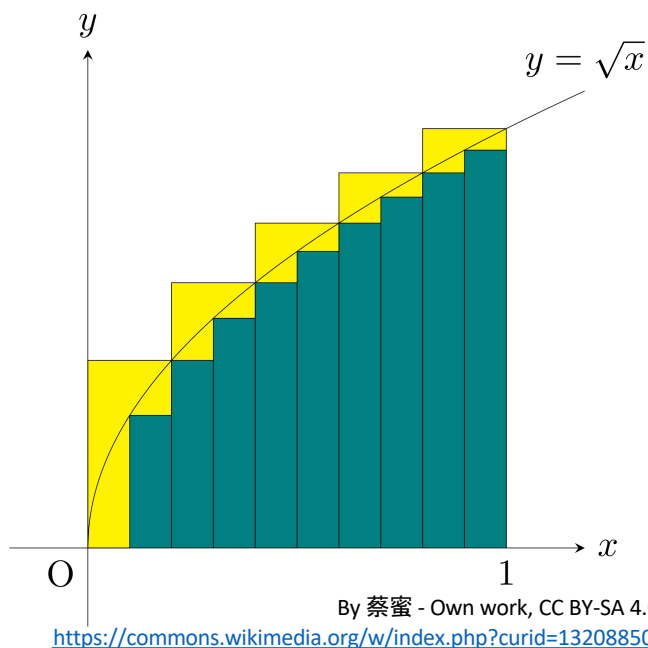
$$f(\mathbf{x}) = \mathbf{w}^{(2),T} \sigma(\mathbf{W}^{(1)} \mathbf{x}) \approx \sum_{j: \mathbf{w}_j^{(1),T} \mathbf{x} > 0} w_j^{(2)}$$

... but here's the cool thing. A long time ago, Isaac Newton proved that any nonlinear function can be approximated by a piece-wise constant function arbitrarily well, if you have enough pieces.



The universal approximation theorem of two-layer neural networks

... and therefore, any nonlinear function can be approximated by a two-layer neural network arbitrarily well, if you have enough hidden nodes.

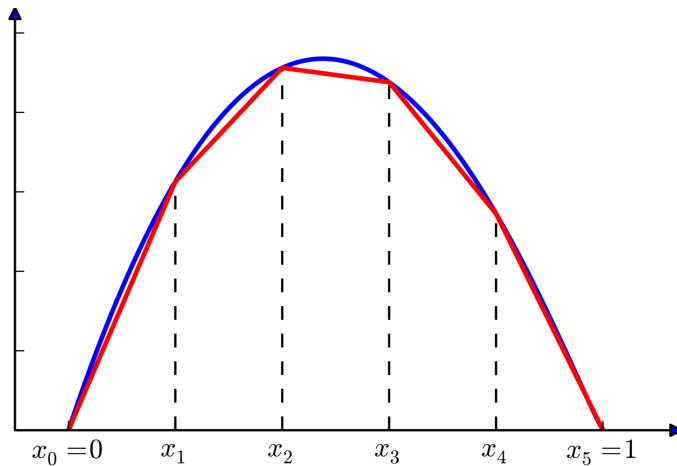


Outline

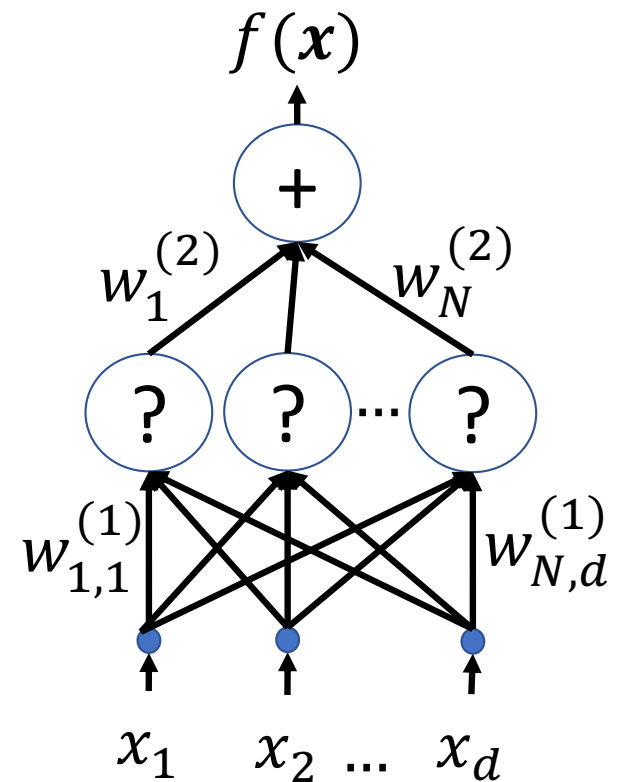
- From linear to nonlinear regression
- Rectified linear units (ReLU)
- Training a two-layer network: Back-propagation

How about piece-wise linear approximations?

What if we want a piece-wise linear output function, instead of piece-wise constant?



Public domain image, Krishnavedala, 2011



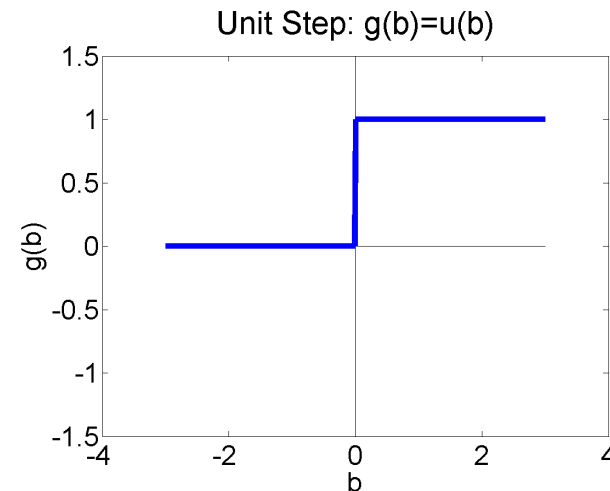
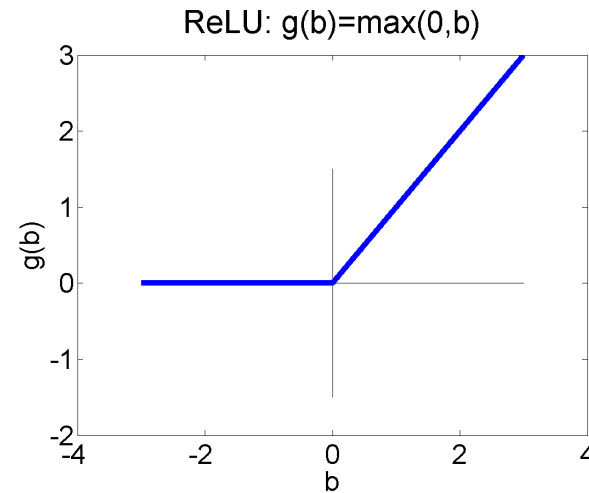
For a PWL neural net, the hidden nodes are ReLU

If the goal is PWL classification boundaries, we can achieve that by using hidden nodes that are the simplest possible PWL function: a Rectified Linear Unit, or ReLU:

$$\text{ReLU}(z) = \max(0, z)$$

This is differentiable everywhere except $z=0$; its derivative is the unit step function:

$$\frac{\partial \text{ReLU}(z)}{\partial z} = u(z) = \begin{cases} 1 & z > 0 \\ 0 & z < 0 \end{cases}$$



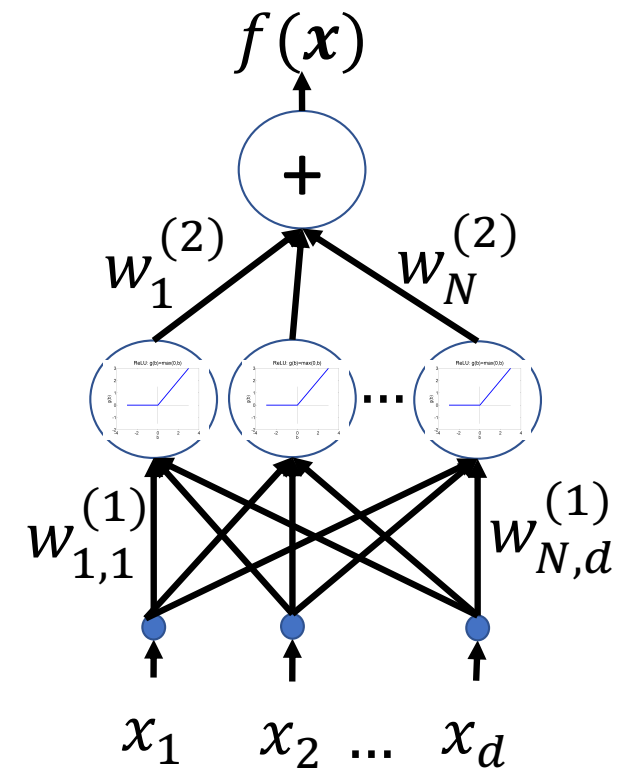
2-layer NN with ReLU hidden nodes

With ReLU hidden nodes, the hidden node vector is

$$\text{ReLU}(\mathbf{W}^{(1)}\mathbf{x}) = \begin{bmatrix} \text{ReLU}(\mathbf{w}_1^{(1),T}\mathbf{x}) \\ \vdots \\ \text{ReLU}(\mathbf{w}_N^{(1),T}\mathbf{x}) \end{bmatrix}$$

Each of its elements is

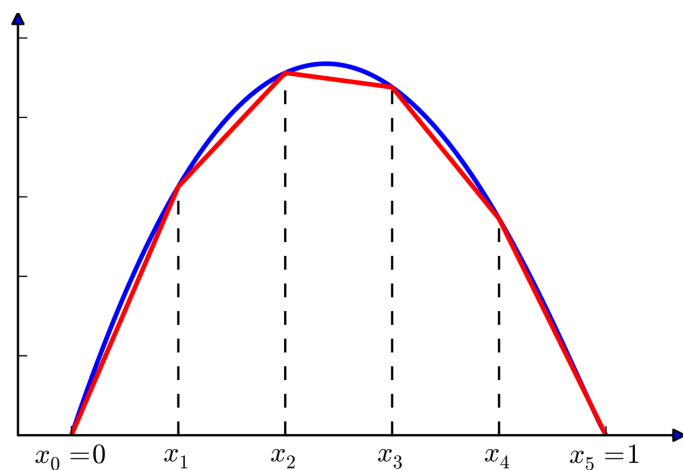
$$\text{ReLU}(\mathbf{w}_j^{(1),T}\mathbf{x}) = \begin{cases} 0 & \text{if } \mathbf{w}_j^{(1),T}\mathbf{x} \leq 0 \\ \mathbf{w}_j^{(1),T}\mathbf{x} & \text{if } \mathbf{w}_j^{(1),T}\mathbf{x} \geq 0 \end{cases}$$



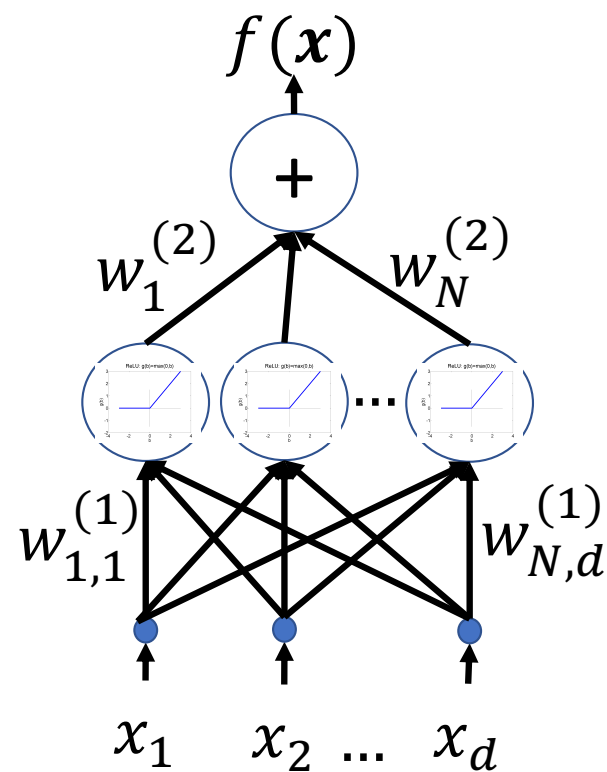
2-layer NN with ReLU hidden nodes

With ReLU hidden nodes, $f(\mathbf{x})$ is a perfectly piece-wise linear function of \mathbf{x} :

$$f(\mathbf{x}) = \mathbf{w}^{(2),T} \text{ReLU}(\mathbf{W}^{(1)}\mathbf{x}) = \sum_{j: \mathbf{w}_j^{(1),T} \mathbf{x} \geq 0} w_j^{(2)} \mathbf{w}_j^{(1),T} \mathbf{x}$$



Public domain image, Krishnavedala, 2011



Outline

- From linear to nonlinear regression
- Rectified linear units (ReLU)
- Training a two-layer network: Back-propagation

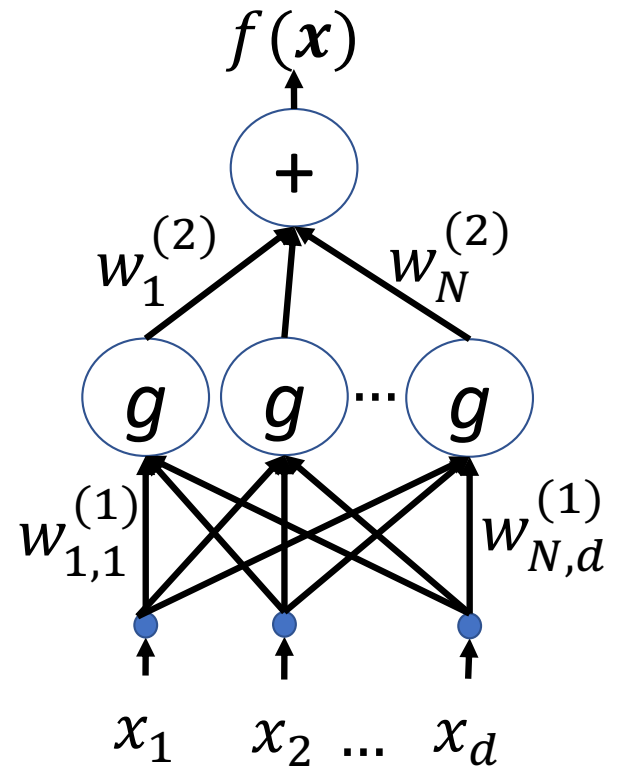
Training a neural net: mean-squared error

Now suppose we have a two-layer NN, with some hidden node nonlinearity $g(\mathbf{z})$ that might be either ReLU or sigmoid:

$$f(\mathbf{x}_i) = \mathbf{w}^{(2),T} g(\mathbf{W}^{(1)} \mathbf{x}_i)$$

Suppose we want to train $\mathbf{w}^{(2)}$ and $\mathbf{W}^{(1)}$ to minimize MSE:

$$\mathcal{L} = \frac{1}{2n} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2$$



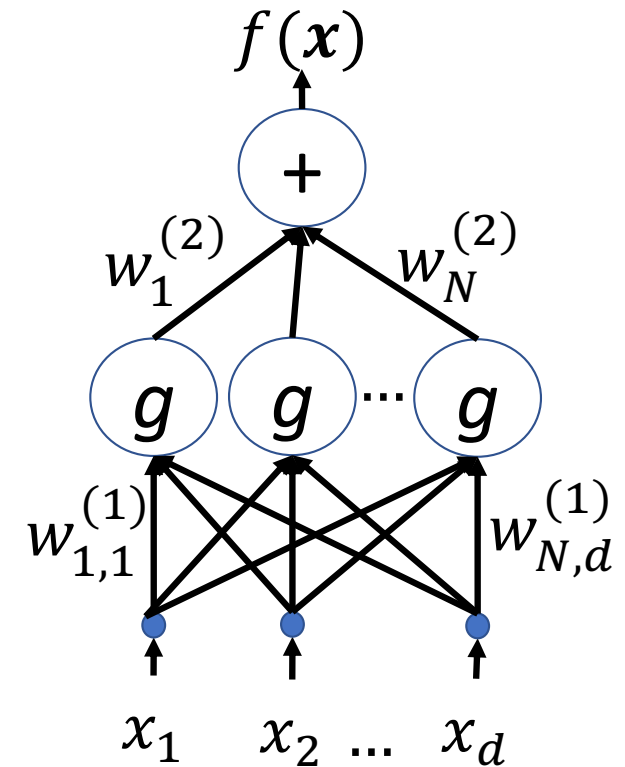
Training a neural net: gradient descent

We can do this using gradient descent:

$$\mathbf{W}^{(1)} \leftarrow \mathbf{W}^{(1)} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(1)}}$$

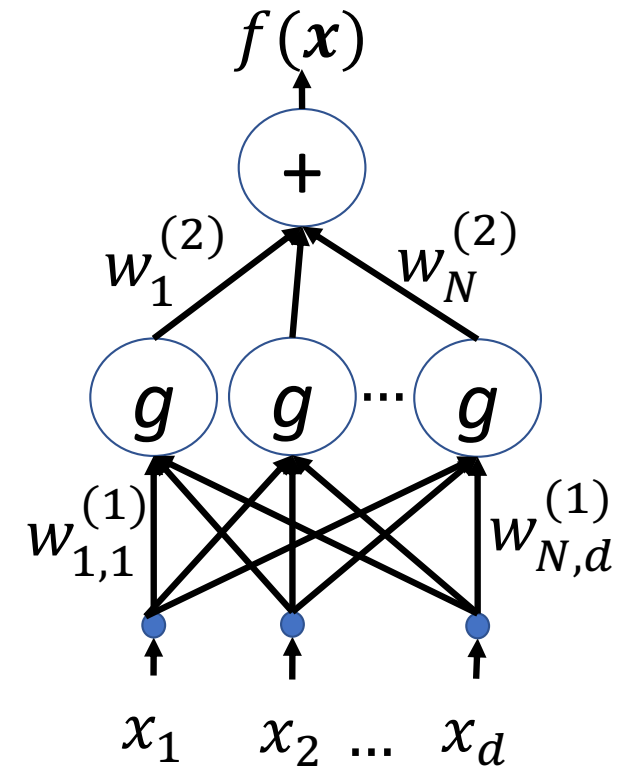
$$\mathbf{w}^{(2)} \leftarrow \mathbf{w}^{(2)} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(2)}}$$

How do we find those derivatives?



Training a neural net: back-propagation

We can find the derivatives using the chain rule of calculus. Since the chain rule propagates backward through the network, this is called “back-propagation.”



Training a neural net: back-propagation

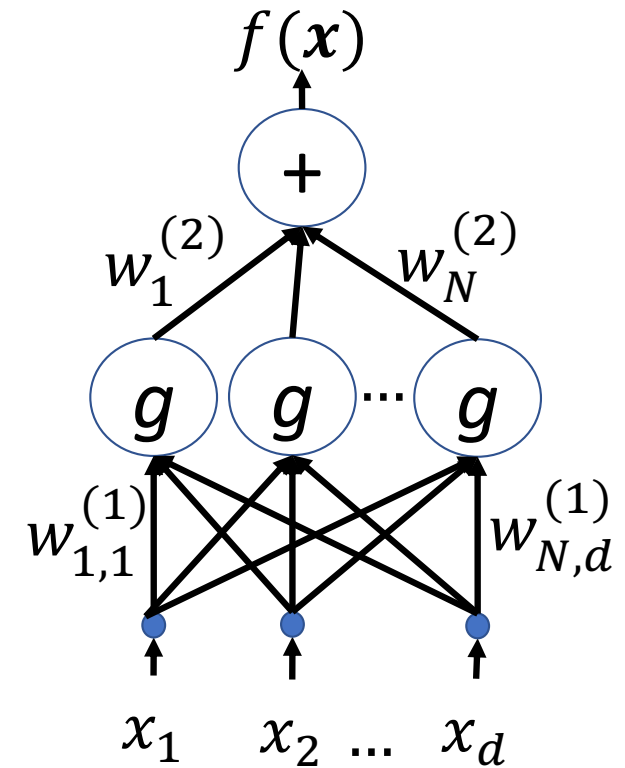
For example:

$$\mathcal{L} = \frac{1}{2n} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2$$

$$\frac{\partial \mathcal{L}}{\partial f(\mathbf{x}_i)} = \frac{1}{n} (f(\mathbf{x}_i) - y_i)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(2)}} = \sum_{i=1}^n \frac{\partial \mathcal{L}}{\partial f(\mathbf{x}_i)} \frac{\partial f(\mathbf{x}_i)}{\partial \mathbf{w}^{(2)}}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(1)}} = \sum_{i=1}^n \frac{\partial \mathcal{L}}{\partial f(\mathbf{x}_i)} \frac{\partial f(\mathbf{x}_i)}{\partial \mathbf{W}^{(1)}}$$

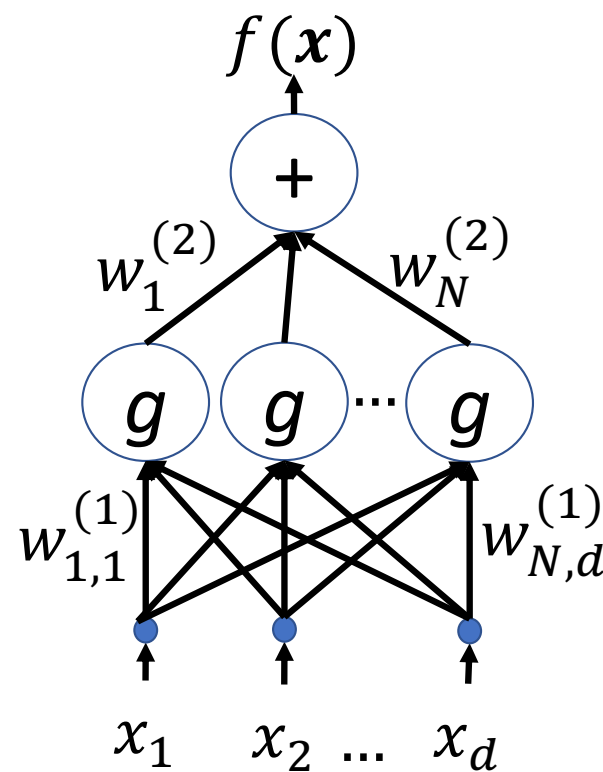


Back-propagation

The derivatives are not always easy to find in matrix form. But since both ReLU and sigmoid are differentiable, we can always find the derivatives in scalar form!

$$f(\mathbf{x}_i) = \mathbf{w}^{(2),T} g(\mathbf{W}^{(1)} \mathbf{x}_i)$$

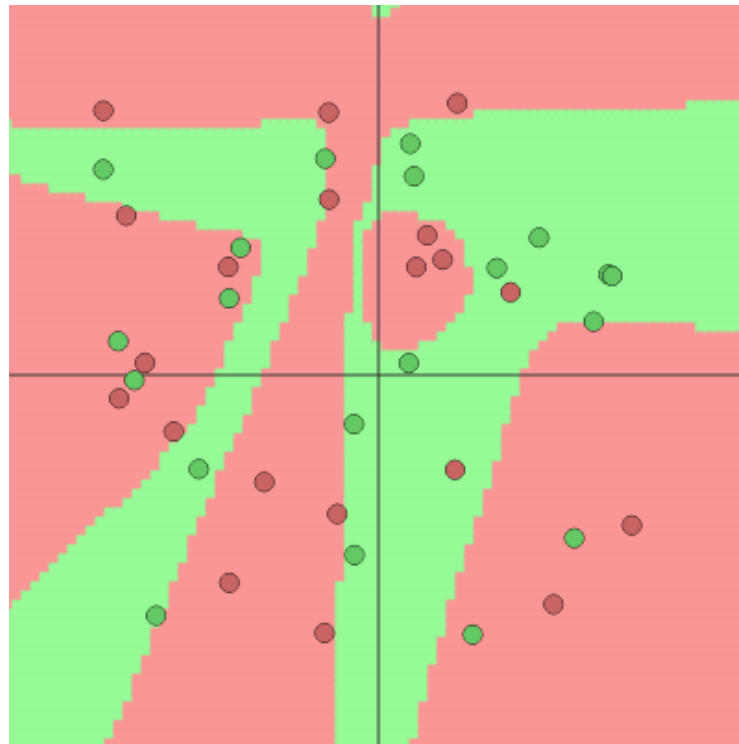
$$\frac{\partial f(\mathbf{x}_i)}{\partial \mathbf{W}^{(1)}} = \begin{bmatrix} \frac{\partial f(\mathbf{x}_i)}{\partial w_{1,1}^{(1)}} & \dots & \frac{\partial f(\mathbf{x}_i)}{\partial w_{1,d}^{(1)}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f(\mathbf{x}_i)}{\partial w_{N,1}^{(1)}} & \dots & \frac{\partial f(\mathbf{x}_i)}{\partial w_{N,d}^{(1)}} \end{bmatrix}$$



Try the quiz!

Go to PrairieLearn, try the quiz!

Approximating an arbitrary nonlinear boundary using a two-layer network



<https://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>

Stochastic gradient descent

You can also reduce computation per step by using SGD:

- From a very large training dataset, randomly choose a training token (\mathbf{x}_i, y_i) . Forward-propagate to find the neural net prediction, $f(\mathbf{x}_i)$, and the loss

$$\mathcal{L}_i = \frac{1}{2n} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2$$

- Back-propagate to find the gradients, then do gradient descent:

$$\mathbf{W}^{(1)} \leftarrow \mathbf{W}^{(1)} - \eta \frac{\partial \mathcal{L}_i}{\partial \mathbf{W}^{(1)}}$$

$$\mathbf{w}^{(2)} \leftarrow \mathbf{w}^{(2)} - \eta \frac{\partial \mathcal{L}_i}{\partial \mathbf{w}^{(2)}}$$

- Repeat!

Summary

- Piece-wise constant nonlinear regression:

$$f(\mathbf{x}) = \mathbf{w}^{(2),T} \sigma(\mathbf{W}^{(1)} \mathbf{x})$$

- Piece-wise linear regression:

$$f(\mathbf{x}) = \mathbf{w}^{(2),T} \text{ReLU}(\mathbf{W}^{(1)} \mathbf{x})$$

- Back-propagation:

$$\mathbf{W}^{(1)} \leftarrow \mathbf{W}^{(1)} - \eta \frac{\partial \mathcal{L}_i}{\partial \mathbf{W}^{(1)}}$$

$$\mathbf{w}^{(2)} \leftarrow \mathbf{w}^{(2)} - \eta \frac{\partial \mathcal{L}_i}{\partial \mathbf{w}^{(2)}}$$