Announcements

Assignments:

- HW4
 - Due date Mon, 2/24, 11:59 pm
- HW5
 - Out tomorrow
 - Due date Thu, 2/27, 11:59 pm

Midterm Conflicts

- See Piazza post
- Due 11:59pm on Wednesday the 19th of February



Plan

Last time

- Decision Boundaries
- Gaussian Generative Models
- Neural Networks

Today

- Neural Networks
 - Universal Approximation
 - Optimization / Backpropagation
 - (Convolutional Neural Networks)

Introduction to Machine Learning

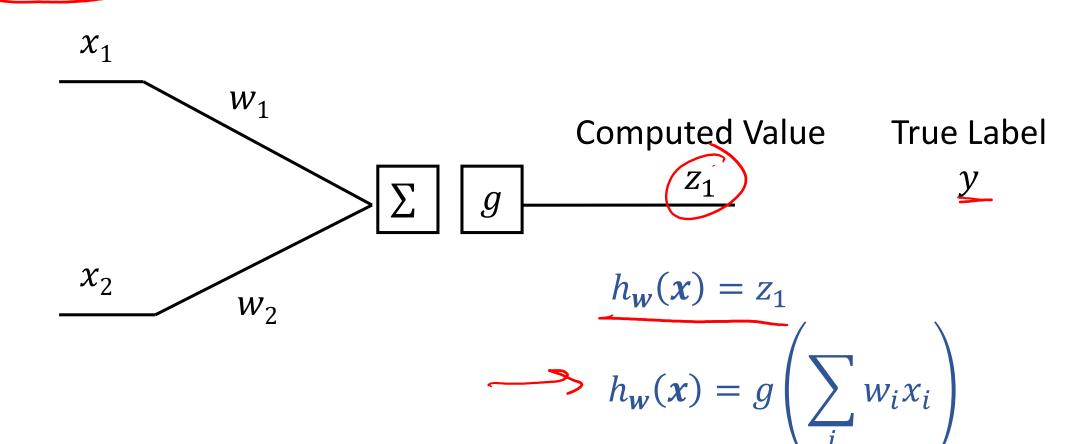
Neural Networks

Instructor: Pat Virtue

Single Neuron

Single neuron system

- Perception (if g is step function)
- Logistic regression (if g is sigmoid)



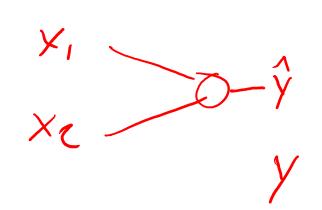
Optimizing

How do we find the "best" set of weights?

$$l(\dot{y}, \dot{\dot{y}}) = \sum_{n=1}^{N} (y_n - \dot{y}_n)^2$$

$$J(u) = l(y, h_w(x))$$

$$= l(y, g(\ddot{w}^T \dot{x}))$$

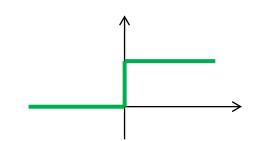


 $h_{w}(x) = g \int_{w}^{w} w$

Activation Functions

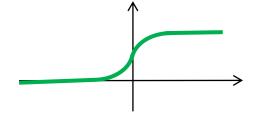
It would be really helpful to have a g(z) that was nicely differentiable

■ Hard threshold:
$$g(z) = \begin{cases} 1 & z \ge 0 \\ 0 & z < 0 \end{cases}$$
 $\frac{dg}{dz} = \begin{cases} 0 & z \ge 0 \\ 0 & z < 0 \end{cases}$



• Sigmoid:
$$g(z) = \frac{1}{1+e^{-z}}$$

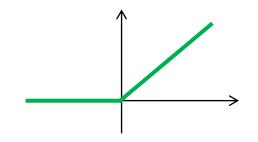
$$\frac{dg}{dz} = g(z) (1 - g(z))$$



(Softmax)

ReLU:

$$g(z) = max(0, z) \qquad \frac{dg}{dz} = \begin{cases} 1 & z \ge 0 \\ 0 & z < 0 \end{cases}$$



Loss Functions

Regression

■ MSE (SSE): $\ell(y, \hat{y}) = ||y - \hat{y}||_2^2 = \sum_n (y_n - \widehat{y_n})^2$

Classification

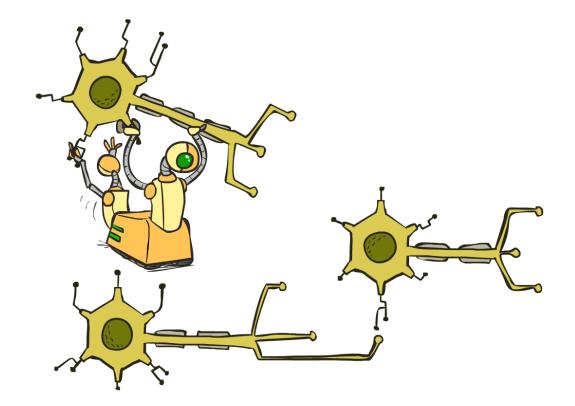
■ Cross entropy: $\ell(\boldsymbol{y}, \widehat{\boldsymbol{y}}) = -\sum_n y_n \log \widehat{y}_n$



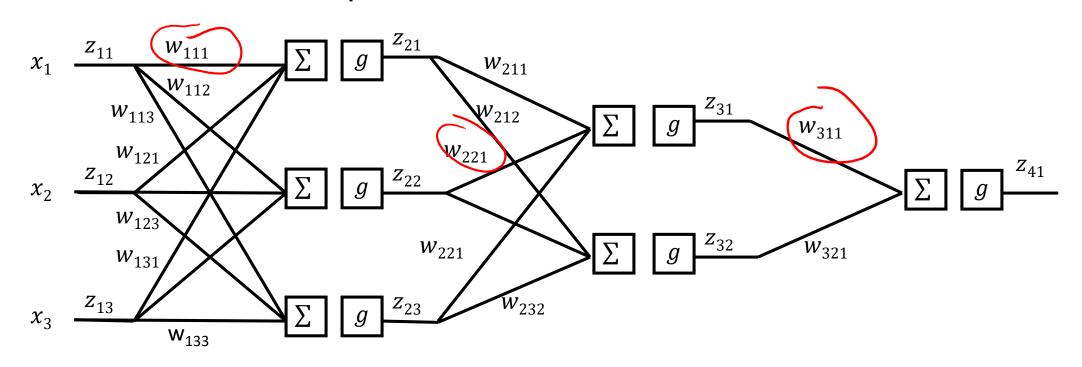
Multilayer Perceptrons

A *multilayer perceptron* is a feedforward neural network with at least one *hidden layer* (nodes that are neither inputs nor outputs)

MLPs with enough hidden nodes can represent any function



Neural Network Equations



$$h_{w}(x) = z_{4,1}$$

$$z_{4,1} = g(\sum_{i} w_{3,i,1} z_{3,i})$$

$$z_{3,1} = g(\sum_{i} w_{2,i,1} z_{2,i})$$

$$z_{d,1} = g(\sum_{i} w_{d-1,i,1} z_{d-1,i})$$

$$h_{w}(x) = g\left(\sum_{k} w_{3,k,1} g\left(\sum_{j} w_{2,j,k} g\left(\sum_{i} w_{1,i,j} x_{i}\right)\right)\right)$$

$$\sum_{k} w_{3,k,1} g\left(\sum_{j} w_{2,j,k} g\left(\sum_{i} w_{1,i,j} x_{i}\right)\right)\right)$$

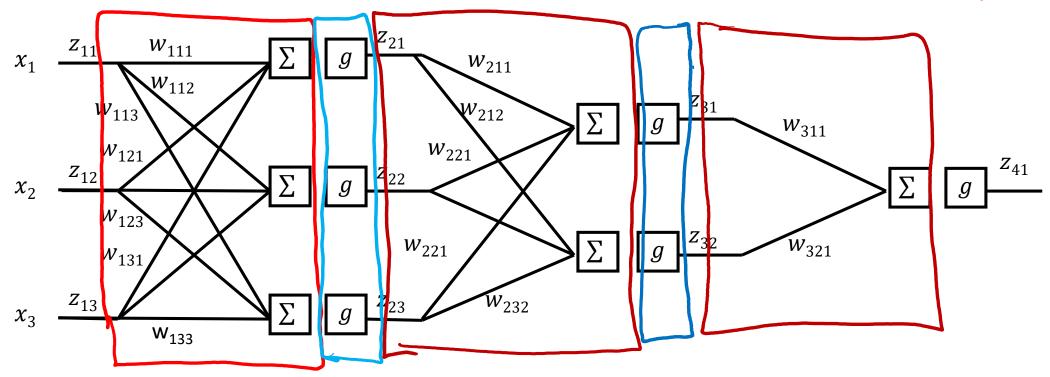
Optimizing

How do we find the "best" set of weights?

$$h_w(x) = g\left(\sum_k w_{3,k,1} \ g\left(\sum_j w_{2,j,k} \ g\left(\sum_i w_{1,i,j} \ x_i\right)\right)\right)$$

Neural Network Equations





How would you represent this specific network in PyTorch?

$$\stackrel{in, out}{\chi} \longrightarrow Linear(3,3) \rightarrow ReLU \rightarrow Linear(3,2) \rightarrow ReLU$$

$$\stackrel{in, out}{\to} Linear(2,1) \rightarrow ReLU$$

Neural Networks Properties

Practical considerations

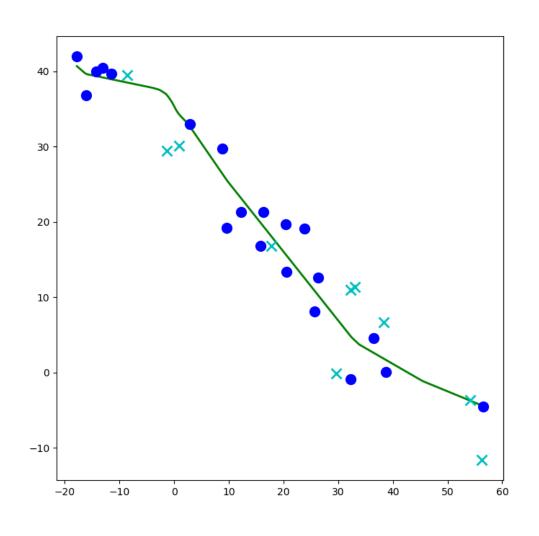
- Large number of neurons
 - Danger for overfitting
- Modelling assumptions vs data assumptions trade-off
- Gradient descent can easily get stuck local optima

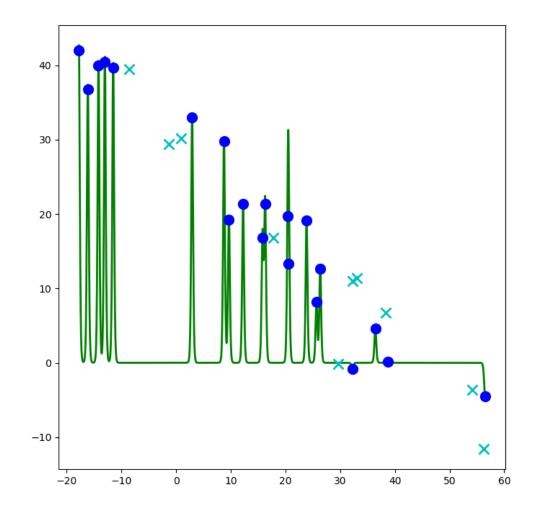
What if there are no non-linear activations?

 A deep neural network with only linear layers can be reduced to an exactly equivalent single linear layer

Universal Approximation Theorem:

 A two-layer neural network with a sufficient number of neurons can approximate any continuous function to any desired accuracy.

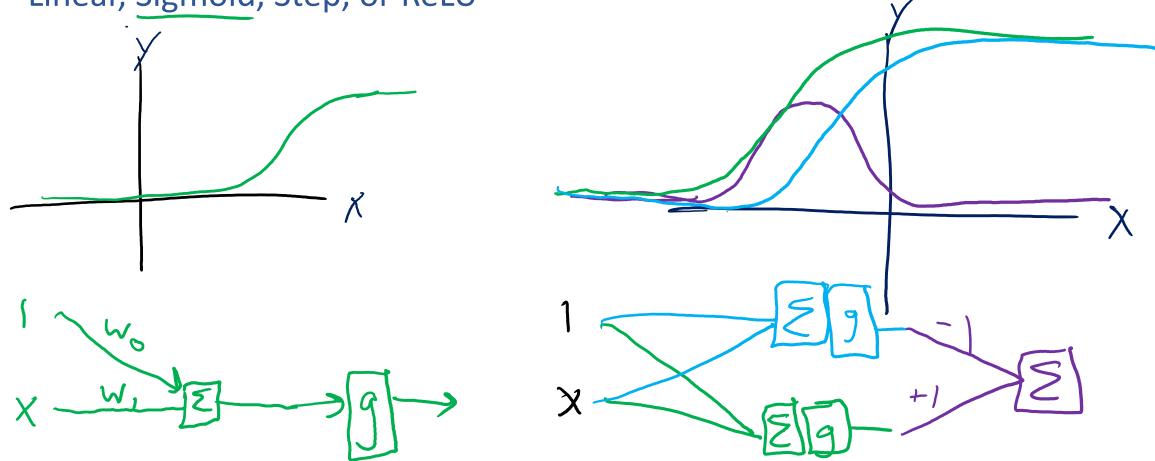


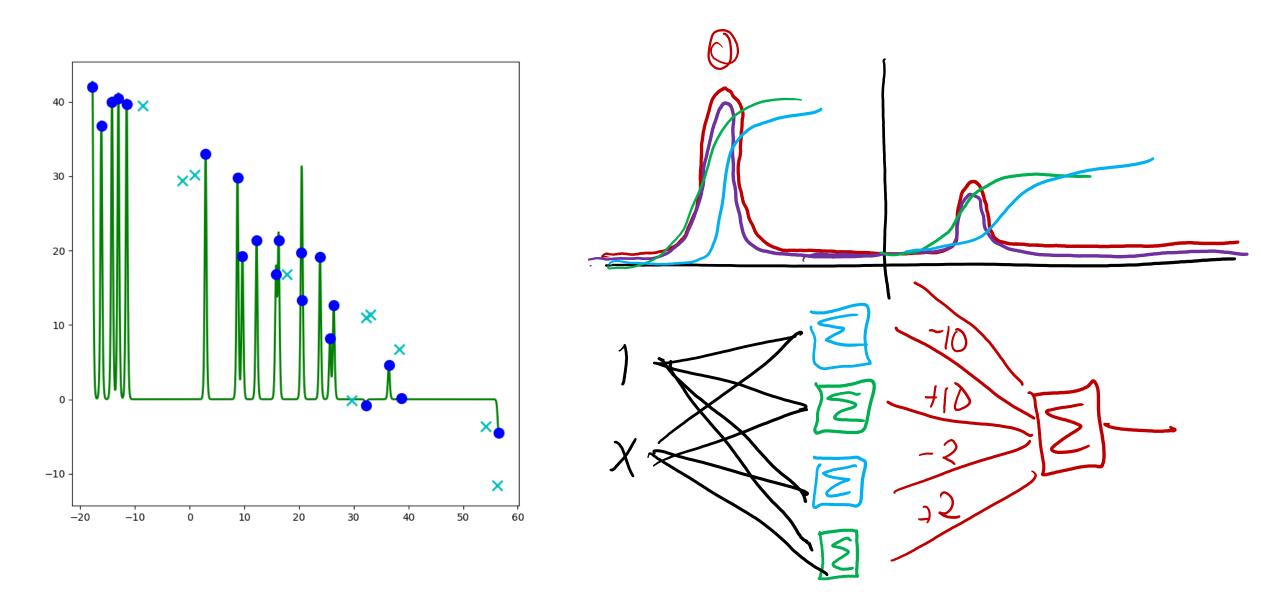


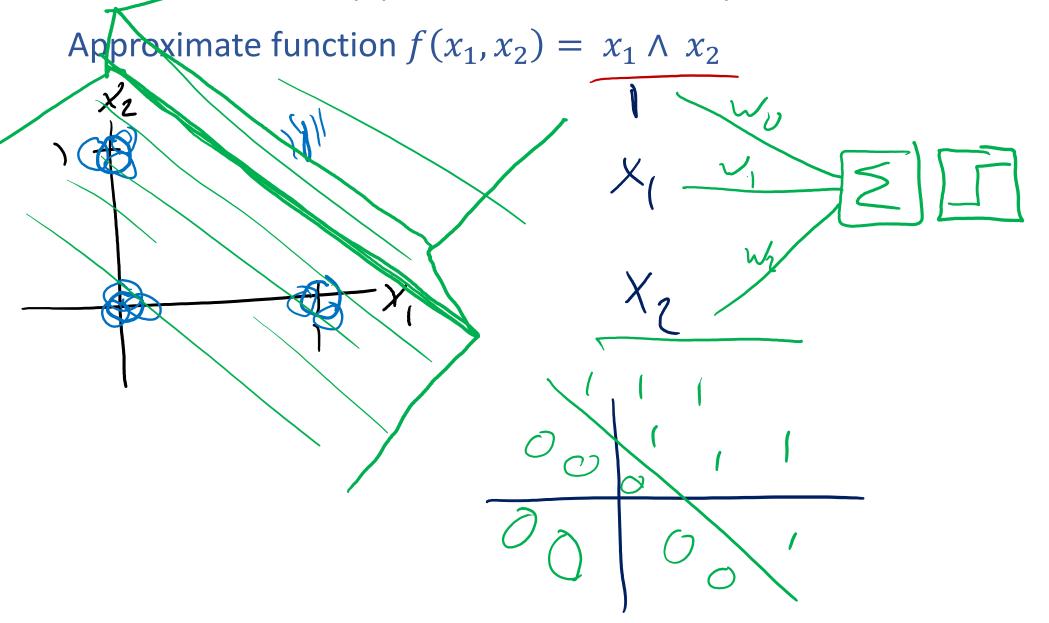
Design a network to approximate this function using: Linear, Sigmoid, Step, or ReLU

Design a network to approximate this function using:

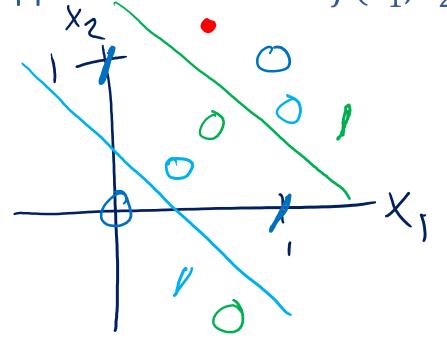
Linear, Sigmoid, Step, or ReLU

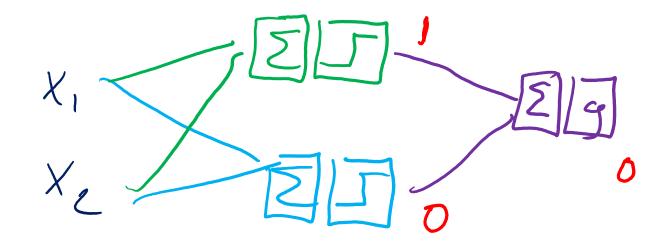






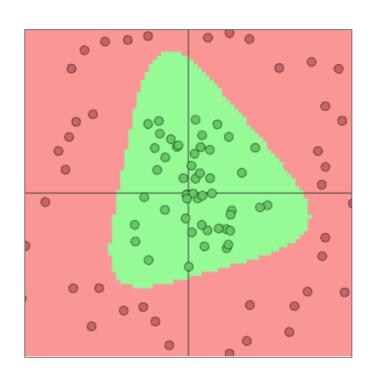
Approximate function $f(x_1, x_2) = x_1 \oplus x_2$

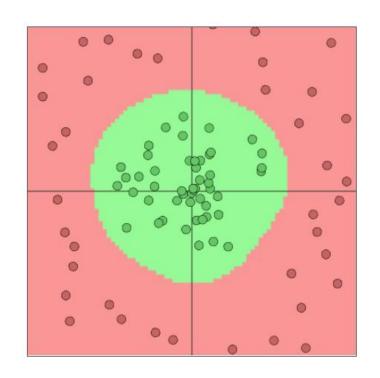


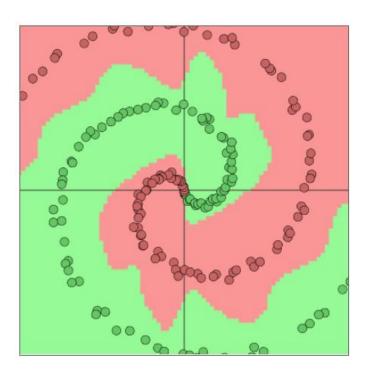


Approximate arbitrary decision boundary

Approximate arbitrary decision boundary







Network Optimization

Reminder: Calculus Chain Rule (scalar version)

$$y = f(z)$$

$$z = g(x)$$

$$y = f(g(x))$$

$$\frac{dy}{dx} = \frac{dy}{dz}\frac{dz}{dx}$$

Network Optimization

$$J(\mathbf{w}) = z_3$$

$$z_3 = f_3(w_3, z_2)$$

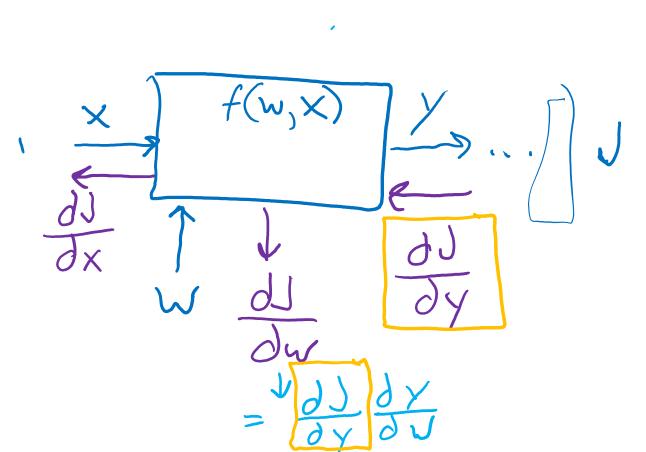
$$z_2 = f_2(w_2, z_1)$$

$$z_1 = f_1(w_1, x)$$

$$\frac{\partial J}{\partial w_3} = \frac{\partial J}{\partial z_3} \frac{\partial z_3}{\partial w_3}$$

$$\frac{\partial J}{\partial w_2} = \frac{\partial J}{\partial z_3} \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial z_3} \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$



Backpropagation (so-far)

Compute derivatives per layer, utilizing previous derivatives

Objective: J(w)

Arbitrary layer: y = f(x, w)

Need: