# 10-301/601: Introduction to Machine Learning Lecture 11 – Neural Networks

Matt Gormley & Henry Chai 2/17/25

#### Front Matter

- Announcements
  - Exam 1 on 2/17 (today!) from 7 PM 9 PM
    - Make sure you check the seating chart (on <u>Piazza</u>)
       so that you know where you're going tonight!
  - Homework 4 released 2/17 (today!), due 2/26 at 11:59 PM

#### Motivation: Regularization

• Occam's Razor: prefer the simplest hypothesis

- What does it mean for a hypothesis (or model) to be simple?
  - small number of features (model selection)
  - 2. small number of "important" features (shrinkage)

#### Regularization

- **Given** objective function:  $J(\theta)$
- Goal is to find:  $\hat{\boldsymbol{\theta}} = \operatorname{argmin} J(\boldsymbol{\theta}) + \lambda r(\boldsymbol{\theta})$
- **Key idea:** Define regularizer  $r(\theta)$  s.t. we tradeoff between fitting the data and keeping the model simple
- Choose form of  $\mathbf{r}(\boldsymbol{\theta})$ :

   Example: q-norm (usually p-norm):  $\|\boldsymbol{\theta}\|_q = \left(\sum_{m=1}^M |\theta_m|^q\right)^{\overline{q}}$

q	$r(oldsymbol{ heta})$	yields parame- ters that are	name	optimization notes
0	$  \boldsymbol{\theta}  _0 = \sum \mathbb{1}(\theta_m \neq 0)$	zero values	Lo reg.	no good computa- tional solutions
$\frac{1}{2}$	$  oldsymbol{ heta}  _1 = \sum   heta_m  \ (  oldsymbol{ heta}  _2)^2 = \sum  heta_m^2$	zero values small values	L1 reg. L2 reg.	subdifferentiable differentiable

#### Regularization Examples

Add an L2 regularizer to Linear Regression (aka. Ridge Regression)

$$J_{\mathsf{RR}}(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \lambda ||\boldsymbol{\theta}||_2^2$$

$$= \frac{1}{N} \sum_{i=1}^{N} \frac{1}{2} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)})^2 + \lambda \sum_{m=1}^{M} \theta_m^2$$

Add an L1 regularizer to Linear Regression (aka. LASSO)

$$J_{\text{LASSO}}(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \lambda ||\boldsymbol{\theta}||_{1}$$

$$= \frac{1}{N} \sum_{i=1}^{N} \frac{1}{2} (\boldsymbol{\theta}^{T} \mathbf{x}^{(i)} - y^{(i)})^{2} + \lambda \sum_{m=1}^{M} |\theta_{m}|$$

#### Regularization Examples

Add an L2 regularizer to Logistic Regression

$$J'(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \lambda ||\boldsymbol{\theta}||_2^2$$

$$= \frac{1}{N} \sum_{i=1}^N -\log p(y^{(i)} \mid \mathbf{x}^{(i)}, \boldsymbol{\theta}) + \lambda \sum_{m=1}^M \theta_m^2$$

Add an L1 regularizer to Logistic Regression

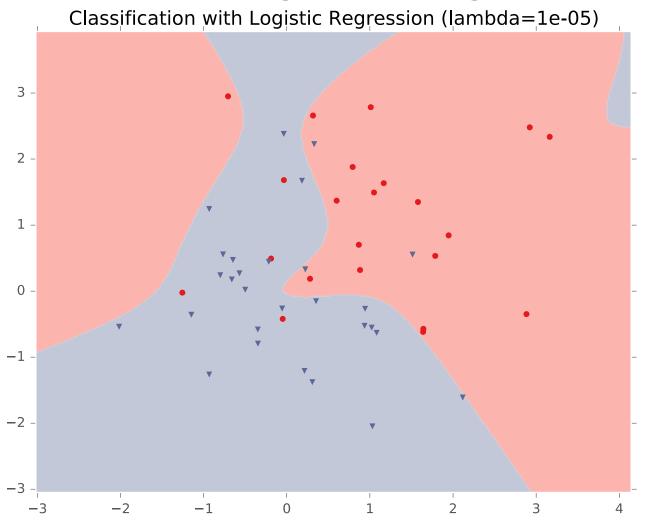
$$J'(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \lambda ||\boldsymbol{\theta}||_{1}$$

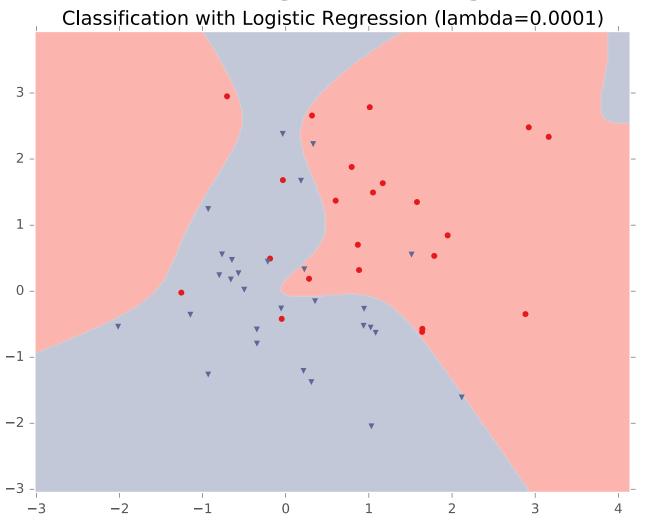
$$= \frac{1}{N} \sum_{i=1}^{N} -\log p(y^{(i)} \mid \mathbf{x}^{(i)}, \boldsymbol{\theta}) + \lambda \sum_{m=1}^{M} |\theta_{m}|$$

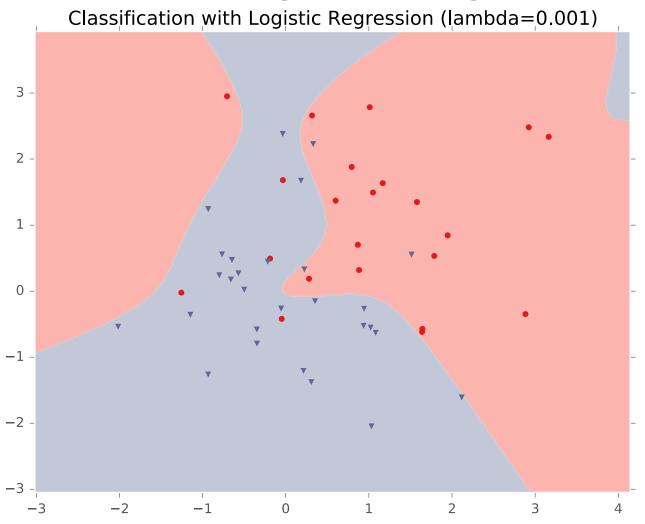
# REGULARIZATION EXAMPLE: LOGISTIC REGRESSION

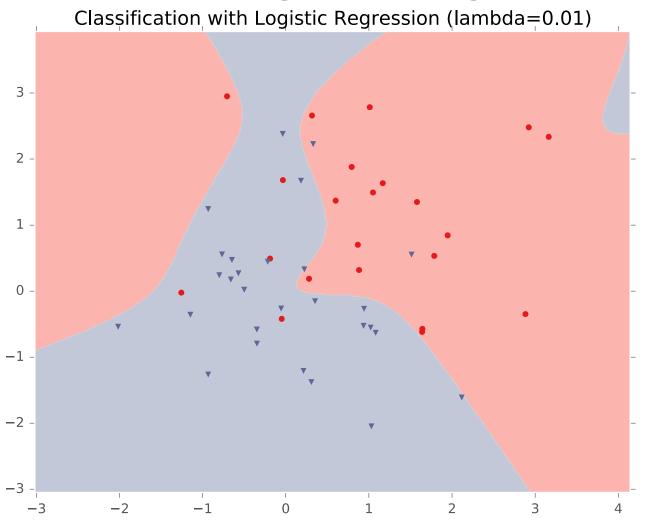


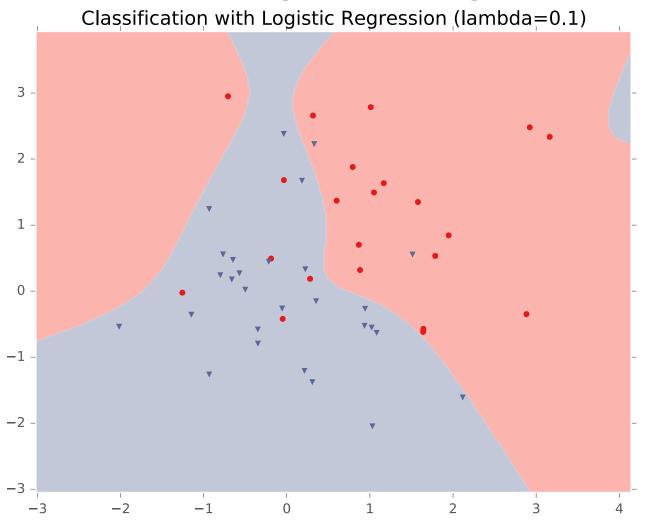
- For this example, we construct nonlinear features
   (i.e. feature engineering)
- Specifically, we add
   polynomials up to order 9 of
   the two original features x<sub>1</sub>
   and x<sub>2</sub>
- Thus our classifier is linear in the high-dimensional feature space, but the decision boundary is nonlinear when visualized in low-dimensions (i.e. the original two dimensions)

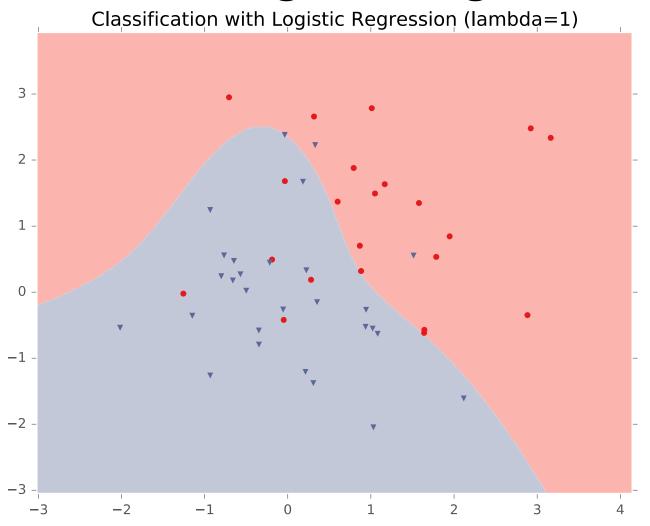


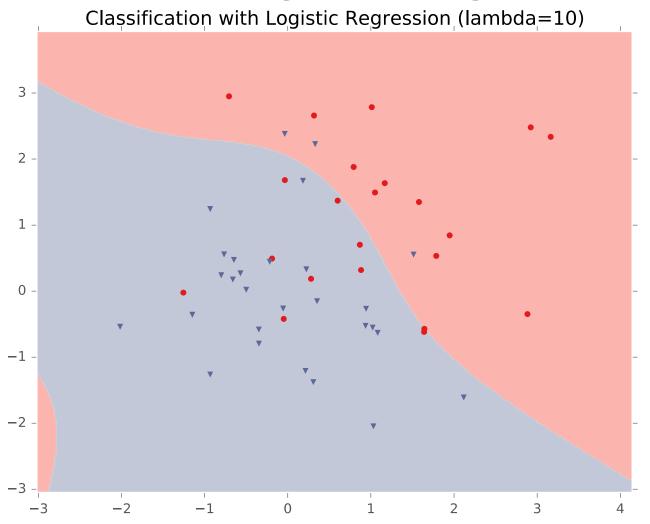


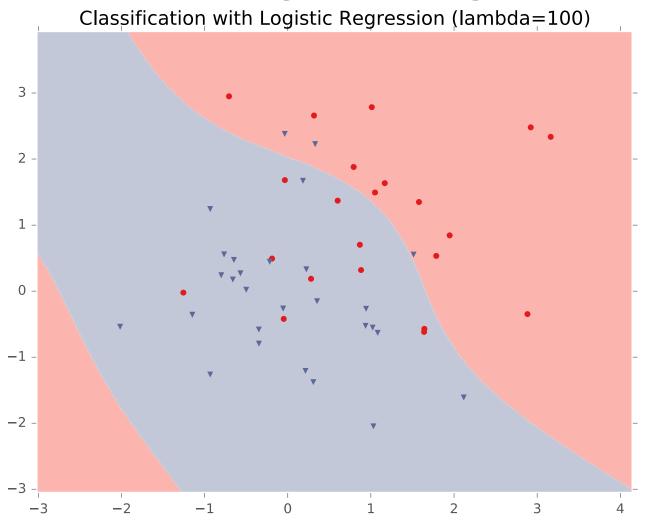


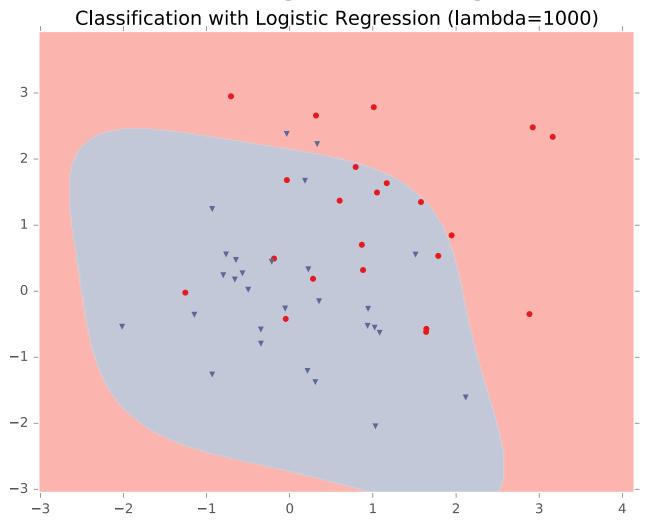


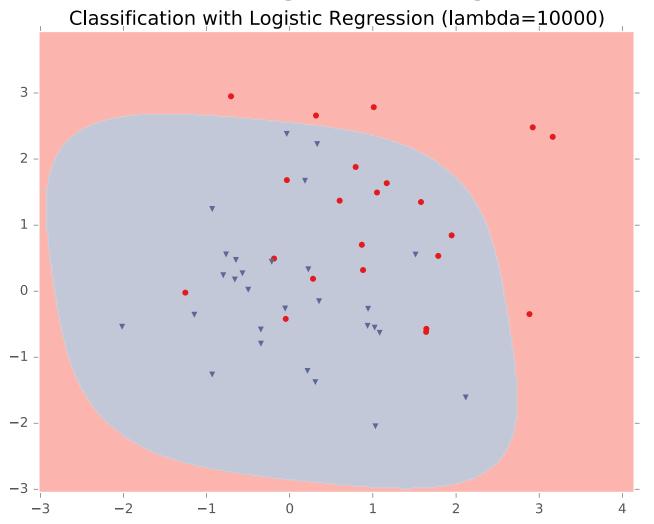


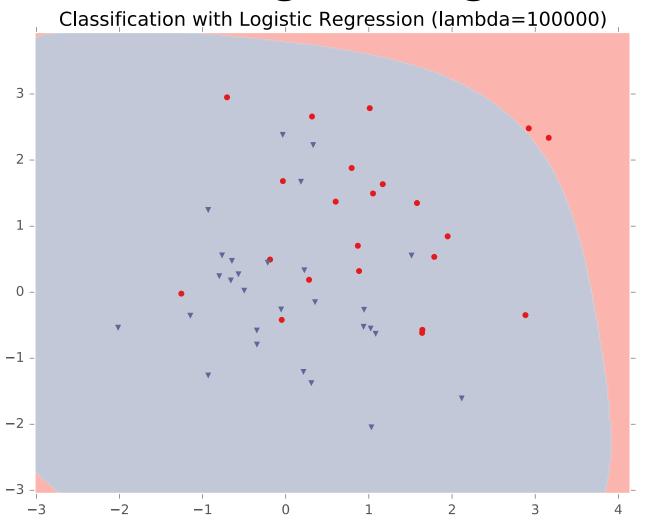


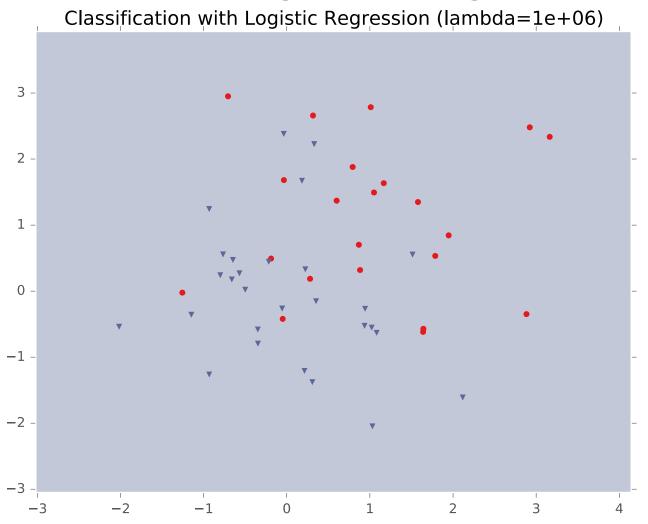


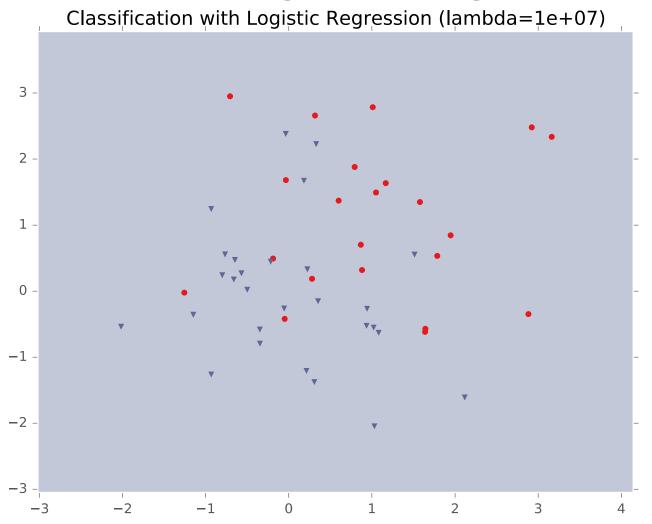


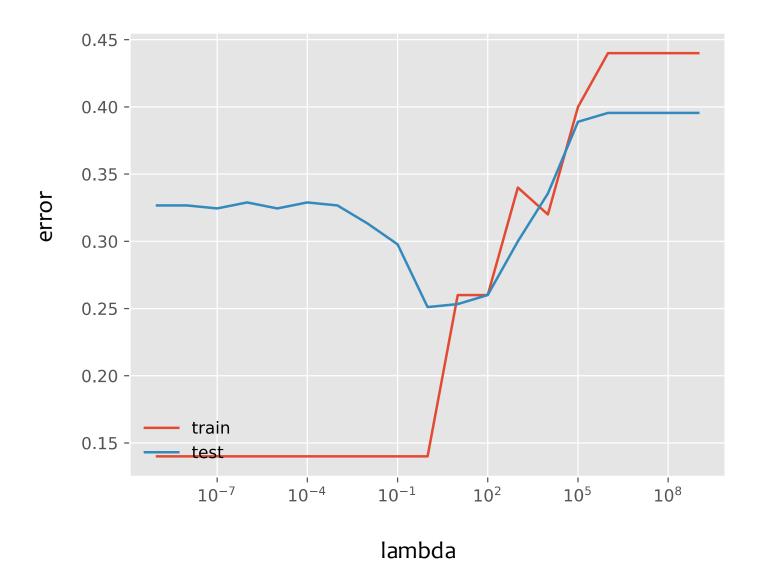












#### Regularization

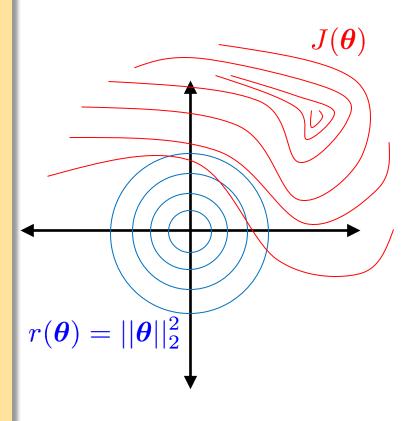
#### Poll Question 1:

Suppose we are minimizing  $J'(\theta)$  where

$$J'(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \lambda r(\boldsymbol{\theta})$$

As  $\lambda$  increases, the minimum of J'( $\theta$ ) will...

- A. ... move towards the midpoint between  $J(\theta)$  and  $r(\theta)$
- B. ... move towards the minimum of  $J(\theta)$
- C. ... move towards the minimum of  $r(\theta)$
- D. ... move towards a theta vector of positive infinities
- E. ... move towards a theta vector of negative infinities
- F. ... stay the same (TOXIC)



#### Regularization: Best Practices

#### Don't Regularize the Bias (Intercept) Parameter!

- In our models so far, the bias / intercept parameter is usually denoted by  $\theta_0$  that is, the parameter for which we fixed  $x_0=1$
- Regularizers always avoid penalizing this bias / intercept parameter
- Why? Because otherwise the learning algorithms wouldn't be invariant to a shift in the y-values

#### **Standardizing Data**

- It's common to standardize each feature by subtracting its mean and dividing by its standard deviation
- For regularization, this helps all the features be penalized in the same units (e.g. convert both centimeters and kilometers to z-scores)

#### Takeaways

- 1. Nonlinear basis functions allow linear models (e.g. Linear Regression, Logistic Regression) to capture nonlinear aspects of the original input
- Nonlinear features require no changes to the model (i.e. just preprocessing)
- 3. Regularization helps to avoid overfitting

#### Feature Engineering / Regularization Objectives

#### You should be able to...

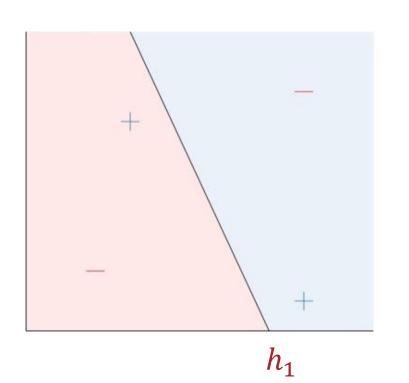
- Engineer appropriate features for a new task
- Use feature selection techniques to identify and remove irrelevant features
- Identify when a model is overfitting
- Add a regularizer to an existing objective in order to combat overfitting
- Explain why we should not regularize the bias term
- Convert linearly inseparable dataset to a linearly separable dataset in higher dimensions
- Describe feature engineering in common application areas

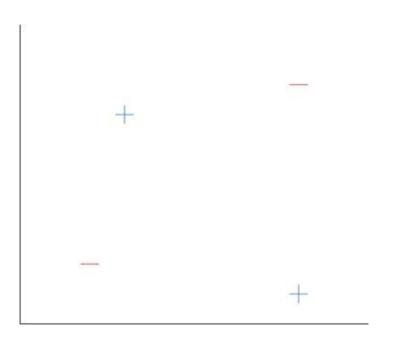
2/17/25

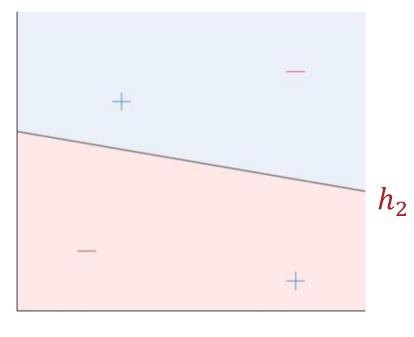
26

## Biological Neural Network



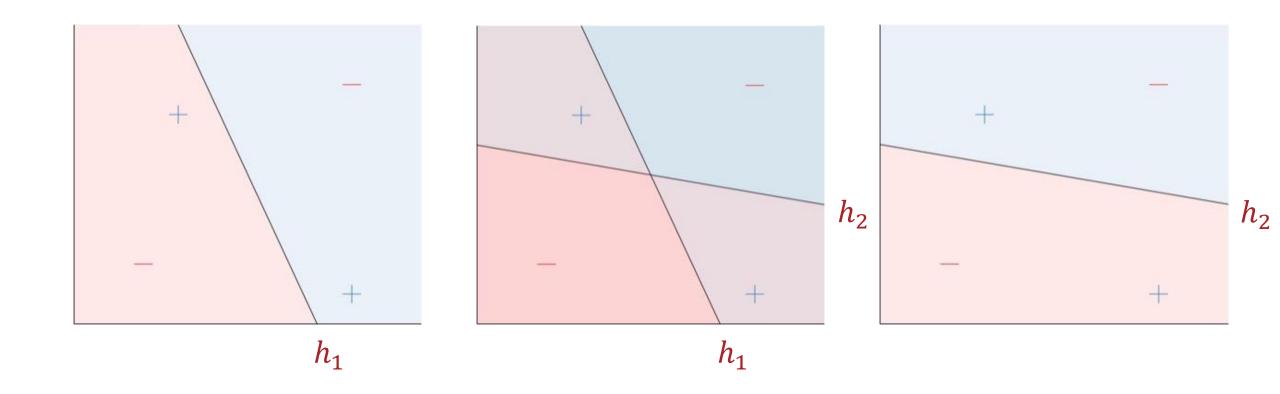




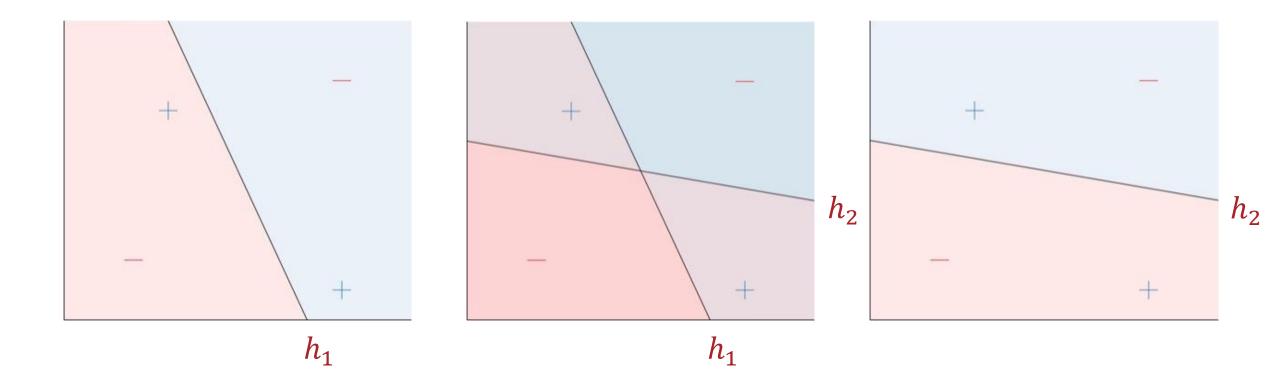


## Perceptrons $h(x) = sign(w^T x)$

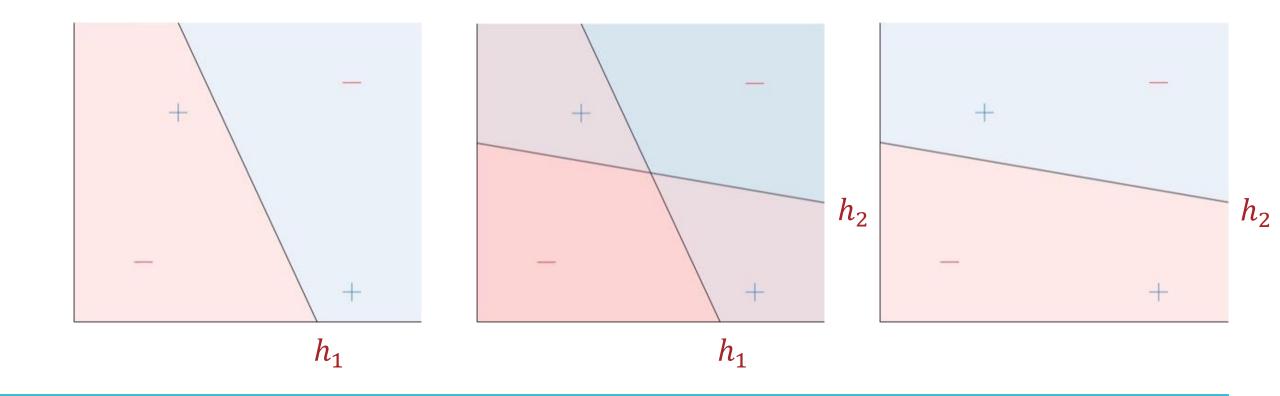
- Linear model for classification
- Predictions are +1 or -1



## **Combining Perceptrons**



$$h(x) = \begin{cases} +1 \text{ if } (h_1(x) = +1 \text{ and } h_2(x) = -1) \text{ or } (h_1(x) = -1 \text{ and } h_2(x) = +1) \\ -1 \text{ otherwise} \end{cases}$$



$$h(\mathbf{x}) = OR\left(AND(h_1(\mathbf{x}), \neg h_2(\mathbf{x})), AND(\neg h_1(\mathbf{x}), h_2(\mathbf{x}))\right)$$

#### Boolean Algebra

- Boolean variables are either +1 ("true") or -1 ("false")
- Basic Boolean operations:
  - Negation:  $\neg z = -1 * z$

• And: 
$$AND(z_1, z_2) = \begin{cases} +1 \text{ if both } z_1 \text{ and } z_2 \text{ equal} + 1 \\ -1 \text{ otherwise} \end{cases}$$

• Or: 
$$OR(z_1, z_2) = \begin{cases} +1 \text{ if either } z_1 \text{ or } z_2 \text{ equals } +1 \\ -1 \text{ otherwise} \end{cases}$$

#### Boolean Algebra

- Boolean variables are either +1 ("true") or -1 ("false")
- Basic Boolean operations
  - Negation:  $\neg z = -1 * z$

• And:  $AND(z_1, z_2) = sign(z_1 + z_2 - 1.5)$ 

• Or:  $OR(z_1, z_2) = sign(z_1 + z_2 + 1.5)$ 

#### Boolean Algebra

- Boolean variables are either +1 ("true") or -1 ("false")
- Basic Boolean operations
  - Negation:  $\neg z = -1 * z$

• And: 
$$AND(z_1, z_2) = \text{sign}\left( [-1.5, 1, 1] \begin{bmatrix} 1 \\ z_1 \\ z_2 \end{bmatrix} \right)$$

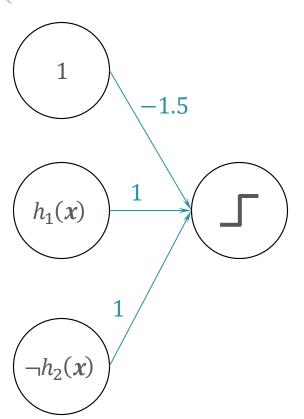
• Or: 
$$OR(z_1, z_2) = sign\left( [1.5, 1, 1] \begin{bmatrix} 1 \\ z_1 \\ z_2 \end{bmatrix} \right)$$

$$h(\mathbf{x}) = OR\left(AND(h_1(\mathbf{x}), \neg h_2(\mathbf{x})), AND(\neg h_1(\mathbf{x}), h_2(\mathbf{x}))\right)$$

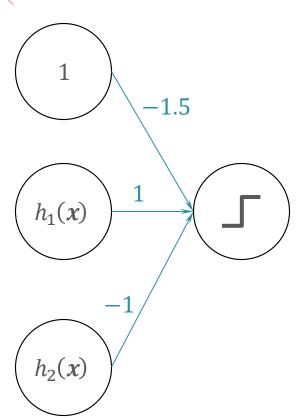
# Building a Network

#### Building a Network

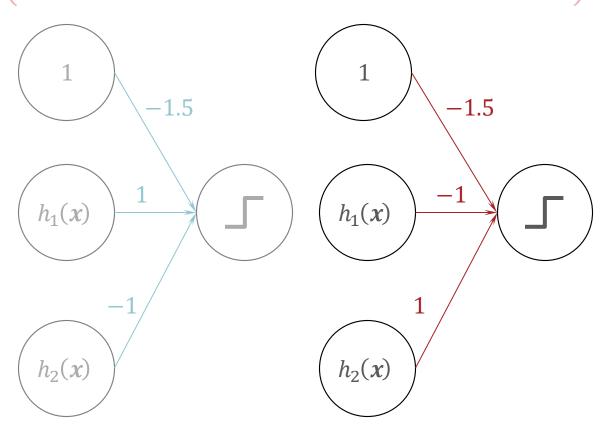
$$h(\mathbf{x}) = OR\left(AND(h_1(\mathbf{x}), \neg h_2(\mathbf{x})), AND(\neg h_1(\mathbf{x}), h_2(\mathbf{x}))\right)$$



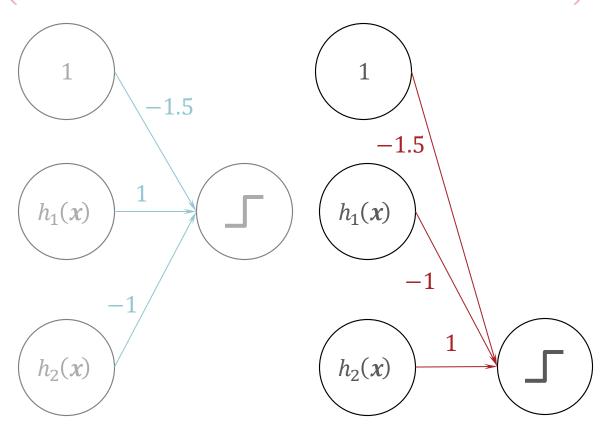
$$h(\mathbf{x}) = OR\left(AND(h_1(\mathbf{x}), \neg h_2(\mathbf{x})), AND(\neg h_1(\mathbf{x}), h_2(\mathbf{x}))\right)$$



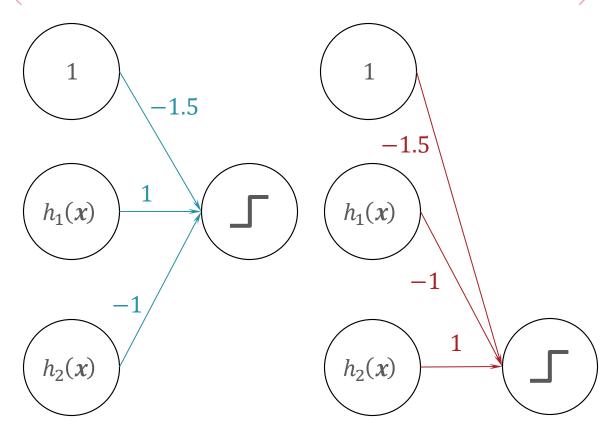
$$h(\mathbf{x}) = OR\left(AND(h_1(\mathbf{x}), \neg h_2(\mathbf{x})), AND(\neg h_1(\mathbf{x}), h_2(\mathbf{x}))\right)$$

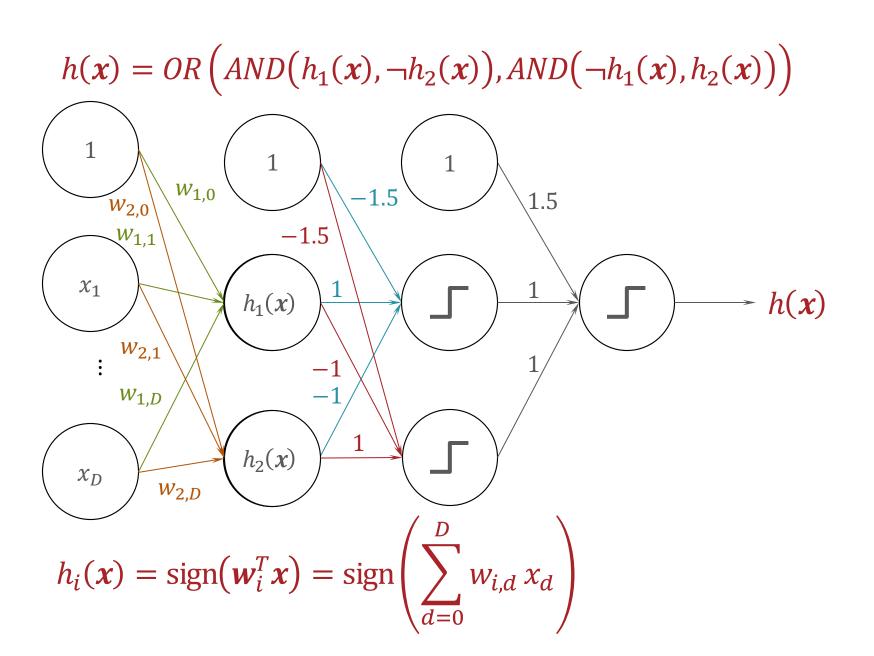


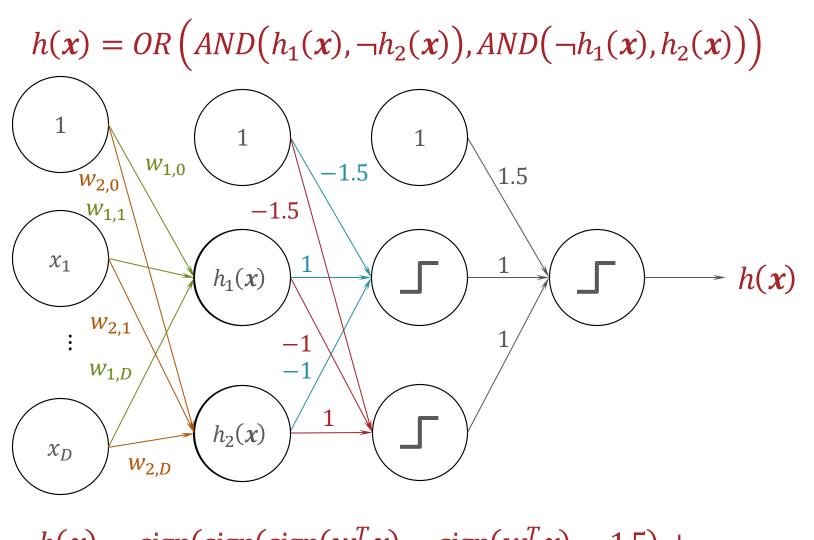
$$h(\mathbf{x}) = OR\left(AND(h_1(\mathbf{x}), \neg h_2(\mathbf{x})), AND(\neg h_1(\mathbf{x}), h_2(\mathbf{x}))\right)$$



$$h(\mathbf{x}) = OR\left(AND(h_1(\mathbf{x}), \neg h_2(\mathbf{x})), AND(\neg h_1(\mathbf{x}), h_2(\mathbf{x}))\right)$$







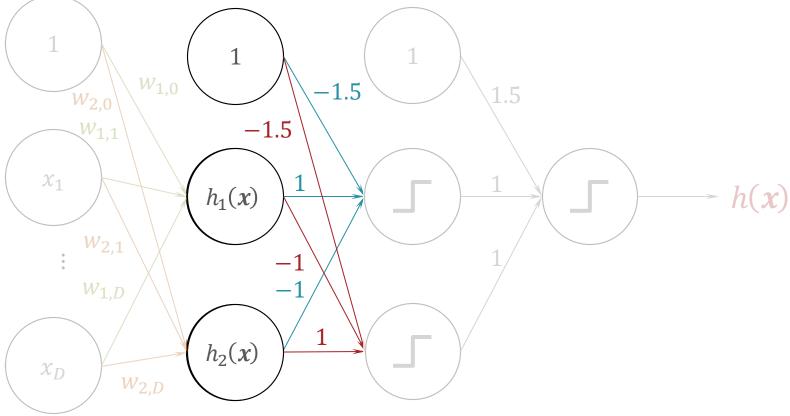
$$h(\mathbf{x}) = \operatorname{sign}(\operatorname{sign}(\mathbf{w}_1^T \mathbf{x}) - \operatorname{sign}(\mathbf{w}_2^T \mathbf{x}) - 1.5) + \\ \operatorname{sign}(-\operatorname{sign}(\mathbf{w}_1^T \mathbf{x}) + \operatorname{sign}(\mathbf{w}_2^T \mathbf{x}) - 1.5) + 1.5)$$

$$h(x) = OR\left(AND(h_1(x), \neg h_2(x)), AND(\neg h_1(x), h_2(x))\right)$$

$$\downarrow w_{2,0} \qquad \downarrow w_{1,0} \qquad \downarrow 1 \qquad \downarrow 1$$

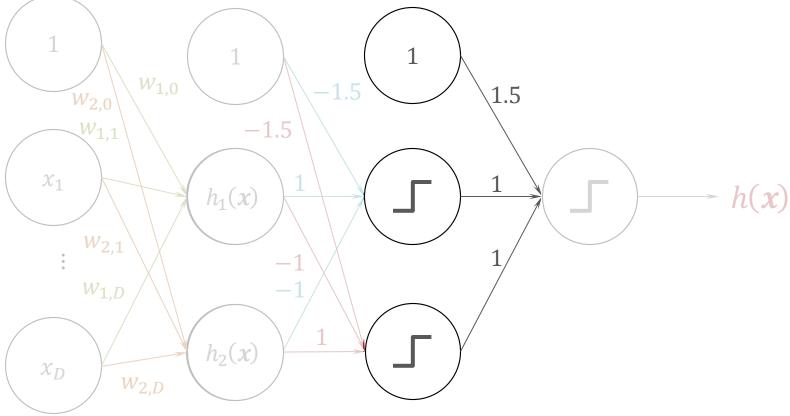
$$h(\mathbf{x}) = \operatorname{sign}(\operatorname{sign}(\mathbf{w}_1^T \mathbf{x}) - \operatorname{sign}(\mathbf{w}_2^T \mathbf{x}) - 1.5) + \\ \operatorname{sign}(-\operatorname{sign}(\mathbf{w}_1^T \mathbf{x}) + \operatorname{sign}(\mathbf{w}_2^T \mathbf{x}) - 1.5) + 1.5)$$

$$h(\mathbf{x}) = OR\left(AND(h_1(\mathbf{x}), \neg h_2(\mathbf{x})), AND(\neg h_1(\mathbf{x}), h_2(\mathbf{x}))\right)$$

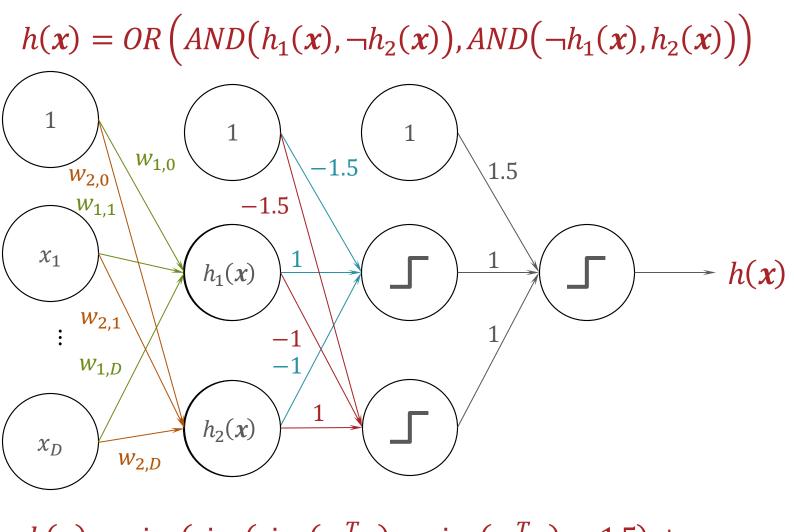


$$h(x) = \operatorname{sign}(\operatorname{sign}(\mathbf{w}_1^T \mathbf{x}) - \operatorname{sign}(\mathbf{w}_2^T \mathbf{x}) - 1.5) + \\ \operatorname{sign}(-\operatorname{sign}(\mathbf{w}_1^T \mathbf{x}) + \operatorname{sign}(\mathbf{w}_2^T \mathbf{x}) - 1.5) + 1.5)$$

$$h(\mathbf{x}) = OR\left(AND(h_1(\mathbf{x}), \neg h_2(\mathbf{x})), AND(\neg h_1(\mathbf{x}), h_2(\mathbf{x}))\right)$$

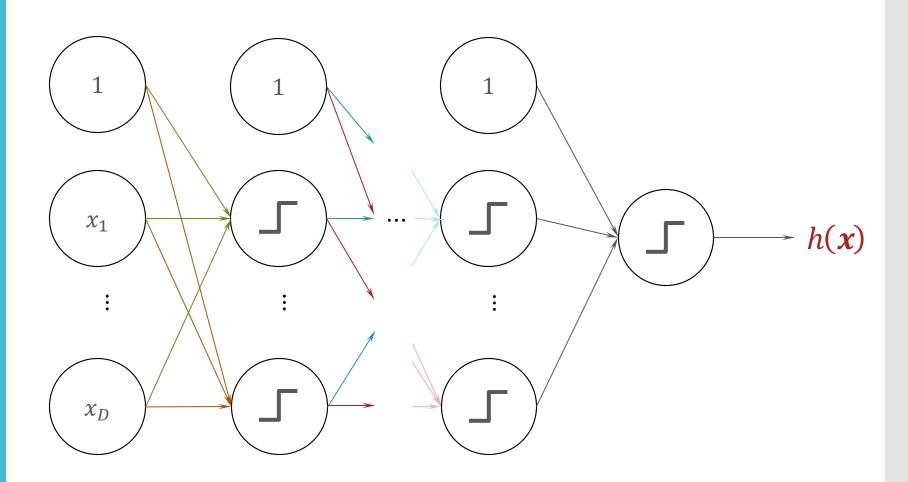


$$h(\mathbf{x}) = \operatorname{sign}(\operatorname{sign}(\mathbf{w}_1^T \mathbf{x}) - \operatorname{sign}(\mathbf{w}_2^T \mathbf{x}) - 1.5) + \\ \operatorname{sign}(-\operatorname{sign}(\mathbf{w}_1^T \mathbf{x}) + \operatorname{sign}(\mathbf{w}_2^T \mathbf{x}) - 1.5) + 1.5)$$

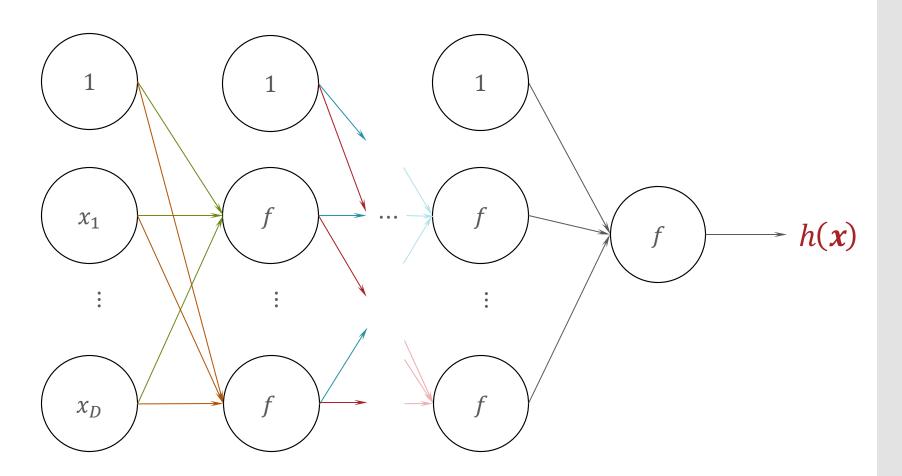


$$h(\mathbf{x}) = \operatorname{sign}(\operatorname{sign}(\mathbf{w}_1^T \mathbf{x}) - \operatorname{sign}(\mathbf{w}_2^T \mathbf{x}) - 1.5) + \\ \operatorname{sign}(-\operatorname{sign}(\mathbf{w}_1^T \mathbf{x}) + \operatorname{sign}(\mathbf{w}_2^T \mathbf{x}) - 1.5) + 1.5)$$

## Multi-Layer Perceptron (MLP)

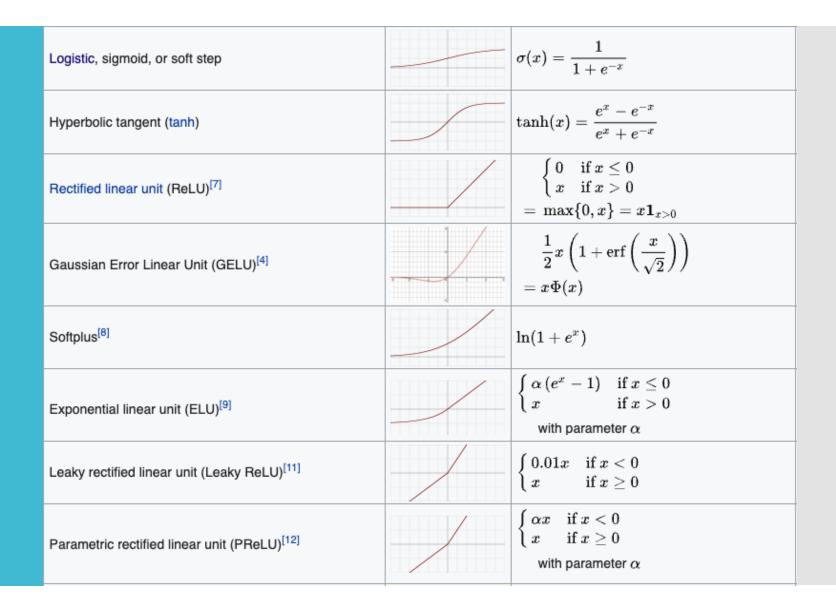


# (Fully-Connected) Feed Forward Neural Network



2/17/25

# Activation Functions



#### Poll Question 2

True or False: Linear and logistic regression models can be expressed as neural networks.

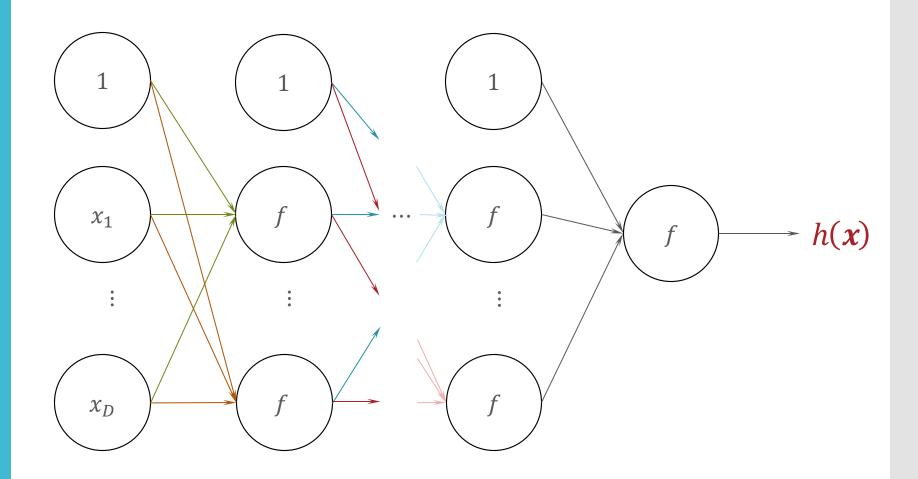
A. Only true for linear regression

B. Only true for logistic regression

C. TOXIC

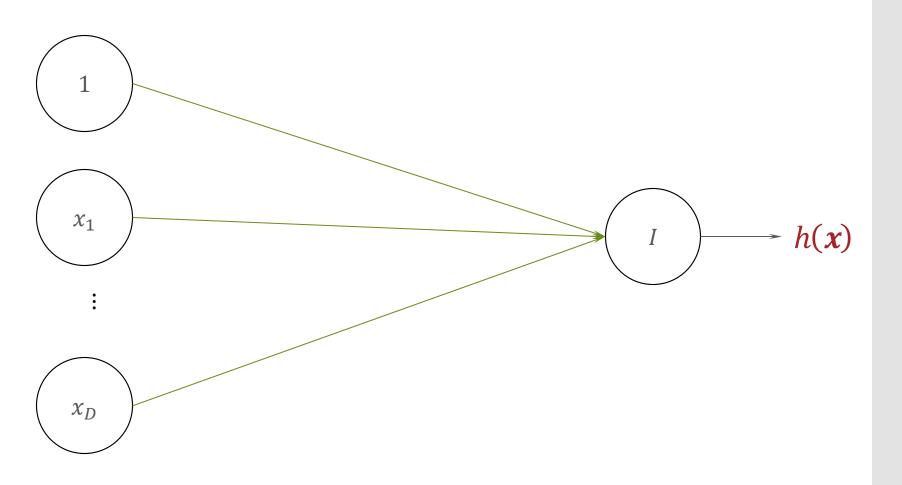
D. True for both

E. False for both

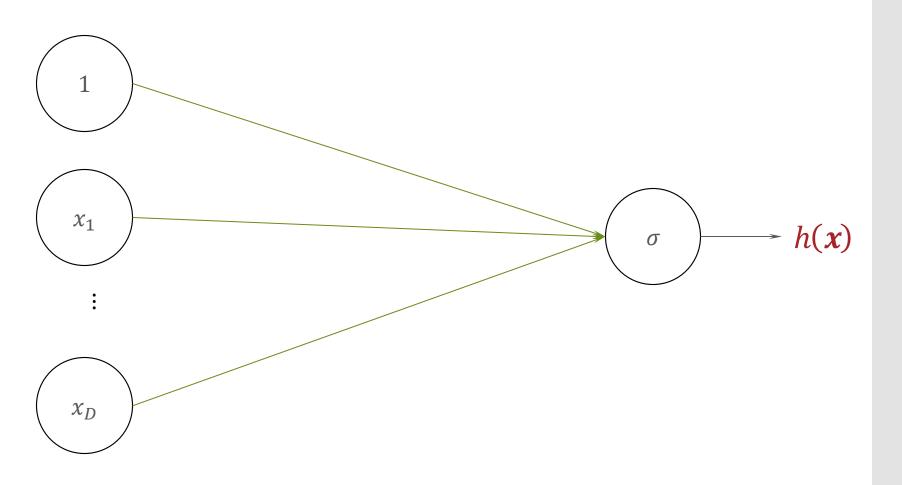


2/17/25

## Linear Regression as a Neural Network



## Logistic Regression as a Neural Network



2/17/25

## (Fully-Connected) Feed Forward Neural Network

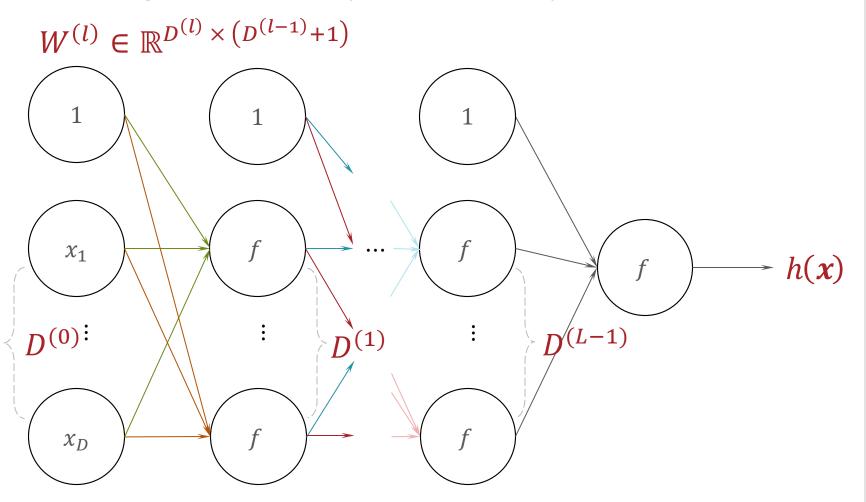
Input layer: Output layer: Hidden layers:  $l \in \{1, \dots, L-1\}$ l = Ll = 0 $x_1$ h(x) $D^{(0)}$ :  $x_D$ 

Layer l has dimension  $D^{(l)} \to \text{Layer } l$  has  $D^{(l)} + 1$  nodes, counting the bias node

2/17/25

## (Fully-Connected) Feed Forward Neural Network

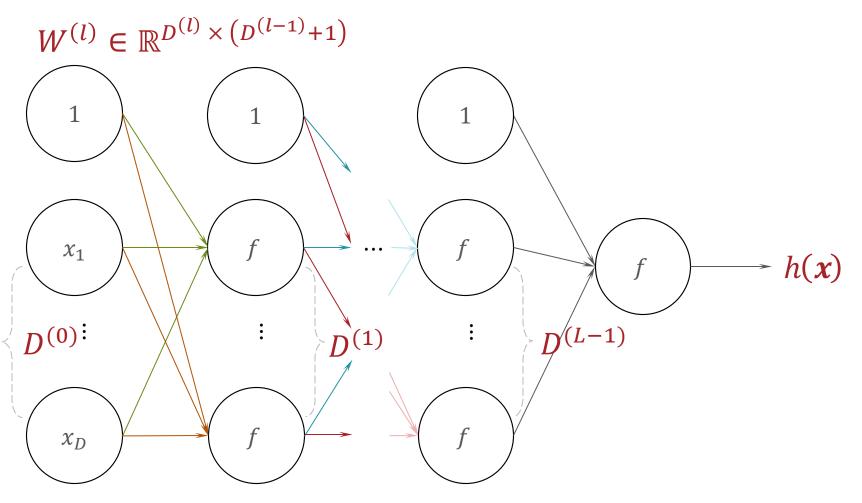
The weights between layer l-1 and layer l are a matrix:



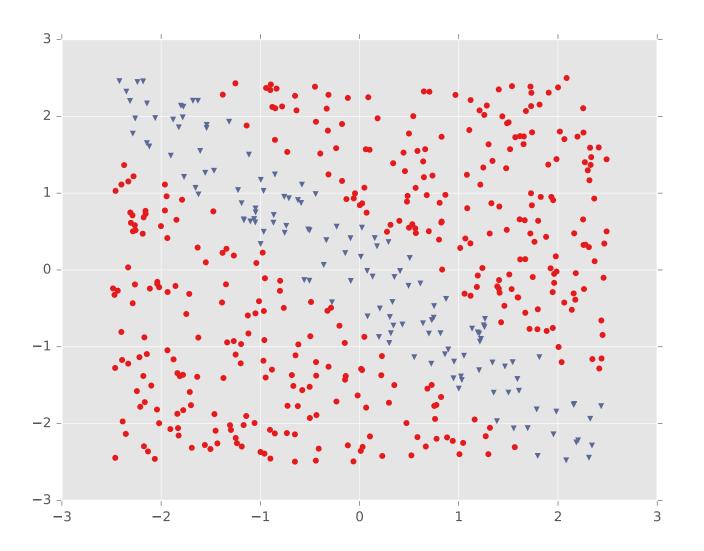
 $w_{j,i}^{(l)}$  is the weight between node i in layer l-1 and node j in layer l

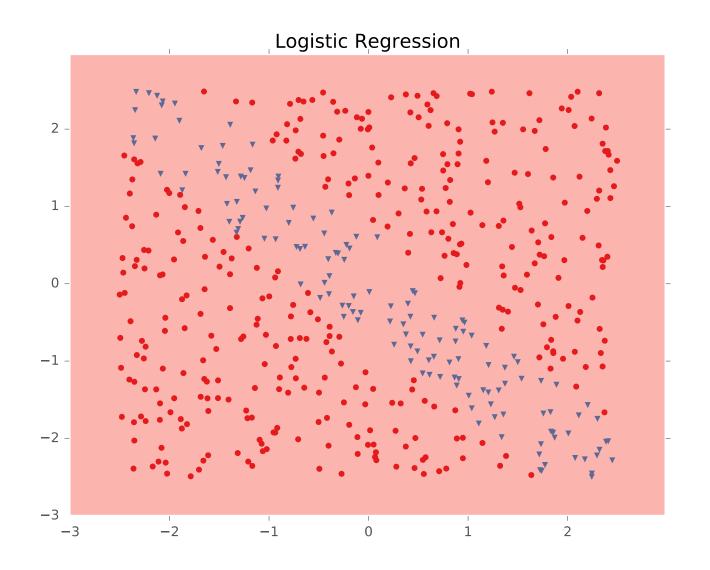
# So what are all these layers doing for us anyway?

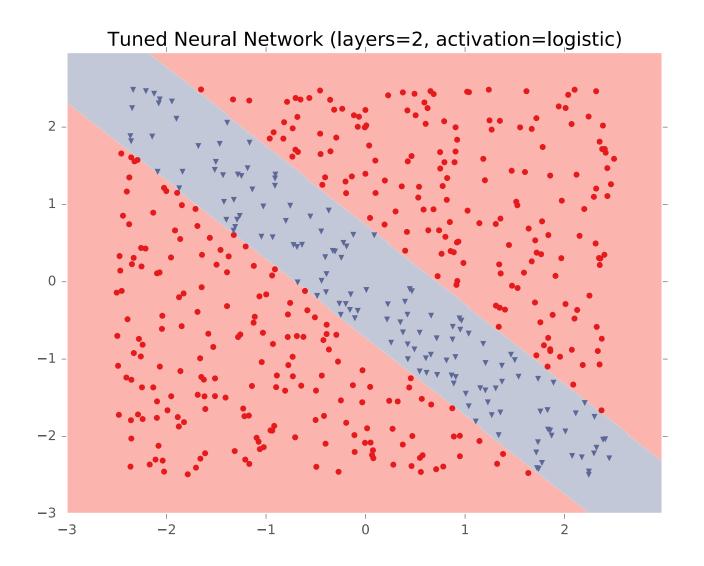
The weights between layer l-1 and layer l are a matrix:

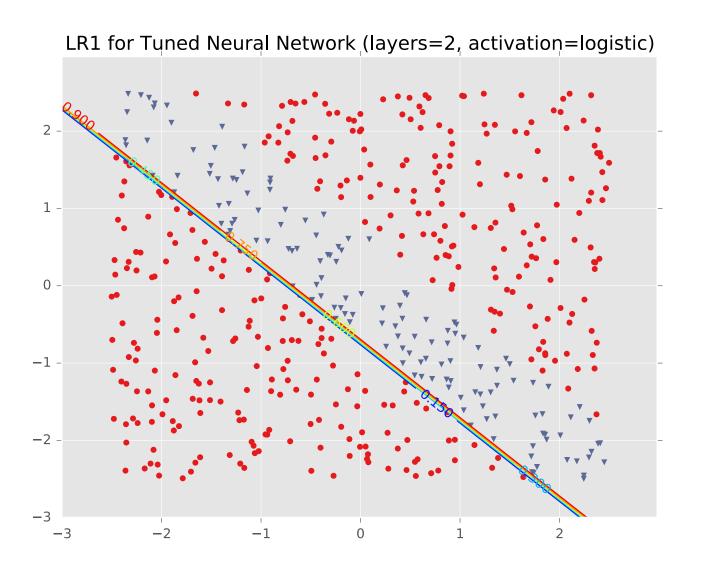


 $w_{j,i}^{(l)}$  is the weight between node i in layer l-1 and node j in layer l

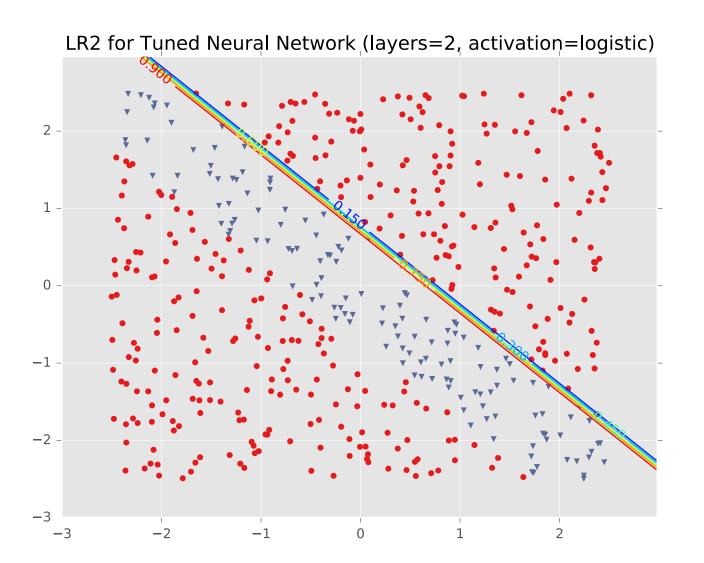


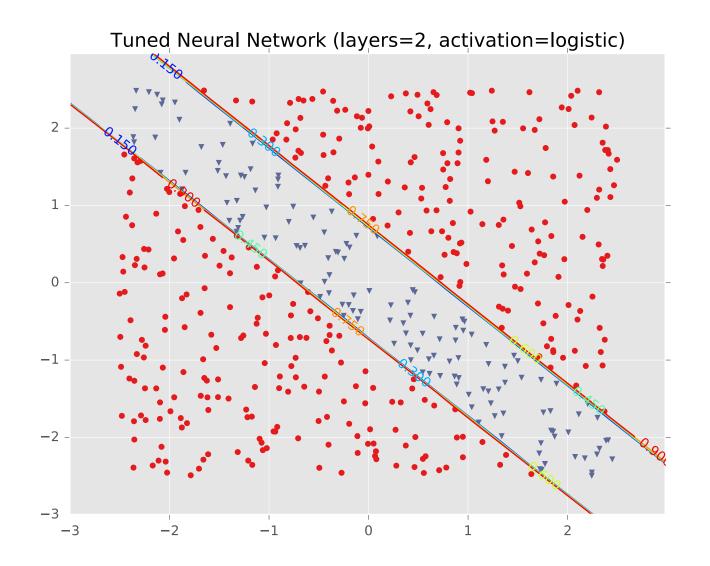


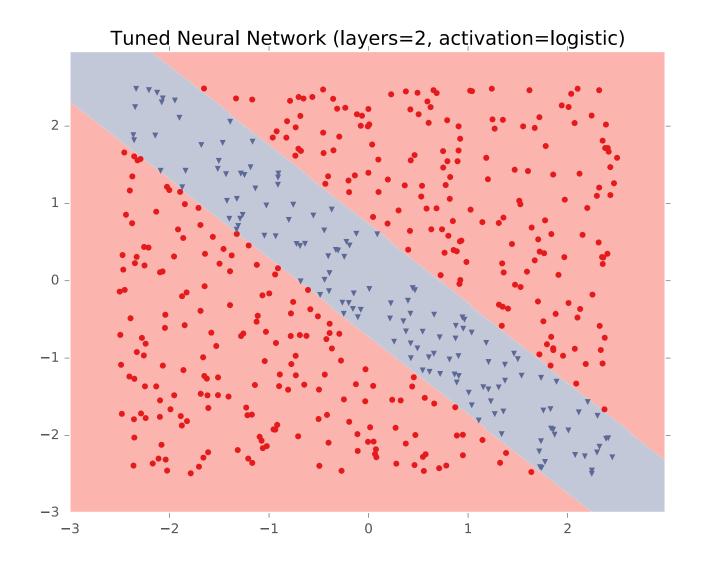


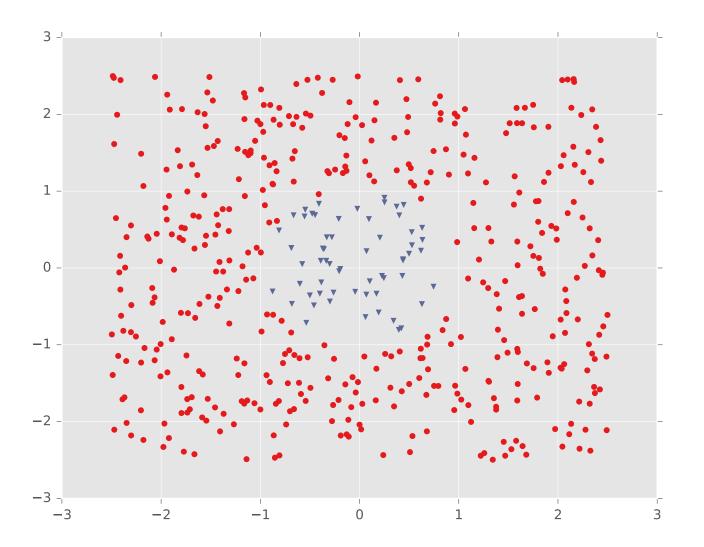


60

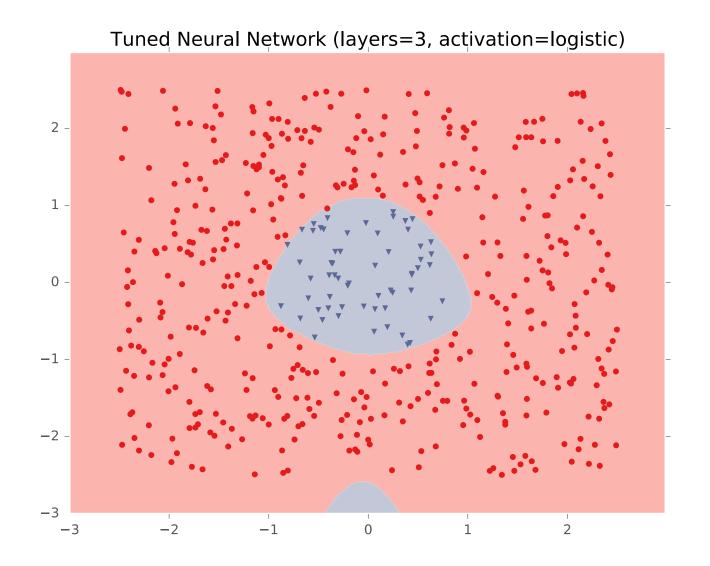


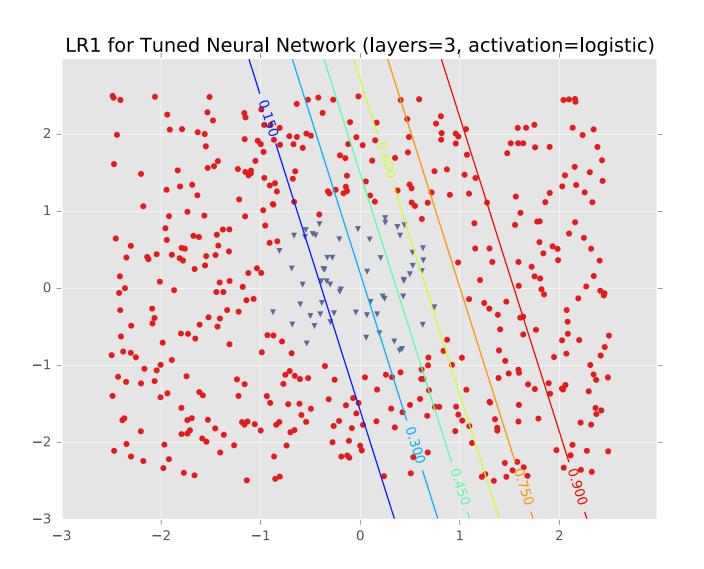


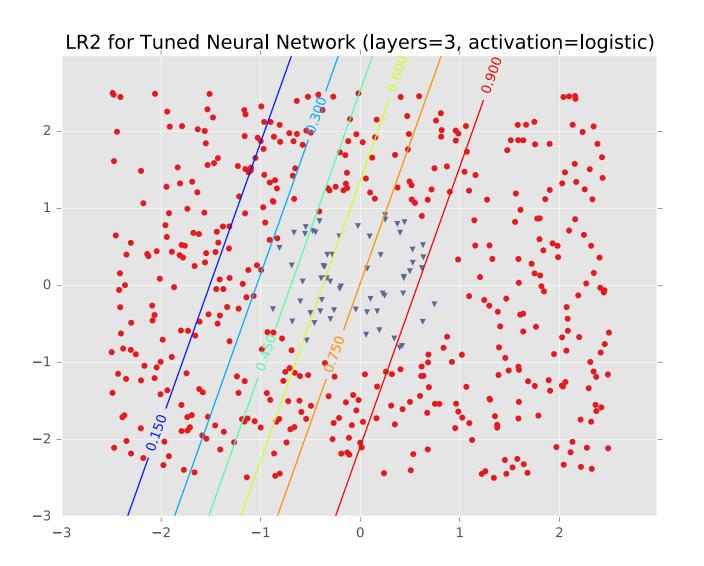


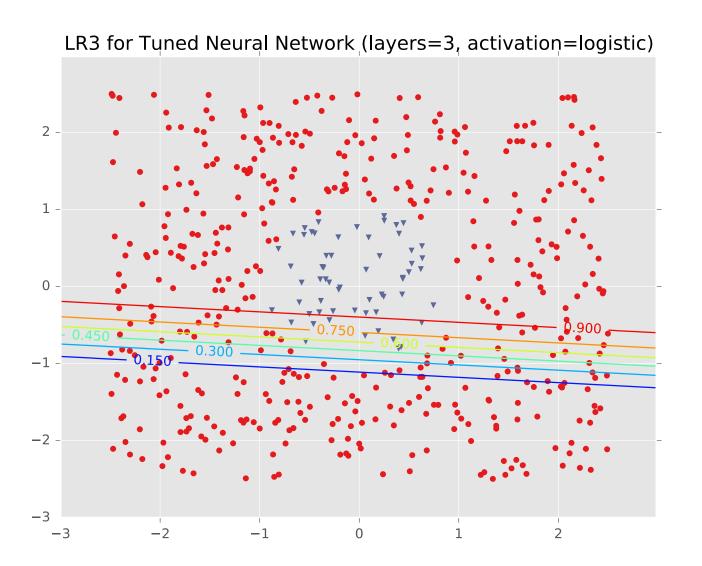


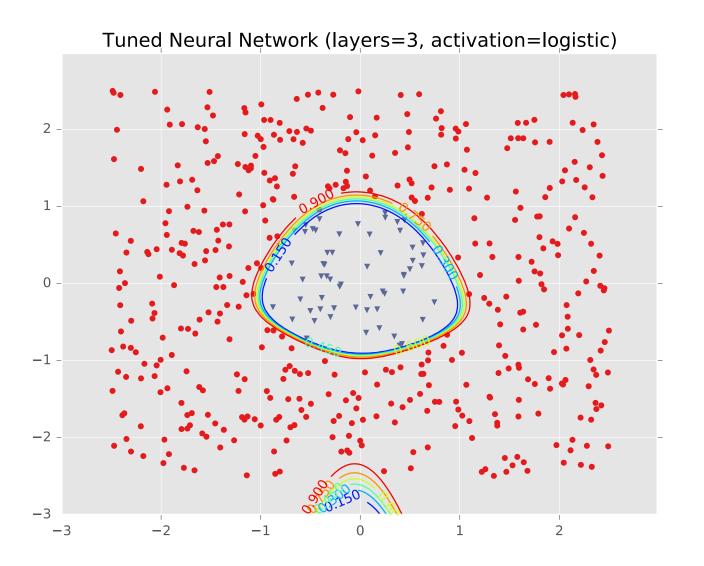


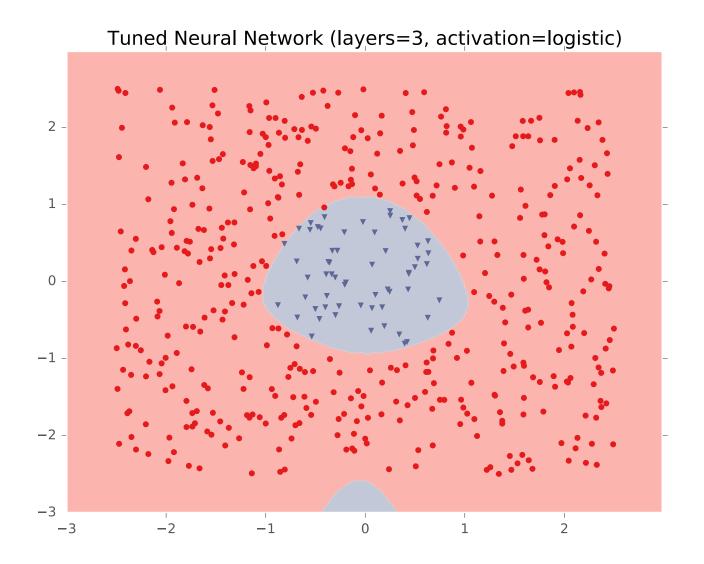






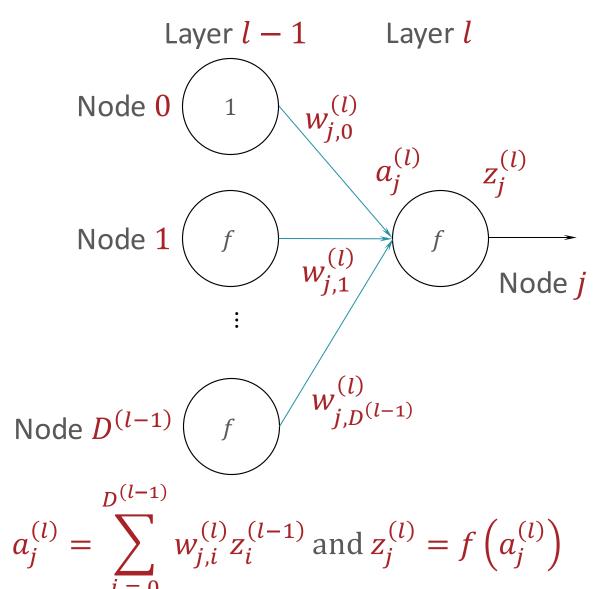






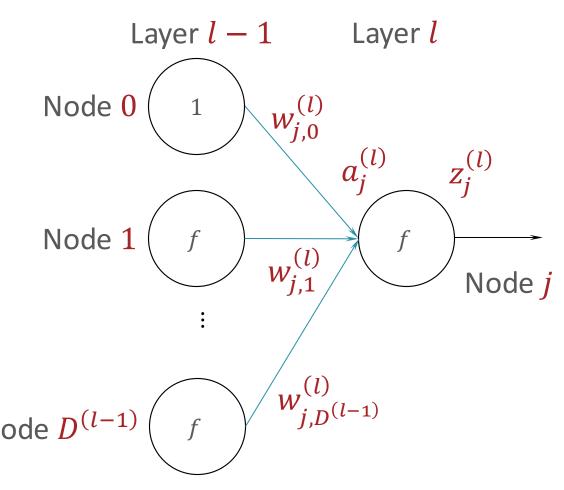
# Signal and Outputs

Every node has an incoming signal and outgoing output



# Signal and Outputs

Every node has an incoming signal and outgoing output



$$\mathbf{a}^{(l)} = W^{(l)} \mathbf{z}^{(l-1)} \text{ and } \mathbf{z}^{(l)} = [1, f(\mathbf{a}^{(l)})]^T$$

# Forward Propagation for Making Predictions

- Input: weights  $W^{(1)}$ , ...,  $W^{(L)}$  and a query data point  $\boldsymbol{x}$
- Initialize  $\mathbf{z}^{(0)} = [1, \mathbf{x}]^T$
- For l = 1, ..., L

• 
$$a^{(l)} = W^{(l)} z^{(l-1)}$$

• 
$$\mathbf{z}^{(l)} = \left[1, f(\mathbf{a}^{(l)})\right]^T$$

• Output:  $h_{W^{(1)},...,W^{(L)}}(x) = z^{(L)}$ 

# Okay, but where do these weights $W^{(1)}, ..., W^{(L)}$ come from?

- Input: weights  $W^{(1)}$ , ...,  $W^{(L)}$  and a query data point  $\boldsymbol{x}$
- Initialize  $\mathbf{z}^{(0)} = [1, \mathbf{x}]^T$
- For l = 1, ..., L

• 
$$a^{(l)} = W^{(l)}z^{(l-1)}$$

$$\cdot \mathbf{z}^{(l)} = \left[1, f(\mathbf{a}^{(l)})\right]^T$$

• Output:  $h_{W^{(1)},...,W^{(L)}}(x) = z^{(L)}$ 

## Neural Network Learning Objectives

You should be able to...

- 1. Explain the biological motivations for a neural network
- 2. Combine simpler models (e.g. linear regression, binary logistic regression, multinomial logistic regression) as components to build up feed-forward neural network architectures
- 3. Explain the reasons why a neural network can model nonlinear decision boundaries for classification
- Compare and contrast feature engineering with learning features
- 5. Identify (some of) the options available when designing the architecture of a neural network

76

6. Implement a feed-forward neural network

2/17/25