CS 333: NLP

Fall 2025

Prof. Carolyn Anderson Wellesley College

#### Reminders

- + HW 6 will be released today:
  - Neural network classifier for song lyrics using contextualized word embeddings
- No class on Tuesday (Tanner)
- CS Colloquium on 10/31
- Quiz 7 next Friday (J&M Chapter 7)
- My next help hours: Monday 4-5:30



Date: Oct 31st, 12:45 - 2:00pm Location: Sci H-105 **RSVP For Lunch** 

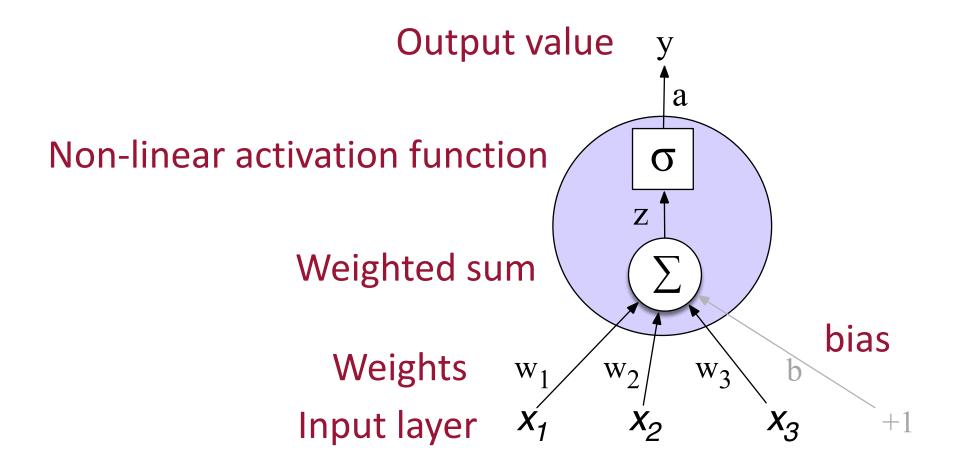


https://bit.ly/CSKawakami

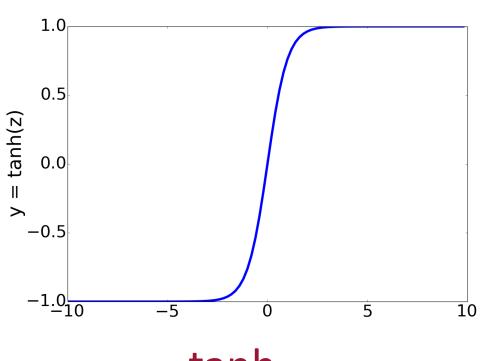
Accessibility and Disability: accessibility@wellesley.edu

?sb129@wellesley.edu

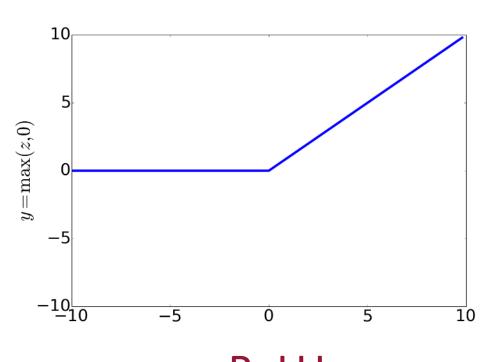
#### **Neural Network Unit**



#### Non-Linear Activation Functions besides sigmoid



#### **Most Common:**



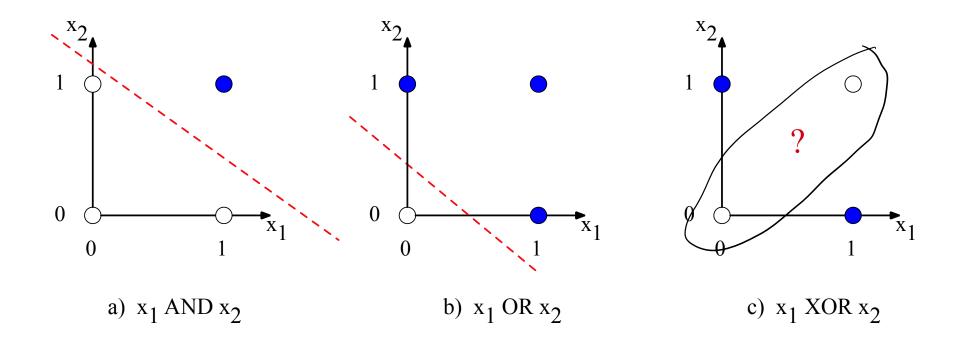
tanh

$$y = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

#### ReLU Rectified Linear Unit

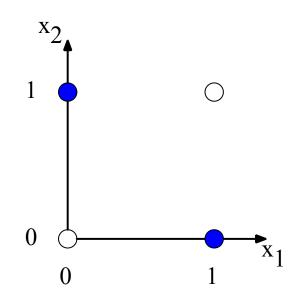
$$y = max(z, 0)$$

#### Decision boundaries

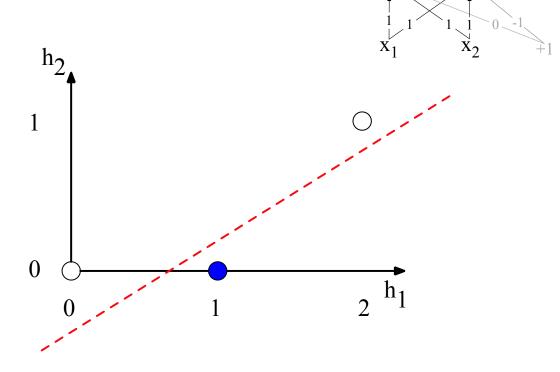


XOR is not a **linearly separable** function!

The hidden representation h



a) The original x space



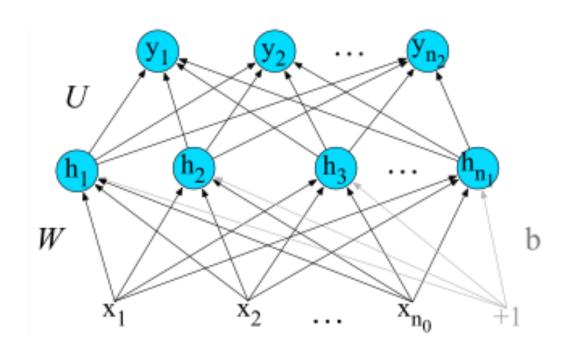
b) The new (linearly separable) h space

(With learning: hidden layers will learn to form useful representations)

## Feedforward Networks

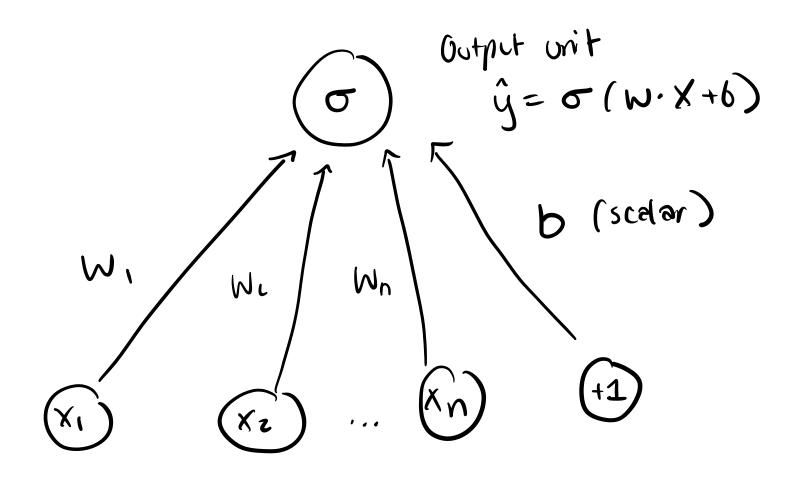
#### Feedforward Neural Networks

Can also be called **multi-layer perceptrons** for historical reasons



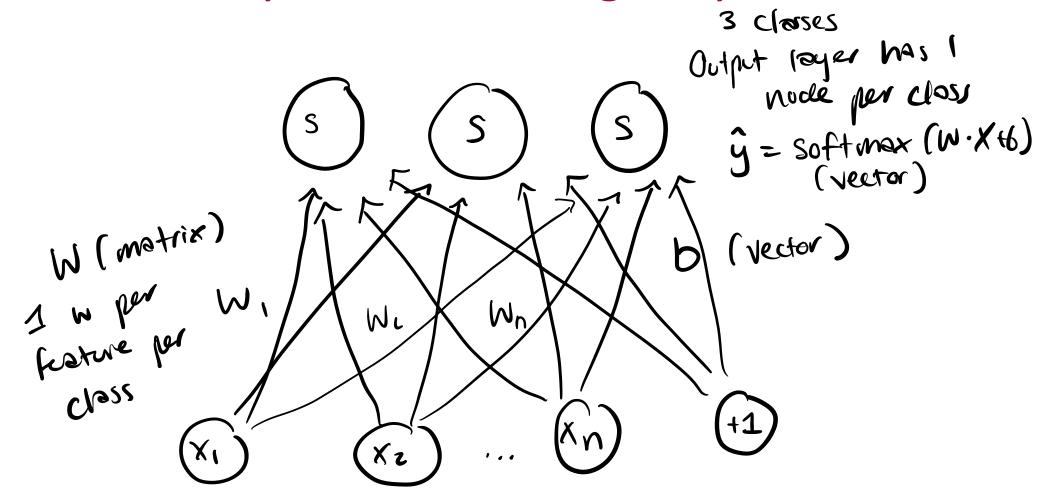
#### Binary Logistic Regression as a 1-layer Network

(we don't count the input layer in counting layers!)



#### Multinomial Logistic Regression as a 1-layer Network

#### Fully connected single layer network



#### Reminder: softmax: a generalization of sigmoid

For a vector z of dimensionality k, the softmax is:

$$\operatorname{softmax}(z) = \left[ \frac{\exp(z_1)}{\sum_{i=1}^k \exp(z_i)}, \frac{\exp(z_2)}{\sum_{i=1}^k \exp(z_i)}, \dots, \frac{\exp(z_k)}{\sum_{i=1}^k \exp(z_i)} \right]$$

$$\operatorname{softmax}(\mathbf{z}_i) = \frac{\exp(\mathbf{z}_i)}{\sum_{i=1}^d \exp(\mathbf{z}_i)} \quad 1 \le i \le d$$

Example:

$$\mathbf{z} = [0.6, 1.1, -1.5, 1.2, 3.2, -1.1],$$

softmax(
$$\mathbf{z}$$
) = [0.055, 0.090, 0.0067, 0.10, 0.74, 0.010].

### Replacing the bias unit

Let's switch to a notation without the bias unit Just a notational change

- 1. Add a dummy node  $a_0=1$  to each layer
- 2. Its weight  $w_0$  will be the bias
- 3. So input layer  $a^{[0]}_0=1$ ,
  - And  $a^{[1]}_0=1$ ,  $a^{[2]}_0=1$ ,...

## Replacing the bias unit

Instead of:

$$x = x_1, x_2, ..., x_{n0}$$

$$y = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$$

$$x = x_0, x_1, x_2, ..., x_{n0}$$

$$y = \sigma(\mathbf{W}\mathbf{x})$$

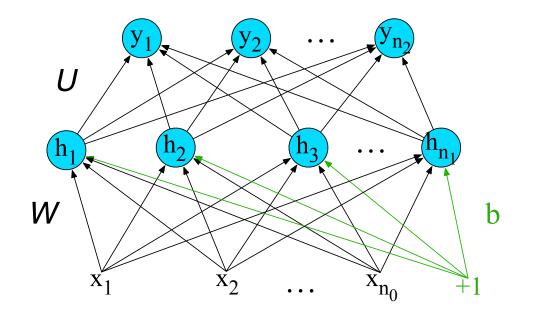
$$\mathbf{y}_j = \sigma \left( \sum_{i=1}^{n_0} \mathbf{W}_{ji} \mathbf{x}_i + \mathbf{b}_j \right)$$

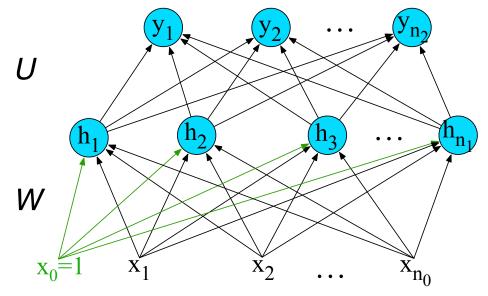
$$\sigma\left(\sum_{i=0}^{n_0}\mathbf{W}_{ji}\mathbf{x}_i\right)$$

## Replacing the bias unit

Instead of:

We'll do this:





# Using feedforward networks

#### Use cases for feedforward networks

Let's reconsider text classification

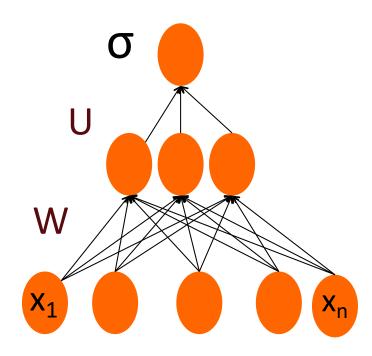
(State-of-the-art systems use more powerful architectures)

## Classification: Sentiment Analysis

We could do exactly what we did with logistic regression

Input layer are binary features as before

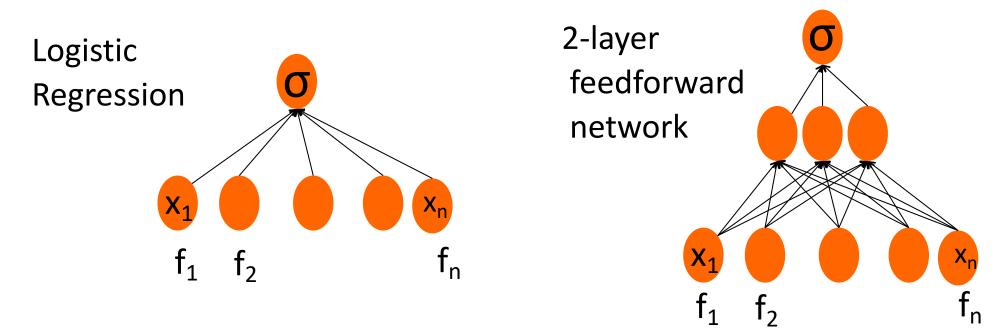
Output layer is 0 or 1



#### Sentiment Features

Var	Definition
$\overline{x_1}$	$count(positive lexicon words \in doc)$
$x_2$	count(negative lexicon words $\in$ doc)
<i>x</i> <sub>3</sub>	<pre> { 1 if "no" ∈ doc</pre>
$x_4$	$count(1st and 2nd pronouns \in doc)$
<i>x</i> <sub>5</sub>	$\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
$x_6$	log(word count of doc)

#### Feedforward nets for simple classification



Just add a hidden layer to logistic regression

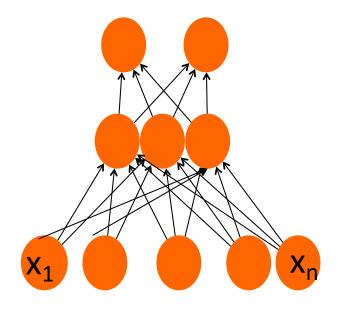
- allows the network to use non-linear interactions between features
- which may (or may not) improve performance.

#### Reminder: Multiclass Outputs

What if you have more than two output classes?

- Add more output units (one for each class)
- And use a "softmax layer"

$$\operatorname{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}} \quad 1 \le i \le D$$



## Training a Feedforward Network

#### Learning in Supervised Classification

#### Supervised classification:

- We know the correct label y (either 0 or 1) for each x.
- But what the system produces is an estimate,  $\hat{y}$ We want to set W to minimize the **distance** between our estimate  $\hat{y}^{(i)}$  and the true  $y^{(i)}$ .
- We need a distance estimator: a loss function or a cost function
- We need an optimization algorithm to update W to minimize the loss.

# Step 1: Compute Loss

1) Perine a loss function 
$$L(\Theta)$$
 to minimize.

We want:

- measure performance

- smooth / convexish

- differentiable

We call use cross-enterpy.

Today: well use Square loss:

 $L = \frac{1}{2}(y - \hat{y})^2$ 

## Hyperparameters

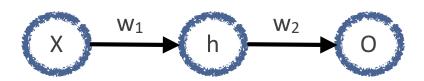
A hyperparameter is a variable in our model that is set beforehand rather than learned.

```
Learning Rate: how much do we grade weight at each step?

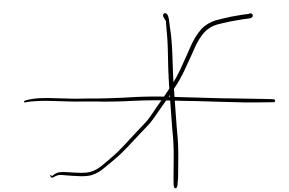
Batch Size: how often we godate weights?

Common congrumise: mini batening
```

## Single Neuron Feedforward Example



tanh is a non-linear function:



to between -1, 1

Inputs: 
$$(x,y)$$
 $h = tonh(w_1 X)$ 
 $0 = tonh(w_2 h)$ 

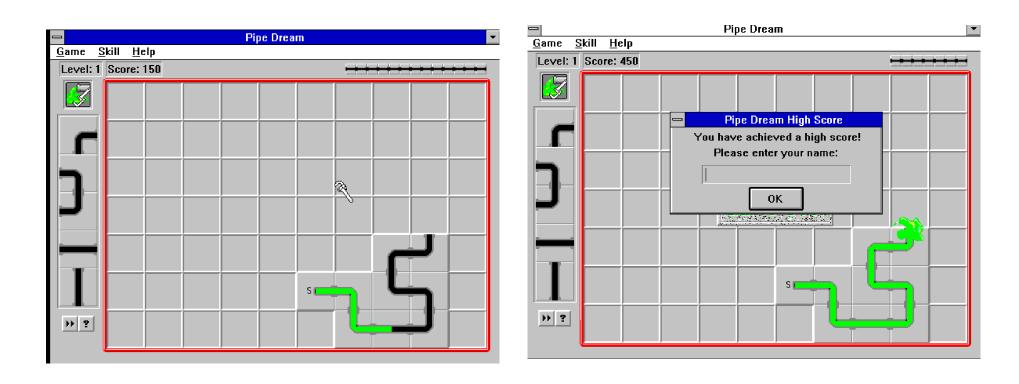
## Step 1: Compute Loss

Step 1: Compute Loss
$$L = \frac{1}{2} (y - \hat{y})^2$$
 Square Loss

#### Step 2: Compute the Gradient of the Loss

Derivative of tanh:

## Data Flows Through the Network



https://archive.org/details/win3\_PipeDr3x

# Step 3: Update the Weights

$$\frac{W_1}{W_2} + \frac{W_2}{W_2} + \frac{W_2}{W_2} = W_2 = W_2 = \frac{\partial L}{\partial w_2}$$

$$W_1 = W_1 = W_1 = W_2 =$$

# Incorporating Embeddings

#### Even better: representation learning

The real power of deep learning comes from the ability to **learn features** from the data, instead of using hand-built humanengineered features for classification.

#### Neural Net Classification with embeddings as input features!

