CS 333: NLP

Fall 2025

Prof. Carolyn Anderson Wellesley College

Reminders

- + HW 5 is due Thursday:
 - Regression with hand-crafted features to classify song lyrics by artist
- No class next Tuesday (Tanner)
- Quiz 7 next Friday (J&M Chapter 7)
- My next help hours: Thursday 4-5



Human and Computational Language Models Can Help Each Other - Focusing on Linguistic Testing in Korean

Dr. Sanghoun Song, Associate Professor of Linguistics and Director of the Research Institute of Language and Information at Korea University



Wednesday, October 22 4:30 - 5:30 pm Science Center - H401

Large Language Models: What can they tell us about language structure and acquisition

1-6pm on Nov. 6th George Sherman Union, Boston University

https://web.sas.upenn.edu/societyforlanguagedevelopment/symposium

Stochastic Gradient Descent Recap

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, ŷ

We want to set w and b to minimize the **distance** between our estimate $\mathfrak{g}^{(i)}$ 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 and b to minimize the loss.

Our goal: minimize the loss

For logistic regression, loss function is convex

- A convex function has just one minimum
- Gradient descent starting from any point is guaranteed to find the minimum
 - (Loss for neural networks is non-convex)

Gradients

The **gradient** of a function of many variables is a vector pointing in the direction of the greatest increase in a function.

Gradient Descent: Find the gradient of the loss function at the current point and move in the **opposite** direction.

Cross-Entropy Loss

$$L_{CE}(\hat{y}, y) = -[y \log \sigma(w \cdot x + b) + (1 - y) \log (1 - \sigma(w \cdot x + b))]$$

Weight Update
$$w^{t+1} = w^t - \eta \frac{d}{dw} L(f(x; w), y)$$

Cross-Entropy $L_{CE}(\hat{y}, y) = -[y \log \sigma(w \cdot x + b) + (1 - y) \log (1 - \sigma(w \cdot x + b))]$ Loss

Weight Update
$$w^{t+1} = w^t - \eta \frac{d}{dw} L(f(x; w), y)$$

Derivative
$$\frac{d}{dw}L(f(x;w),y) = \frac{\partial}{\partial \mathbf{w}_j} - [y\log\sigma(\mathbf{w}\cdot\mathbf{x}+b) + (1-y)\log(1-\sigma(\mathbf{w}\cdot\mathbf{x}+b))]$$
of Loss

Cross-Entropy $L_{CE}(\hat{y}, y) = -[y \log \sigma(w \cdot x + b) + (1 - y) \log (1 - \sigma(w \cdot x + b))]$ Loss

Weight Update
$$w^{t+1} = w^t - \eta \frac{d}{dw} L(f(x; w), y)$$
 Chain Rule $\frac{df}{dx} = \frac{du}{dv} \cdot \frac{dv}{dx}$

Derivative of $\sigma(x)$ $\frac{d\sigma(z)}{dz} = \sigma(z)(1 - \sigma(z))$ Derivative of $\ln(x)$ $\frac{d}{dx} \ln(x) = \frac{1}{x}$

Derivative of Loss
$$\frac{d}{dw}L(f(x;w),y) = \frac{\partial}{\partial w_{j}} - [y\log\sigma(w \cdot x + b) + (1-y)\log(1-\sigma(w \cdot x + b))]$$

$$= -\frac{y}{\sigma(wx+b)} \frac{\partial \sigma(wx+b)}{\partial w_{j}} - \frac{\partial \sigma(wx+b)}{\partial w_{j}} + \frac{\partial \sigma(wx+b)}{\partial w_{j}}$$

$$= -\frac{y}{\sigma(wx+b)} \frac{\partial \sigma(wx+b)}{\partial w_{j}} - \frac{\partial \sigma(wx+b)}{\partial w_{j}}$$

$$= -\frac{y}{\sigma(wx+b)} \frac{\partial \sigma(wx+b)}{\partial w_{j}} - \frac{\partial \sigma(wx+b)}{\partial w_{j}}$$

$$= -\frac{y}{\sigma(wx+b)} \frac{\partial \sigma(wx+b)}{\partial w_{j}} - \frac{\partial \sigma(wx+b)}{\partial w_{j}}$$

Cross-Entropy $L_{CE}(\hat{y}, y) = -[y \log \sigma(w \cdot x + b) + (1 - y) \log (1 - \sigma(w \cdot x + b))]$ Loss

Weight Update
$$w^{t+1} = w^t - \eta \frac{d}{dw} L(f(x; w), y)$$
 Chain Rule $\frac{df}{dx} = \frac{du}{dv} \cdot \frac{dv}{dx}$
Derivative of $\sigma(x)$ $\frac{d\sigma(z)}{dz} = \sigma(z)(1 - \sigma(z))$ Derivative of $\ln(x)$ $\frac{d}{dx} \ln(x) = \frac{1}{x}$

Derivative of Loss
$$\frac{\partial}{\partial \mathbf{w}_{j}} L(f(x; \mathbf{w}), \mathbf{y}) = \frac{\partial}{\partial \mathbf{w}_{j}} - [y \log \sigma(\mathbf{w} \cdot \mathbf{x} + b) + (1 - y) \log (1 - \sigma(\mathbf{w} \cdot \mathbf{x} + b))]$$

$$= -\left[\frac{\partial}{\partial \mathbf{w}_{j}} y \log \sigma(\mathbf{w} \cdot \mathbf{x} + b) + \frac{\partial}{\partial \mathbf{w}_{j}} (1 - y) \log [1 - \sigma(\mathbf{w} \cdot \mathbf{x} + b)]\right]$$

$$= -\frac{y}{\sigma(\mathbf{w} \cdot \mathbf{x} + b)} \frac{\partial}{\partial \mathbf{w}_{j}} \sigma(\mathbf{w} \cdot \mathbf{x} + b) - \frac{1 - y}{1 - \sigma(\mathbf{w} \cdot \mathbf{x} + b)} \frac{\partial}{\partial \mathbf{w}_{j}} 1 - \sigma(\mathbf{w} \cdot \mathbf{x} + b)$$

$$= -\left[\frac{y}{\sigma(\mathbf{w} \cdot \mathbf{x} + b)} - \frac{1 - y}{1 - \sigma(\mathbf{w} \cdot \mathbf{x} + b)}\right] \frac{\partial}{\partial \mathbf{w}_{j}} \sigma(\mathbf{w} \cdot \mathbf{x} + b)$$

$$= -\left[\frac{y - \sigma(\mathbf{w} \cdot \mathbf{x} + b)}{\sigma(\mathbf{w} \cdot \mathbf{x} + b)[1 - \sigma(\mathbf{w} \cdot \mathbf{x} + b)]}\right] \sigma(\mathbf{w} \cdot \mathbf{x} + b) [1 - \sigma(\mathbf{w} \cdot \mathbf{x} + b)] \frac{\partial(\mathbf{w} \cdot \mathbf{x} + b)}{\partial \mathbf{w}_{j}}$$

$$= [\sigma(\mathbf{w} \cdot \mathbf{x} + b) - y]\mathbf{x}_{j}$$

Deriving cross-entropy loss for multi-label classification

Goal: maximize probability of the correct label p(y|x)

$$L_{\text{CE}}(\hat{\mathbf{y}}, \mathbf{y}) = -\sum_{k=1}^{K} \mathbf{y}_k \log \hat{\mathbf{y}}_k$$

$$= -\log \hat{p}(\mathbf{y}_c = 1 | \mathbf{x}) \quad \text{(where } c \text{ is the correct class)}$$

$$= -\log \frac{\exp(\mathbf{w}_c \cdot \mathbf{x} + b_c)}{\sum_{j=1}^{K} \exp(\mathbf{w}_j \cdot \mathbf{x} + b_j)} \quad (c \text{ is the correct class)}$$

How much do we move in that direction?

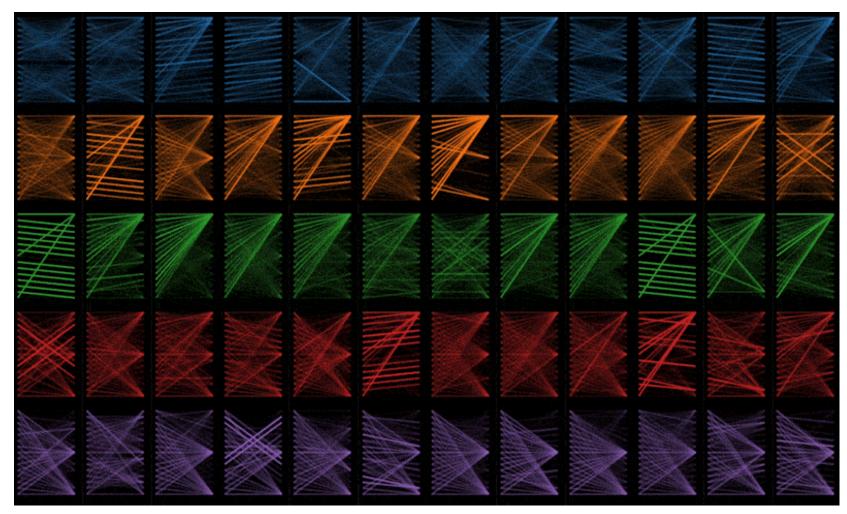
- The value of the gradient (slope in our example) $\frac{d}{dw}L(f(x;w),y)$ weighted by a learning rate η
- Higher learning rate means that we make bigger adjustments to the weights

$$w^{t+1} = w^t - \eta \frac{d}{dw} L(f(x; w), y)$$

www.i-am.ai/gradient-descent.html

Neural Networks

This is not your brain

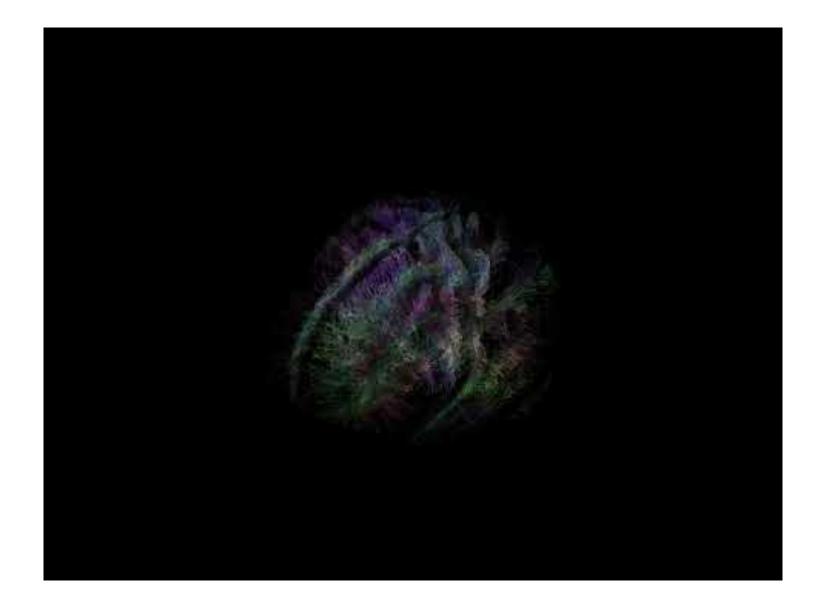


https://github.com/jessevig/bertviz

It's a large language model (neural network)

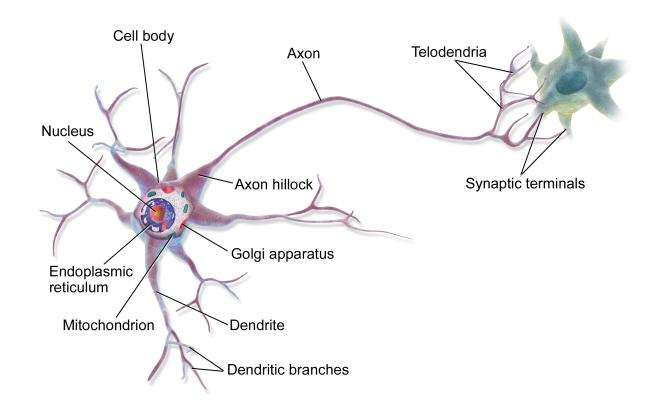
Most contemporary models have billions of parameters (GPT3: 175 billion); some may have trillions

This is someone's brain

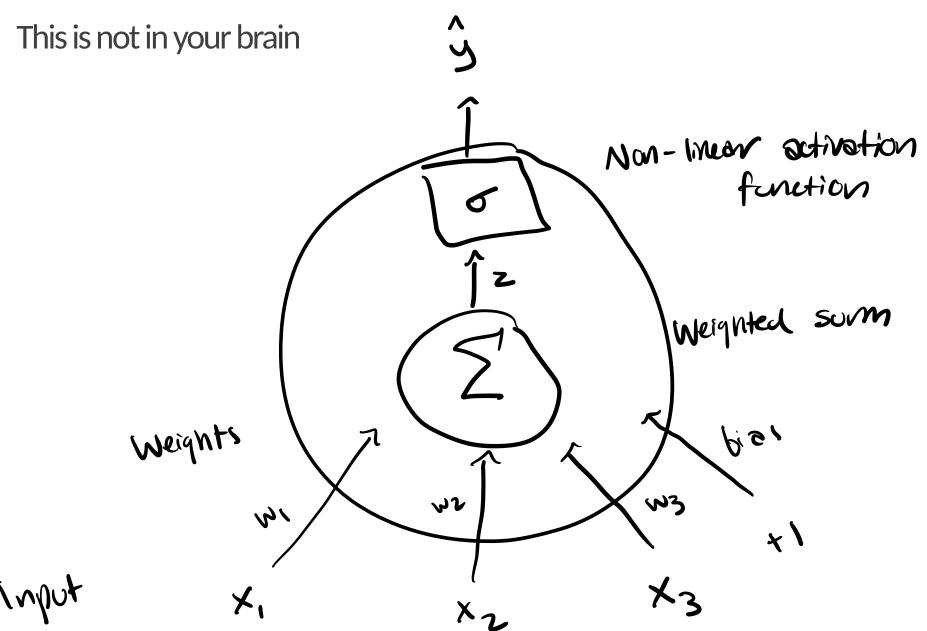


estimated to have 100 trillion synapses connecting 100 billion neurons

This is in your brain



Neural Network Unit



Units in Neural Networks

Neural unit

Takes a weighted sum of inputs

(provided Sum of inputs

$$Z = b + Ewix;$$
 $= w \cdot x + b$
 $G = a = f(z)$

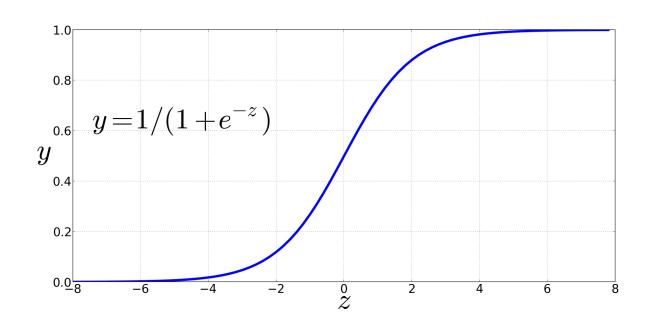
(could be sigmoid)

Non-Linear Activation Functions

We're already seen the sigmoid for logistic regression:

Sigmoid

$$\sigma(z) = \frac{1}{1 + e^z}$$



Final function of a unit:

Spot the differences

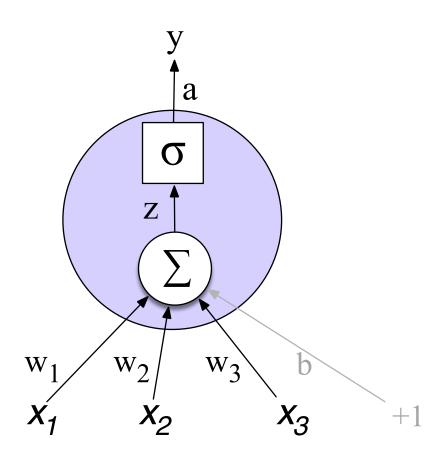
Neural Network Unit

Logistic Regression

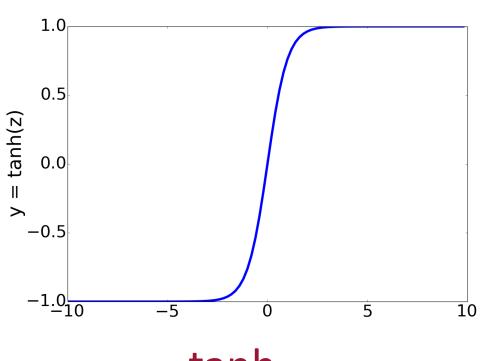
$$Z = \left(\sum_{i=1}^{n} w_{i} x_{i} \right) + b$$

$$p(y=1|x) = \sigma(w \cdot x + b)$$

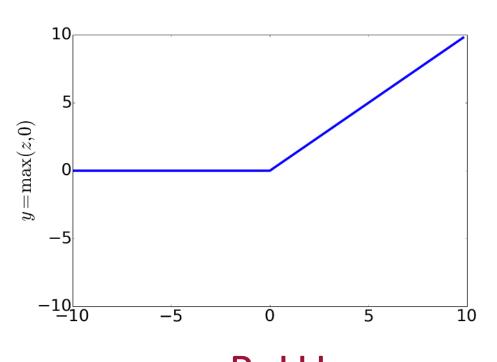
Final unit again



Non-Linear Activation Functions besides sigmoid



Most Common:



tanh

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

ReLU Rectified Linear Unit

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

Example: XOR

The XOR problem

Can neural units compute simple functions of input?

	AND			OR		2	XOR	
x 1	x2	у	x 1	_ x2	у	x 1	x2	у
		0	0			0		1
	1	0	0	1		0	1	1
	0	0	1	0	1		0	1
1	1	1	1	1	1	1	1	0

Perceptrons

A very simple neural unit

- Binary output (0 or 1)
- No non-linear activation function

$$y = \begin{cases} 0, & \text{if } \mathbf{w} \cdot \mathbf{x} + b \le 0 \\ 1, & \text{if } \mathbf{w} \cdot \mathbf{x} + b > 0 \end{cases}$$

Solving AND

Deriving AND

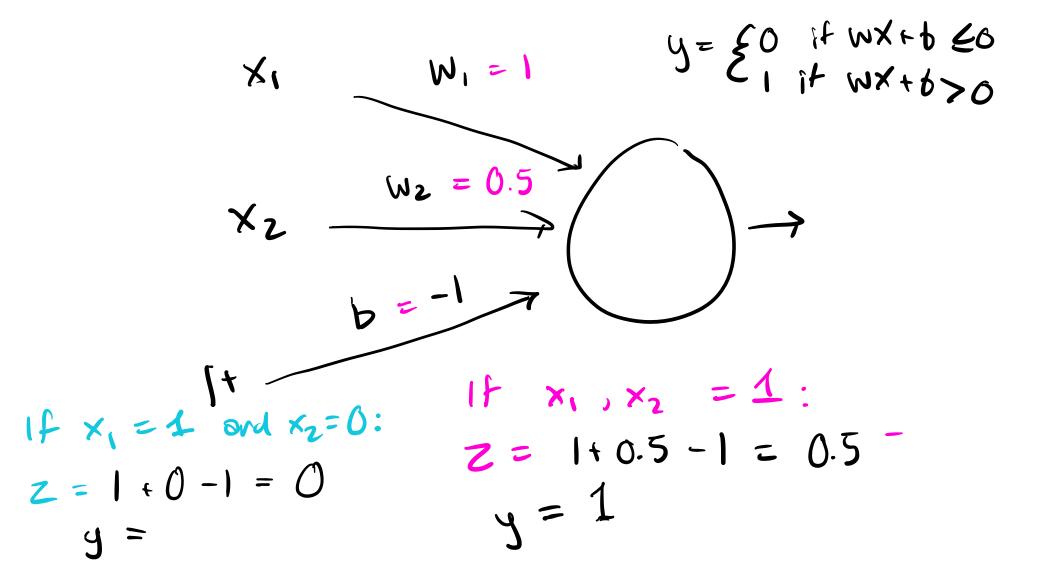
$$y = \begin{cases} 0, & \text{if } \mathbf{w} \cdot \mathbf{x} + b \le 0 \\ 1, & \text{if } \mathbf{w} \cdot \mathbf{x} + b > 0 \end{cases}$$

Goal: return 1 if x1 and x2 are 1

AND				
x 1	x 2	у		
0	0	0		
0	1	0		
1	0	0		
1	1	1		

Deriving AND

Goal: return 1 if x1 and x2 are 1



Exercise: solving OR

OR			
x 1	x 2	у	
0	0	0	
0	1	1	
1	0	1	
1	1	1	

Deriving OR

$$y = \begin{cases} 0, & \text{if } \mathbf{w} \cdot \mathbf{x} + b \le 0 \\ 1, & \text{if } \mathbf{w} \cdot \mathbf{x} + b > 0 \end{cases}$$

Goal: return 1 if either input is 1

OK				
x2	у			
0	0			
1	1			
0	1			
1	1			
	x2 0 1 0			

 ΔD

Deriving OR

-0.5

Goal: return 1 if either input is 1

solving XOR

	XOR	
x 1	x2	у
0	0	0
0	1	1
1	0	1
1	1	0

Trick question! It's not possible to capture XOR with perceptrons



Why? Perceptrons are linear classifiers

Perceptron equation is the equation of a line

$$w_1 x_1 + w_2 x_2 + b = 0$$

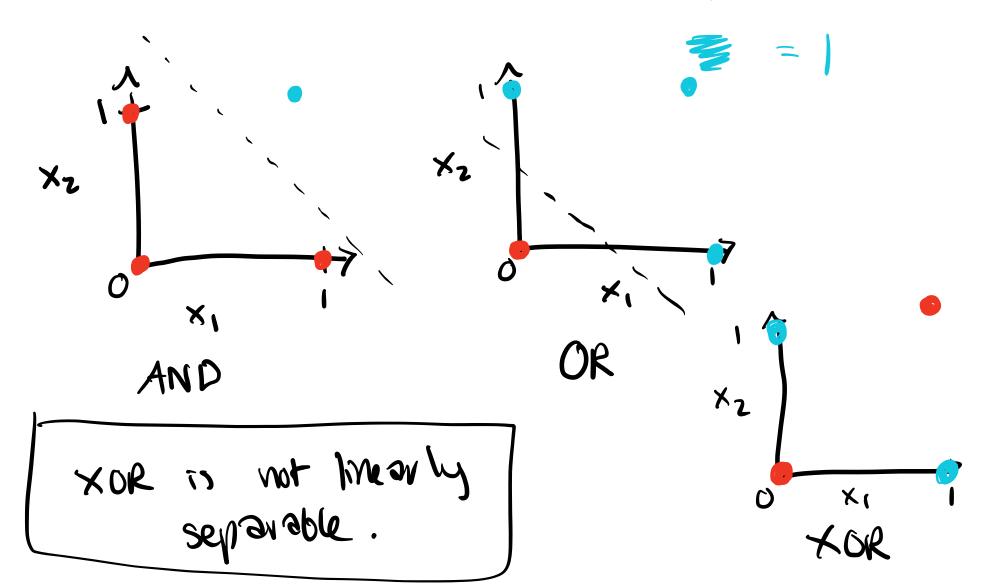
(in standard linear format: $x_2 = (-w_1/w_2)x_1 + (-b/w_2)$

This line acts as a decision boundary

- 0 if input is on one side of the line
- 1 if on the other side of the line

Decision boundaries





Solution to the XOR problem

XOR can't be calculated by a single perceptron XOR can be calculated by a layered network of units.

	XOR		ReLU y_1 -2 0
x 1	x2	у	ReLU h_1 h_2 $+1$
0	0	0	
0	1	1	\mathbf{X}_{1} \mathbf{X}_{2} $+1$
1	0	1	
1	1	0	