

# Neural Networks

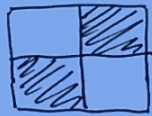
## CS506 Lecture (3/29/21)

\* ~~inherent~~ limitations make it very far from workings of brain activity

\* Approach NOT from biological perspective

↳ Recall: LOGISTIC REGRESSION

• Given a 2x2 grid → Find function,  $f$ , that can identify diagonal patterns (specifically, this one for now)



$$\blacksquare = c_1 = 1$$

$$\square = c_2 = -1$$

-1	1
1	-1

assume this is our only def. of diagonal

$$f\left(\begin{bmatrix} a_{00} & a_{01} \\ a_{10} & a_{11} \end{bmatrix}\right) = \begin{cases} \text{yes, if } \begin{bmatrix} a_{00} & a_{01} \\ a_{10} & a_{11} \end{bmatrix} = \begin{bmatrix} \blacksquare & \square \\ \square & \blacksquare \end{bmatrix} \\ \text{no otherwise} \end{cases}$$

→ assign weights to each cell of the grid

$$w_1 a_{00} + w_2 a_{01} + w_3 a_{10} + w_4 a_{11} =$$

$$\begin{cases} = b & \text{if diagonal found} \\ < b & \text{otherwise} \end{cases}$$

(if below some threshold, it's not a diagonal)

Ex]

$w_1$	$w_2$
$w_3$	$w_4$

-2	2
2	-2

-1	1
1	-1

→ apply these weights,

$$-2(-1) + 2(1) + 2(1) + -2(-1) = 8$$

strange having this arbitrary #

this is b  
the solution for a diagonal

\* incorporate  $b$  into the eqn, so you can evaluate w reference to 0

$$w_1 a_{00} + w_2 a_{01} + w_3 a_{10} + w_4 a_{11} + b = \begin{cases} 0, & \text{diagonal} \\ < 0, & \text{NOT diagonal (otherwise)} \end{cases}$$

↳ translate vagueness of spec. + reflect that in the uncertainty!

\* say, instead of  $c_1$  and  $c_2$ , we want a continuum of colors →  $[c_1, c_2]$

↳ what methods do we already have to apply to this case??



translate grid into vector

$$\begin{bmatrix} a_{00} & a_{01} & a_{10} & a_{11} \end{bmatrix}^T \rightarrow \begin{bmatrix} 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} c_1 & c_2 \end{bmatrix}$$

$$\begin{bmatrix} c_1 & c_2 \end{bmatrix}$$

$$\dots$$

$$\begin{bmatrix} 0 & 1 \end{bmatrix}$$

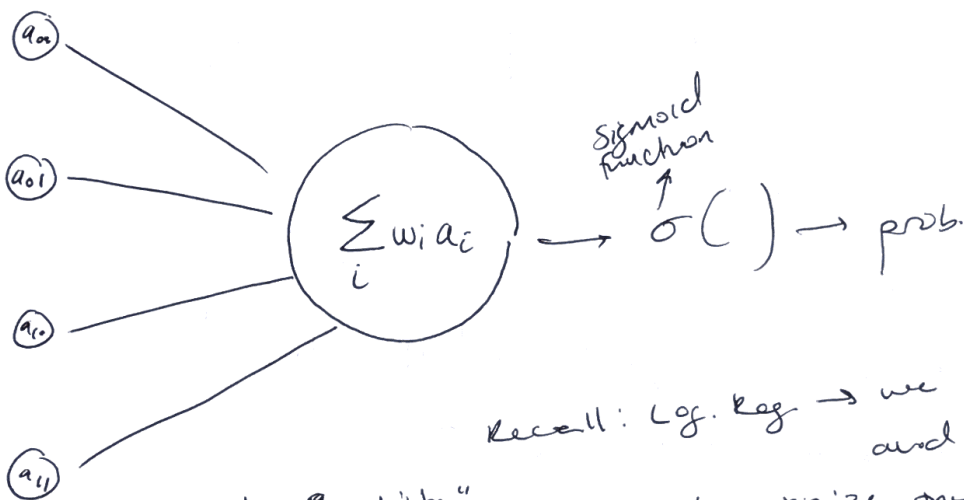
↳ prob. of being diagonal  
 $P(\text{diagonal})$

BINARY CLASSIFICATION

\* LOGISTIC REGRESSION

$$f\left(\begin{bmatrix} \oplus \end{bmatrix}\right) = \frac{1}{1 + e^{-\vec{w}\vec{a} + b}} \begin{matrix} \uparrow [a_{00}, a_{01}, \dots] \\ \downarrow [w_0, w_1, w_2, \dots] \end{matrix} \xrightarrow{\text{BIAS}}$$

goal:



Recall: Log. Reg.  $\rightarrow$  we want the weights,  $\vec{w}$ , and bias,  $b$ .

that maximize the probs of having seen the data we saw

"product of the probability"

$$\max \prod_{i=1}^n P(y_i = 1 | x_i)$$

# data points

"smart trick" = multiply by  $y_i$  to transform into single eq

$$= \min -\frac{1}{n} \sum_{i=1}^n \left[ \underbrace{y_i \log\left(\frac{1}{1 + e^{-\vec{w}x_i + b}}\right)}_{\substack{\text{prob} = 1 \\ y_i \\ \hookrightarrow P(y_i = 1)}} + \underbrace{(1 - y_i) \log\left(1 - \frac{1}{1 + e^{-\vec{w}x_i + b}}\right)}_{\substack{\text{prob. } y_i = 0 \\ \hookrightarrow P(y_i = 0)}} \right]$$

hide complexity under the "cost"  
(which we want to minimize)

$\hookrightarrow$  what is cost a function of?

$$= \min \text{Cost}(\vec{w}, b)$$

Cost( $\vec{w}, b$ )

\* want a numerical process for obtaining this minimum  $\rightarrow$  start @ random  $\vec{w}$  and random  $b$

$\hookrightarrow$  make steps in the right direction.

$$f\left(\begin{bmatrix} + \\ + \\ + \end{bmatrix}\right) = \frac{1}{1 + e^{-\vec{w}\vec{a} + b}} \quad \begin{matrix} [a_{00}, a_{01}, \dots] \\ \text{BIAS, accounts for consistent/systemic} \\ \text{dimming/brightening of the grid} \\ [w_1, w_2, w_3, \dots] \end{matrix}$$

(something consistently less bright)

Goal: find weights that lead to right amount of activation of  $f$  that we either meet a threshold or not

→ Suppose we start @  $w_0$  = random weight

some more negative than others



→ Look @  $\text{Cost}()$  around  $w_0$  (all directions surrounding)

+ (pos) - (neg)

→ pick best direction to minimize cost

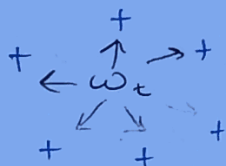
→ REPEAT ...

want cost to be as small as possible

⇒ clearly the best

"nudge" to give  $w_0$  & improve the cost

→ Want to reach a point  $w_c$  that looks like this

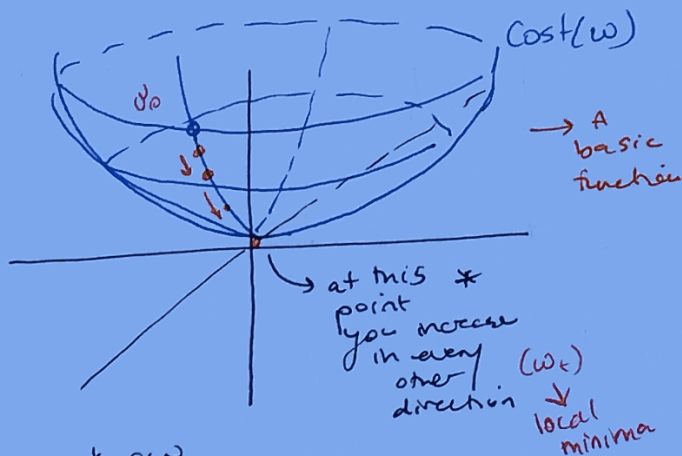


increases everywhere you go ...

cannot further minimize

(find some sort of min, maybe a local min)

### EXAMPLE



### EXAMPLE more complicated



Q How do you know direction of the best nudge?

→ highest downward rate -  $F$  change.  
(need the derivative)

⇒ discussion of GRADIENTS

(derivative in each direction)

→ at best, you find a local minimum.

- and there may be many local minima

→ the one you find depends where you start

Gradients: Best "nudge" should be in direction of steepest rate of change (decreasing)

Note: rate of change = derivative

"gradient"

$$\nabla f(x, y, z) = \frac{df}{dx} \vec{i} + \frac{df}{dy} \vec{j} + \frac{df}{dz} \vec{k}$$

→ what is  $\vec{i}, \vec{j}, \vec{k}$ ?  
unit vectors for each direction in space



↳ global rate of change is the rate of change of

each orthonormal component (rate of change in each direction)

gradient = rate of change ~~in direction~~ of each coordinate

$$(a\vec{i} + b\vec{j} + c\vec{k}) = [a, b, c]$$

\* unit length, orthonormal vectors

Example 1

$$f(x, y) = \frac{3}{2}x^2 - 2y$$

\* can already tell  $\Delta x$  affects  $f$  more than a  $\Delta y$

$$\nabla f = \frac{\partial f}{\partial x} \vec{i} + \frac{\partial f}{\partial y} \vec{j}$$

$$= 3x\vec{i} + (-2)\vec{j} = 3x\vec{i} - 2\vec{j}$$

evaluate  $\nabla f @ p = (0, 0)$

$$\nabla f_{(0,0)} = 3 \cdot 0 \vec{i} - 2 \vec{j} = -2 \vec{j}$$

↳ gradient @ origin

"move two units away in  $j^{\text{th}}$  direction"

$$p' = \alpha \cdot \nabla f_{(0,0)} + p$$

$$= \alpha \cdot \begin{bmatrix} 0 \\ -2 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ -2\alpha \end{bmatrix}$$

$$f(p') = \frac{3}{2}(0)^2 - 2(-2\alpha) = 4\alpha > f(p) = 0$$

↳ sanity check that we increased  $f$

↳ move  $\alpha$  steps in direction of gradient, increase the value of function

gradient moves in most pos. direction of change

recall: we want opp.

so add - to  $\nabla f$  and we'll find min

$$\text{say } p'' = \alpha \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} \alpha \\ 0 \end{bmatrix}$$

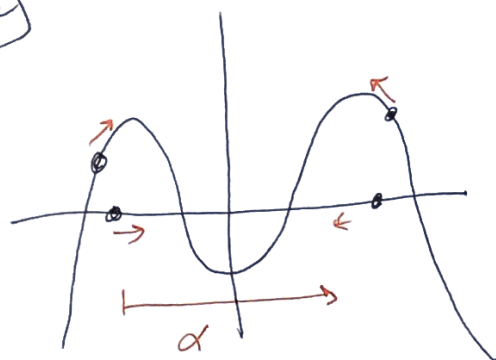
$$f(p'') = \frac{3}{2}\alpha^2$$

$$\left[ 4\alpha > \frac{3}{2}\alpha^2 \right] ?$$

↳ need to take SMALL steps so you can capture ~~start~~ it all   
 ~~reason~~ ↳ don't want to overshoot direction of min



Example



\* ~~and~~  
 $\alpha$  too big... end up ping ponging back and forth. (miss the min)

don't do this! too big jump  $\rightarrow$  missed the min

## Gradient Descent $\star \rightarrow$ Important

$\nabla f \rightarrow$  direction of steepest increase  
 $-\nabla f \rightarrow$  direction of steepest decrease

- 1) Define  $\alpha$  (not too small, not too big)  $\rightarrow$  size of the step we want to take (how much of a nudge you want to give)  
<sup>tuning parameter</sup>
- 2) Initialize  $p$  to be random
- 3)  $p_{\text{new}} = -\alpha \nabla f(p) + p$   $\rightarrow$  update  $p$   $\leftarrow \alpha$  steps in direction of gradient  $\rightarrow$  update  $p$
- 4)  $p \leftarrow p_{\text{new}}$
- 5) Repeat 3 & 4 until  $p \approx p_{\text{new}}$   $\rightarrow$  or gradient is zero  
 (much easier constraint, numerically)

Projects: Early Insights Presentation - lightning talk! (43 min @ most)

$\Rightarrow$  motivation/background  
 $\star$  progress  $\star$  next steps

$\star$  limitations  
 what challenges have we encountered?

practice your pitch so it's FAST!!

$\Rightarrow$  powerpoint / visualizations

- slides - copy/paste from Desmos/GeoGebra

more info this weekend

"Why is your project awesome?"