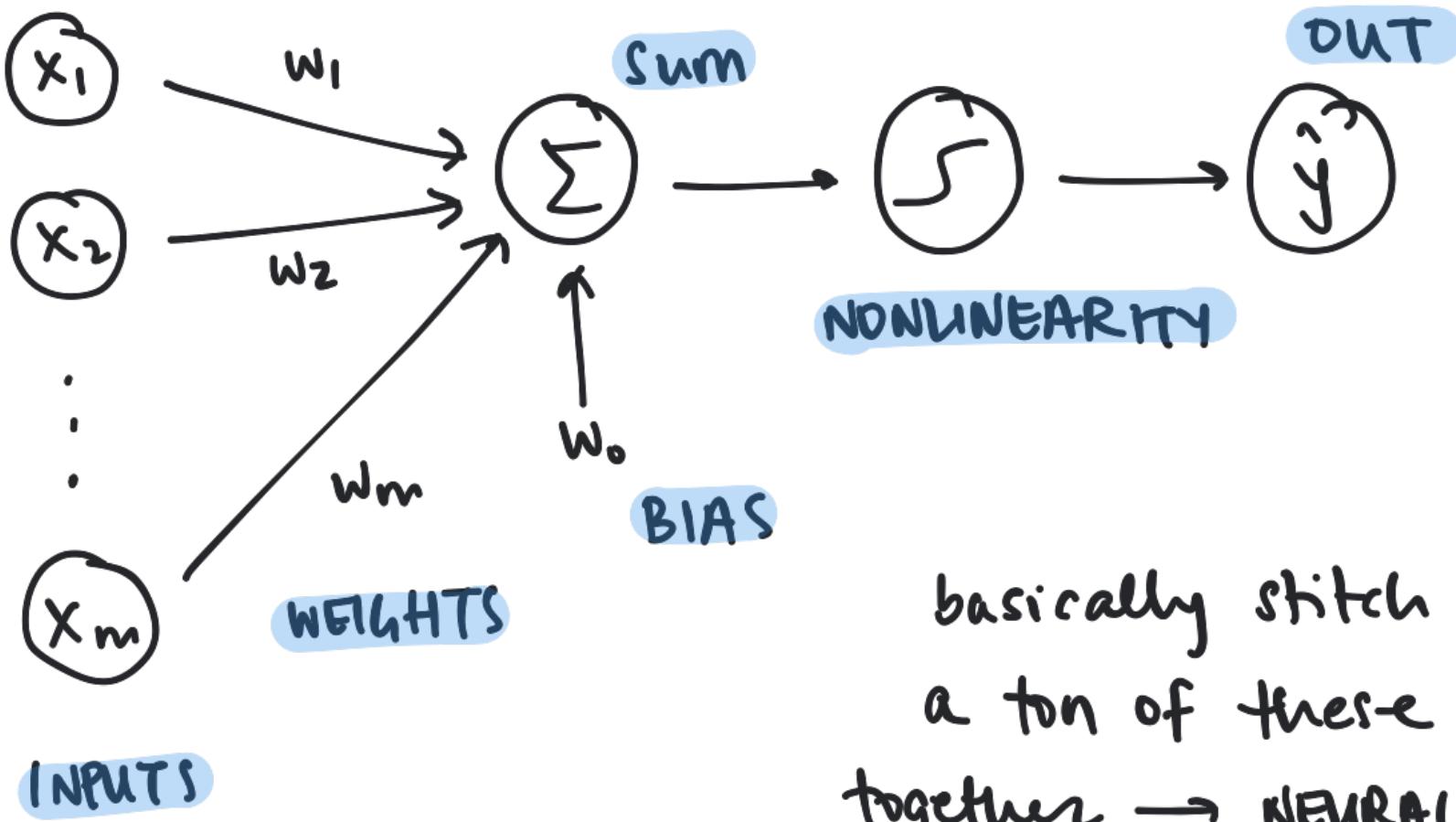


NEURAL NETS

building block is the neuron:



basically stitch
a ton of these
together \rightarrow NEURAL NET!
(we'll do this later)

can write this mathily:

$$\begin{aligned}\hat{y} &= g(w_0 + \sum_{i=1}^m x_i w_i) = g(z) \\ &= g(w_0 + X^T W), \quad X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}, \quad W = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}\end{aligned}$$

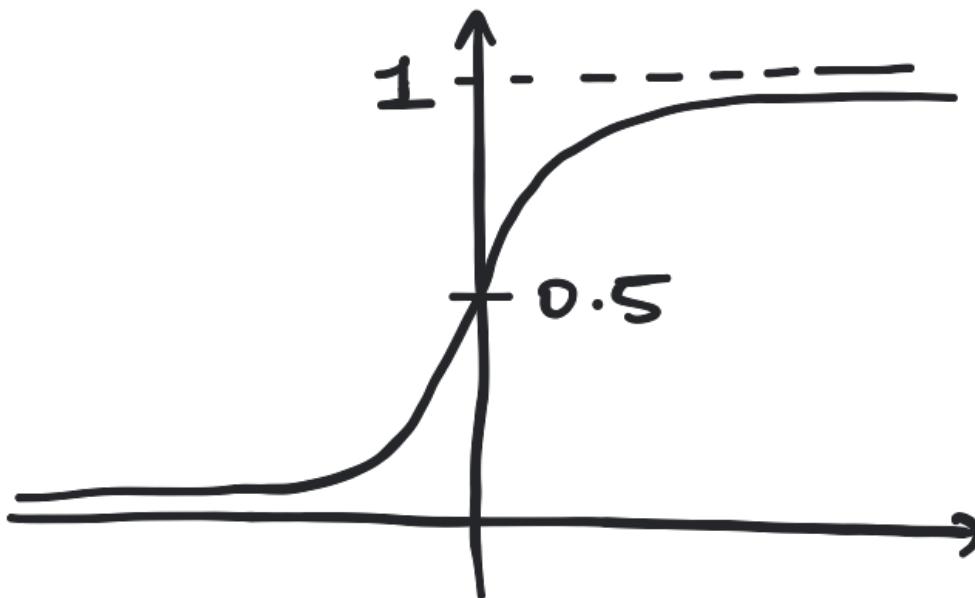
what is g ? **ACTIVATION FUNCTION** \Rightarrow

a popular choice of g is the SIGMOID:

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

$$\sigma: \mathbb{R} \rightarrow (0, 1)$$

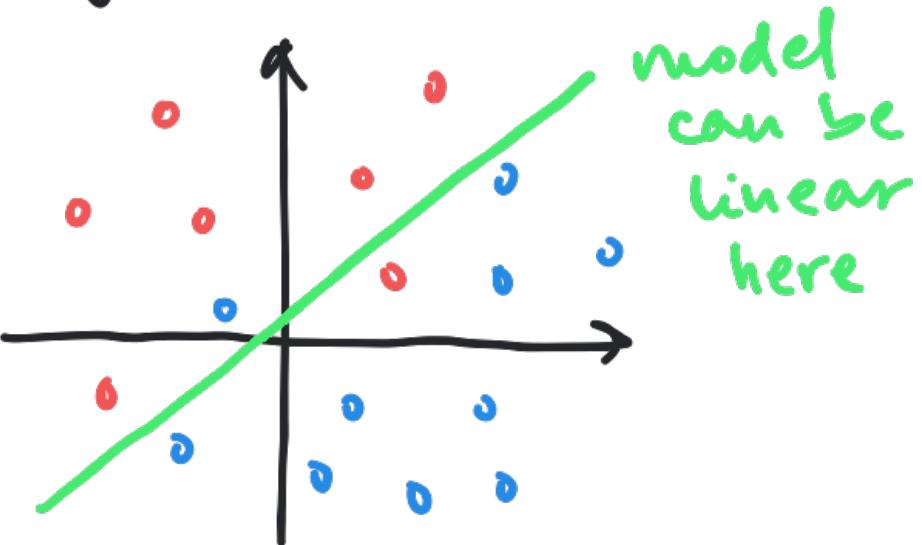
↳ can be interpreted as a probability!



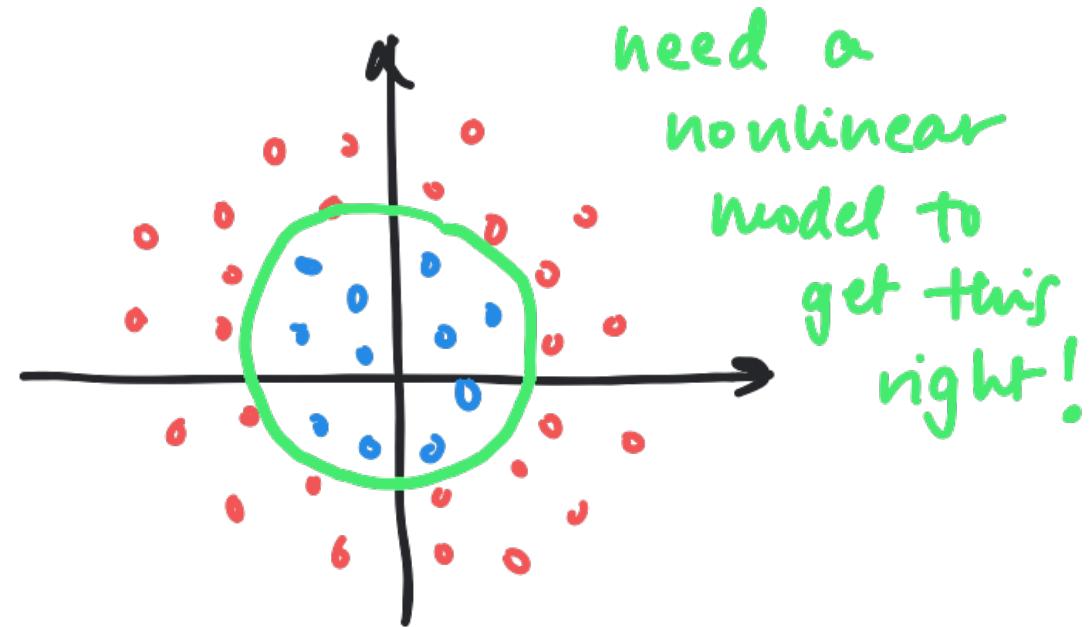
the choice of activation function is pretty important!

↳ builds essential nonlinearities into the model.

e.g. classification problems:



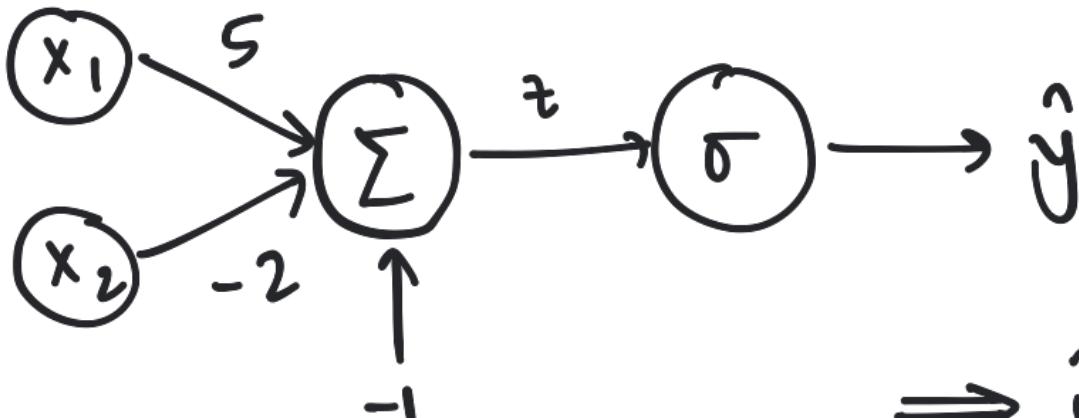
model
can be
linear
here



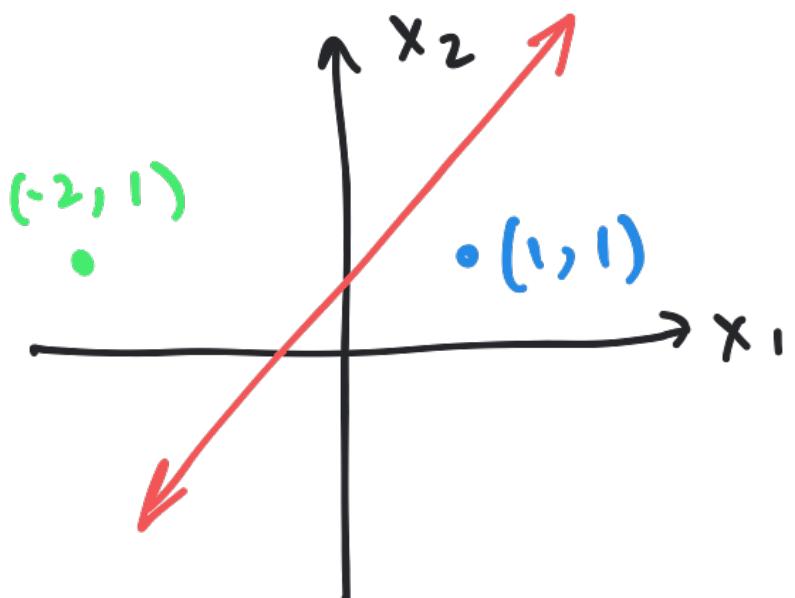
need a
nonlinear
model to
get this
right!

how does classification work?

Simplest case. say we have the right weights.
using just one neuron for now.



$$\Rightarrow \hat{y} = \sigma(1 + 5x_1 - 2x_2)$$



this is a **LINEAR LOGISTIC CLASSIFIER**.

$$\begin{aligned}\hat{y} &= \sigma(1 + 5 - 2) \\ &= \sigma(4) = 0.982 \\ &> 0.5\end{aligned}$$

$$\begin{aligned}\hat{y} &= \sigma(1 - 10 - 2) \\ &= \sigma(-11) = 1.7 \times 10^{-5} \\ &< 0.5\end{aligned}$$

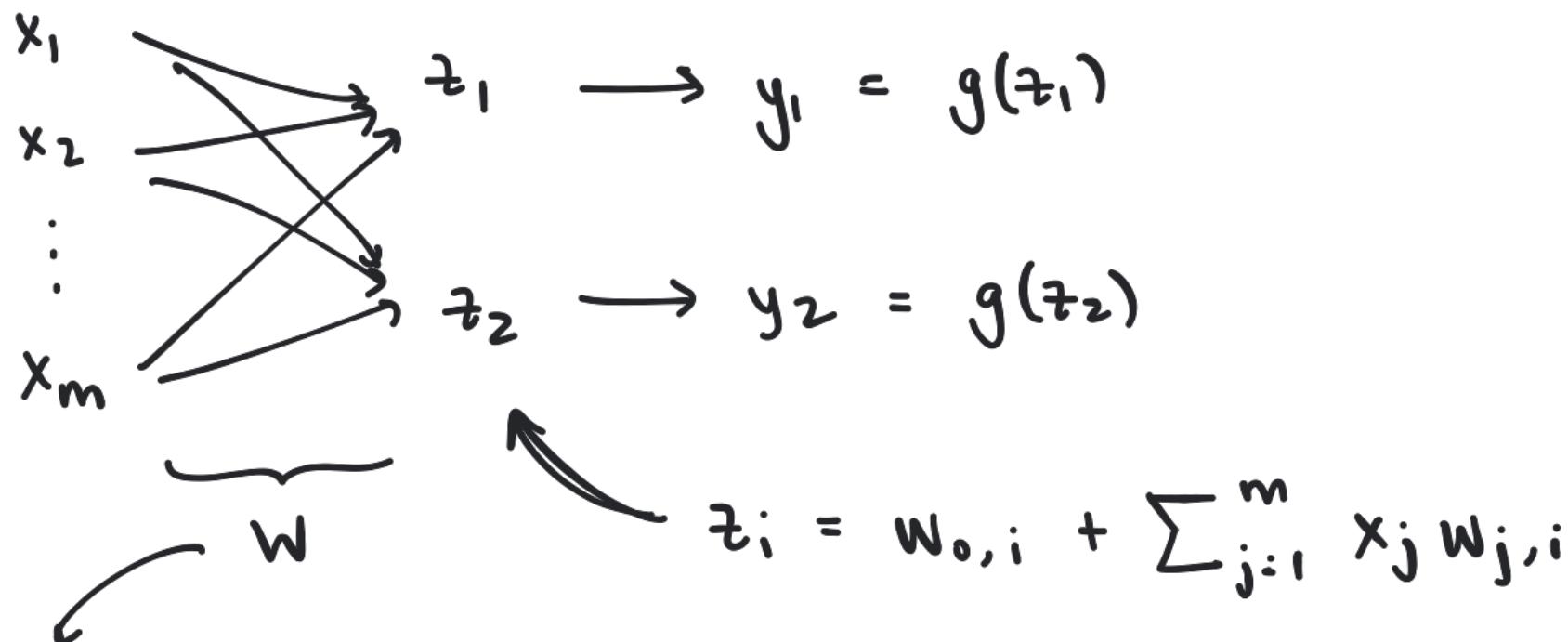
↳ predicts A if $\hat{y} > 0.5$
 B if $\hat{y} < 0.5$

} we can set the threshold!
 \hat{y} gives us probabilities basically.

putting it together...

so that's ONE neuron. put a bunch together
and you can get all kinds of fancy behavior!

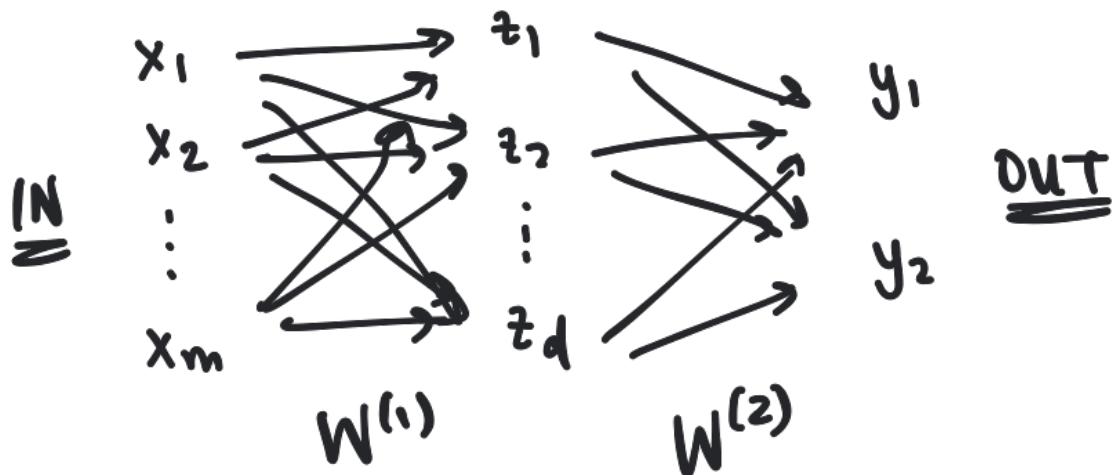
↳ hard to draw out graphically bc LOTS of dimensions.



these weights define a LAYER of the neural net.

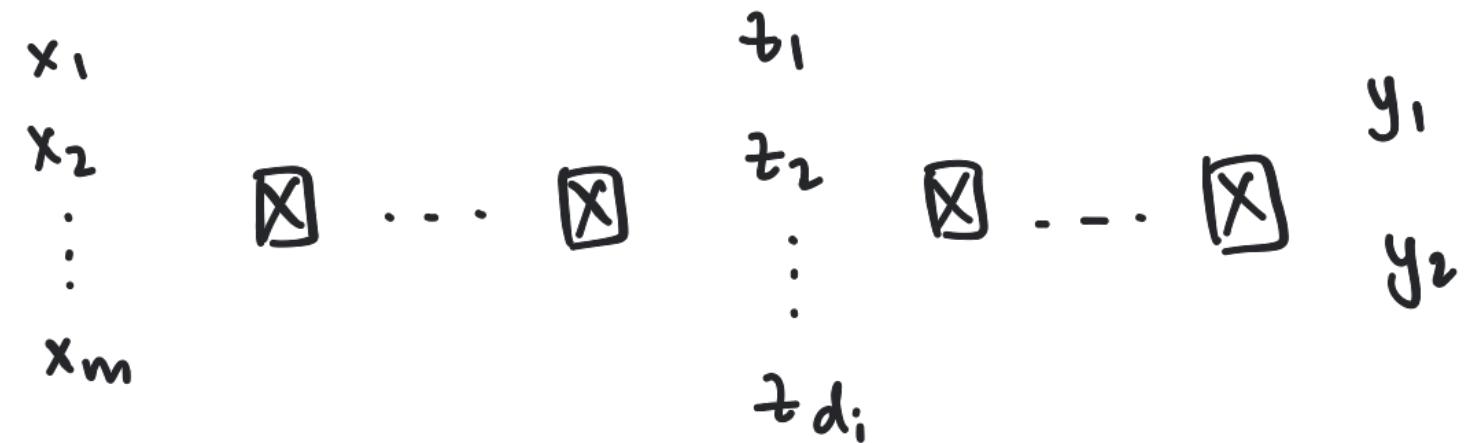
because it's so connected, call it a DENSE LAYER.

a **SINGLE LAYER NEURAL NET**:



- the z 's form a **HIDDEN LAYER**.
↳ neither IN or OUT.
- 2 weight matrices $W^{(1)}$

can stack these! (\otimes means dense layer)



okay so you have a bunch of layers. now what???

↳ need to TRAIN the neural net to optimize W .

have to quantify how wrong its predictions are.

↳ use a LOSS FUNCTION.

$$L(f(x^{(i)}; W), y^{(i)})$$

$f(x, W)$: prediction
 y : actual (from training data)

e.g. for classification problem:

$$J(W) = \frac{1}{n} \sum_{i=1}^n y^{(i)} \cdot \log(f(x^{(i)}; W)) + (1 - y^{(i)}) \cdot \log(1 - f(x^{(i)}; W))$$

↓
"objective function"

⇒ here, $y^{(i)}$ is 0/1. $f(x, W)$ is $[0, 1]$.

how to train your dragon

LOSS OPTIMIZATION is the goal.

↳ want to be as correct as you can* on the training data

→ want to minimize $J(w)$.

* barring overfitting!

to get: $w^* = \underset{w}{\operatorname{argmin}} J(w)$ ← best weights for the NN

can use a minimization algorithm like...

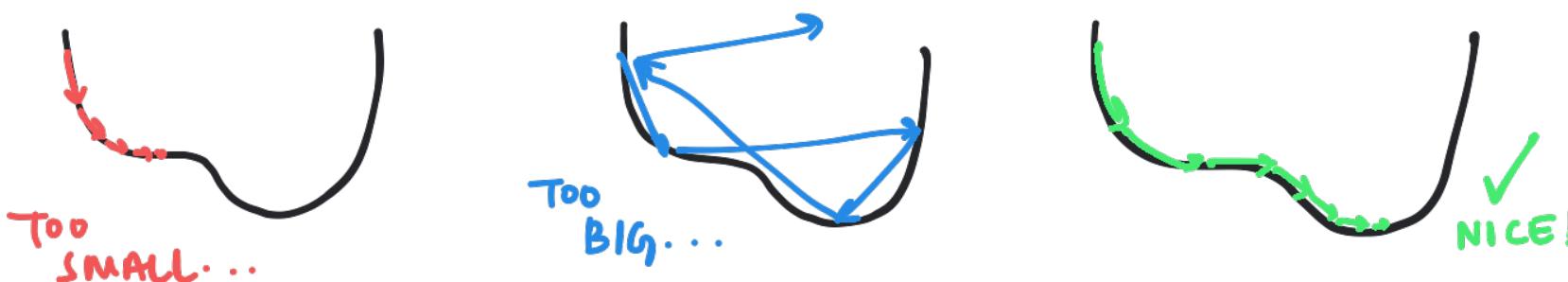
GRADIENT DESCENT

calculate gradient + update weights.

$$w_{\text{new}} = w - \eta \cdot \frac{\partial J(w)}{\partial w} \quad \eta: \text{learning rate}$$

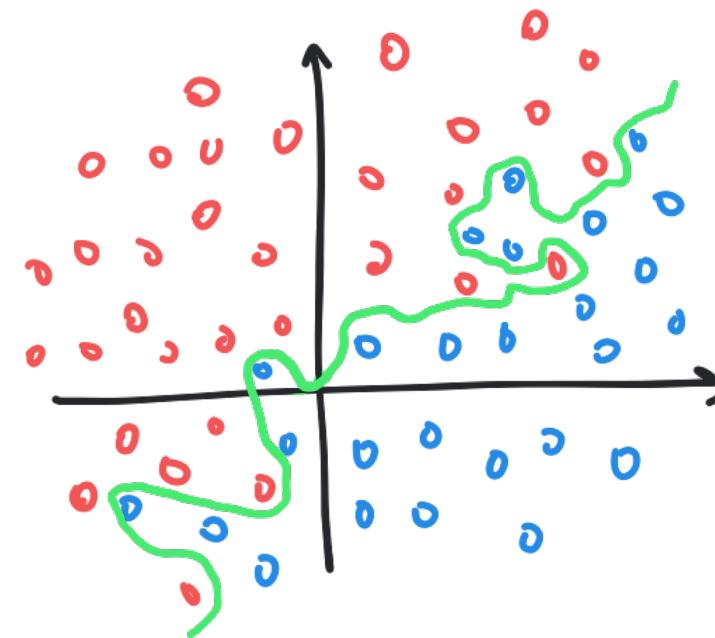
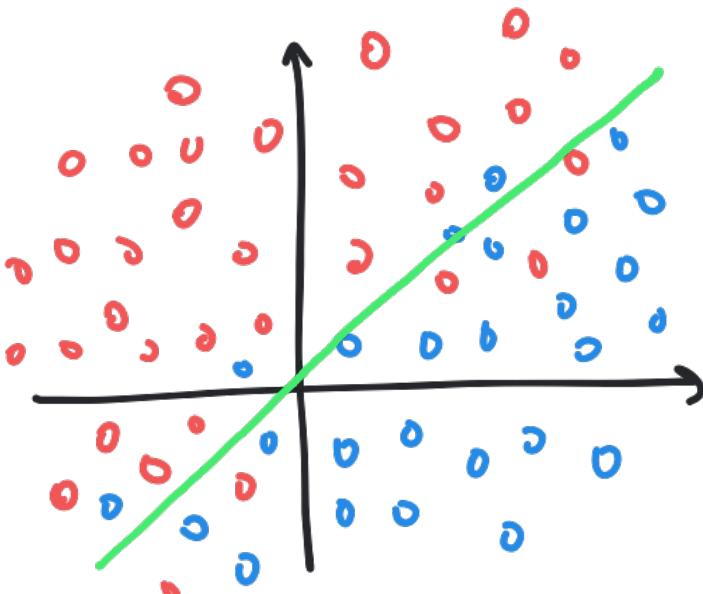
finding $\partial J / \partial w$ is nontrivial (take 036 to understand!)

setting η can make a huge difference too!



can have adaptive learning rates that scale in various helpful ways.

overfitting?



which looks more reasonable?

↳ want to prevent overcomplicated models.

good way to tackle this is not to overtrain.

↳ "early stopping"



✗ STOP TRAINING HERE!

let's check out
some examples!

TO THE NOTEBOOK!