# INTRO TO DEEP LEARNING

INPUT LAYER → HIDDEN LAYER → OUTPUT LAYER

AREA
#ROOMS  → FAMILY SIZE
LOCATION → WALK ABILITY → ○ — PRICE
WEALTH → SCHOOL QUAL

## SUPERVISED LEARNING

| INPUT: X | OUTPUT: Y | NN TYPE |
|---|---|---|
| HOME FEATURES | PRICE | STANDARD NN |
| AD + USER INFO | WILL CLICK ON AD (0/1) | |
| IMAGE | OBJECT (1...1000) | CONV. NN (CNN) |
| AUDIO | TEXT TRANSCRIPT | RECURRENT NN (RNN) |
| ENGLISH | CHINESE | |
| IMAGE/RADAR | POS OF OTHER CARS | CUSTOM/HYBRID |

$X_1$
$X_2$  → ○○○ → ○ → $\hat{y}$
$X_3$

STANDARD NN

× → ☐ → ☐ → ○ → $\hat{y}$

CONVOLUTIONAL NN

## NETWORK ARCHITECTURES

$\hat{y}$ $\hat{y}$ $\hat{y}$
(a)(a)(a)
x  x  x

RECURRENT NN

## NNs CAN DEAL WITH BOTH STRUCTURED & UNSTRUCTURED DATA

STRUCTURED

"THE QUICK BROWN FOX" UNSTRUCTURED

HUMANS ARE GOOD AT THIS

## WHY NOW?

LOTS OF DATA
HARDWARE
OPTIMIZED ALGOS

LARGE NN
MED NN
SMALL NN
CLASSIC ML

PERFORM. vs AMT. OF DATA (LABELED)

ONE OF THE BIG BREAKTHROUGHS HAS BEEN MOVING FROM SIGMOID TO RELU FOR FASTER GRADIENT DESCENT

SIGMOID          RELU

IDEA → CODE → EXPERIM. → (loop)

FASTER COMPUTATION IS IMPORTANT TO SPEED UP THE ITERATIVE PROCESS

@TessFerrandez

# LOGISTIC REGRESSION
## AS A NEURAL NET

## BINARY CLASSIFICATION

$X \rightarrow$ 1 : CAT
0 : NOT CAT

$y$

LOGISTIC REGRESSION

$= $ LINEAR REGRESSION $+$ $\sigma$ SIGMOID

$\hat{y} = \sigma(wx+b)$

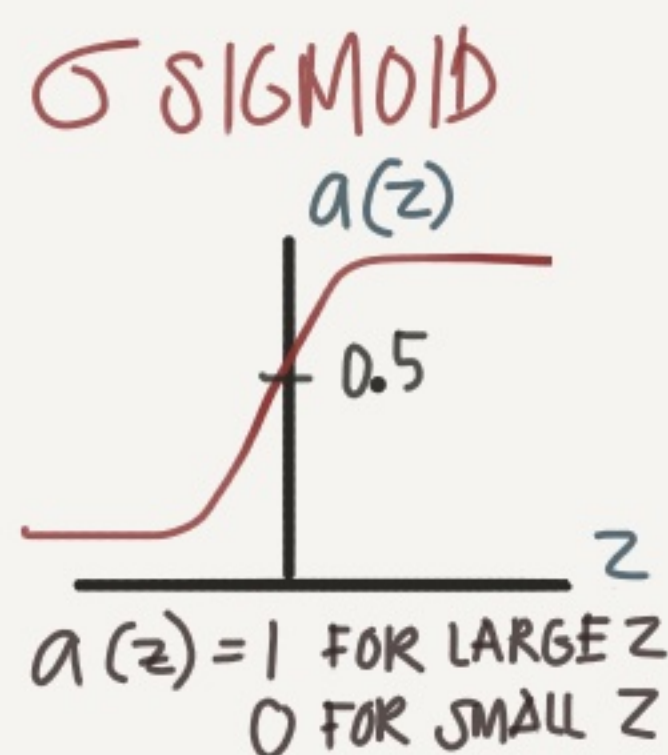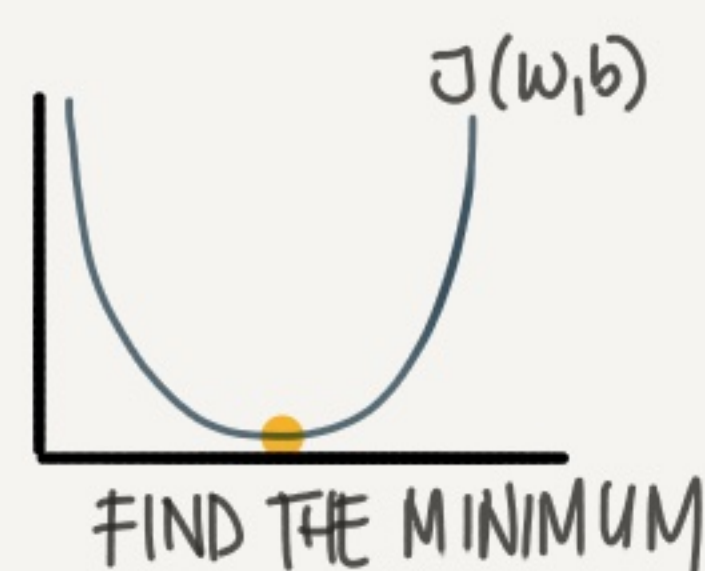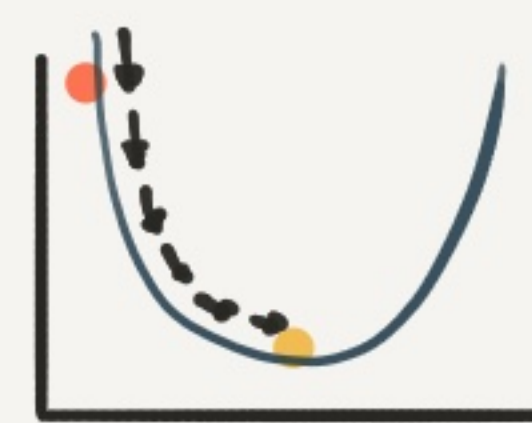$\hat{y} = wx+b$

$a(z)$
0.5
$a(z) = 1$ FOR LARGE $z$
0 FOR SMALL $z$

## THE TASK IS TO LEARN $w$ & $b$ BUT HOW?

A : OPTIMIZE HOW GOOD THE GUESS IS BY MINIMIZING THE DIFF BETWEEN GUESS ($\hat{y}$) AND TRUTH ($y$)

$LOSS = \mathcal{L}(\hat{y}, y)$

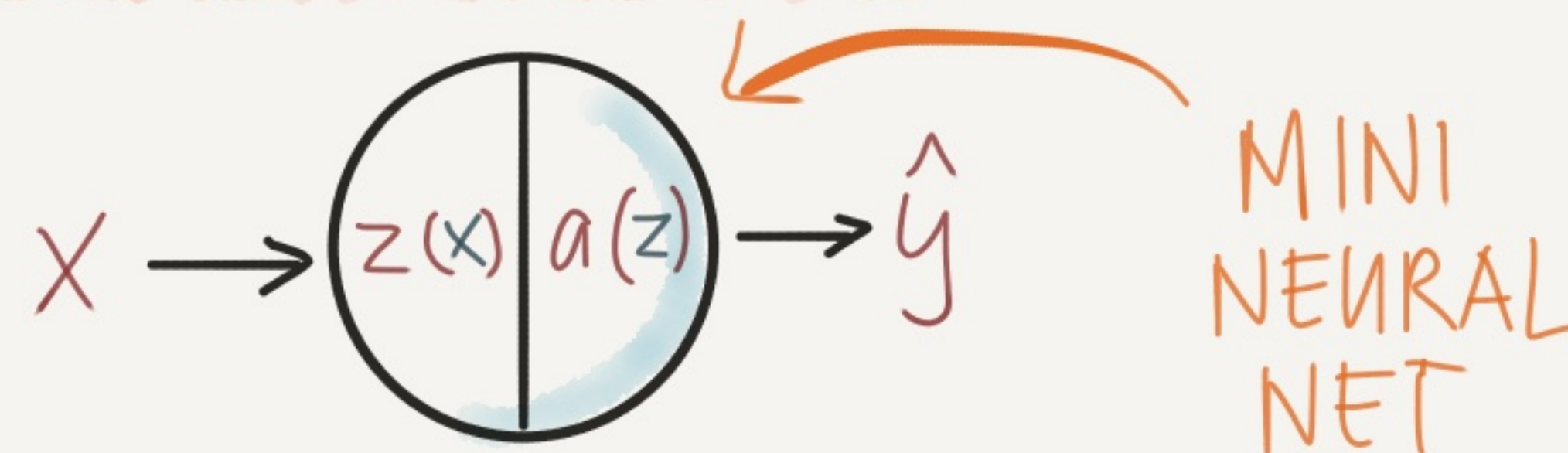$COST = J(w,b) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$

COST = LOSS FOR THE ENTIRE DATASET

$J(w,b)$

FIND THE MINIMUM

## FINDING THE MINIMUM WITH GRADIENT DESCENT

1. FIND THE DOWNHILL DIRECTION (USING DERIVATIVES)
2. WALK (UPDATE $w$ & $b$) AT A $\alpha$ LEARNING RATE

REPEAT UNTIL YOU REACH BOTTOM (CONVERGE)

## PUTTING IT ALL TOGETHER

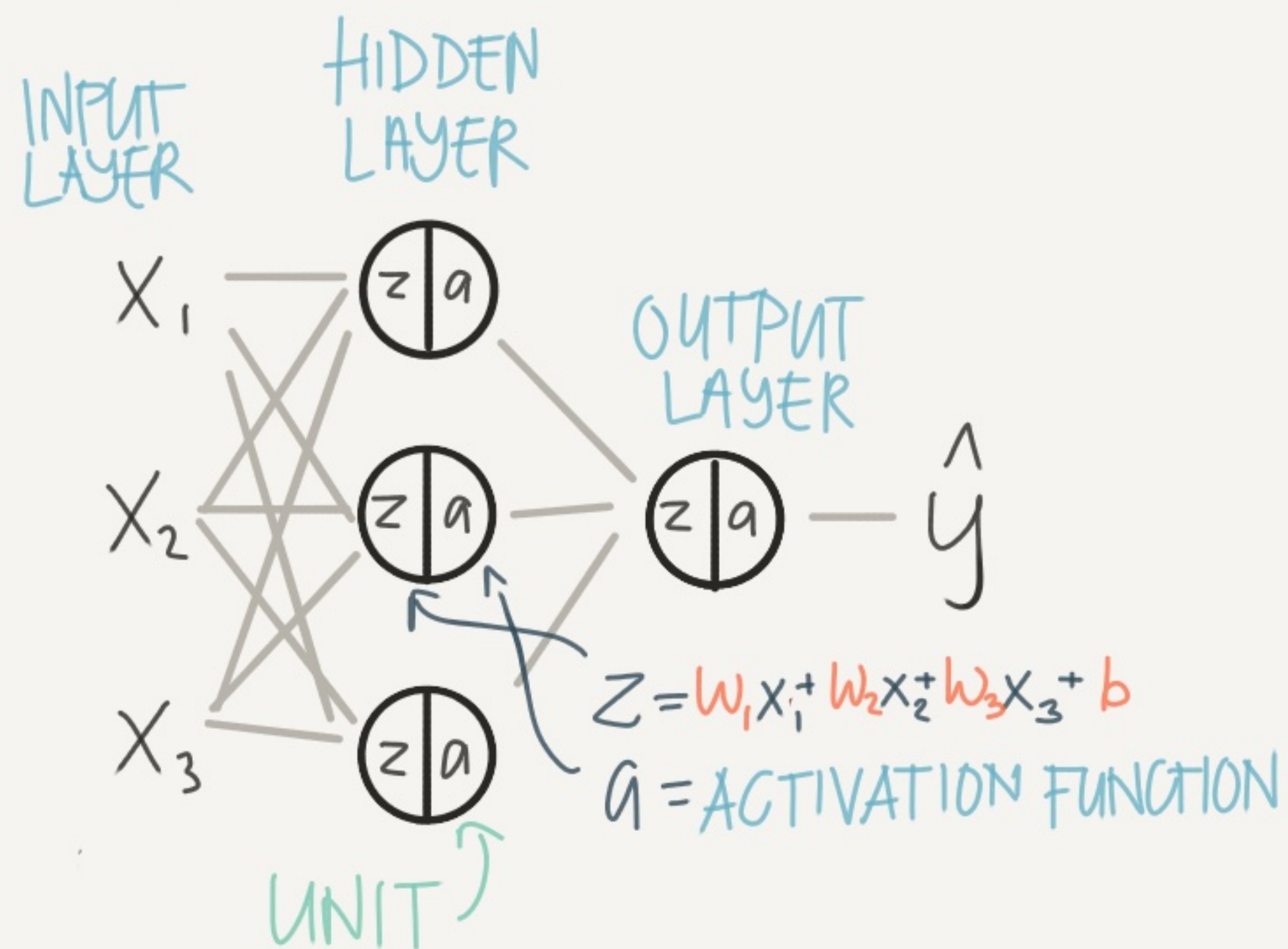$X \rightarrow z(x) | a(z) \rightarrow \hat{y}$

MINI NEURAL NET

$z(x) = wx + b$
$\hat{y} = a(z) = \sigma$ SIGMOID $(z)$

1. FORWARD PROPAGATION • CALCULATE $\hat{y}$
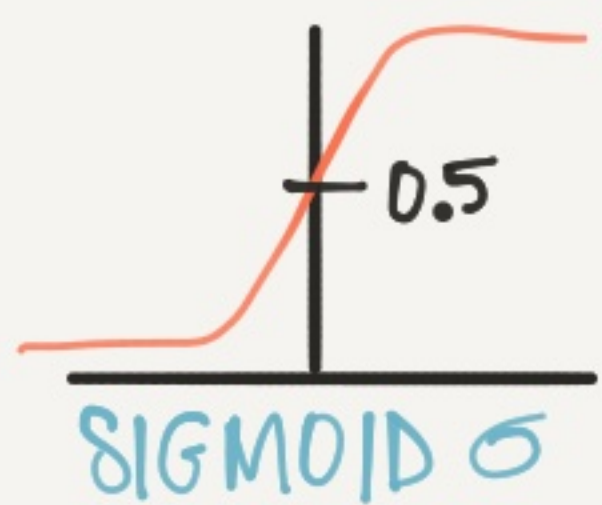2. BACKWARD PROPAGATION • GRADIENT DESCENT + UPDATE $w$ & $b$
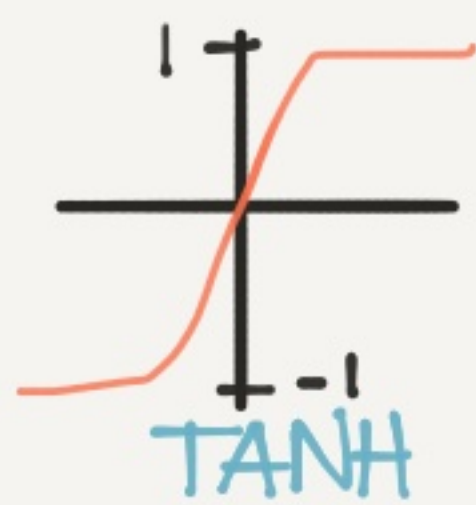
REPEAT UNTIL IT CONVERGES

@TessFerrandez

# 2 LAYER NEURAL NET

INPUT LAYER
HIDDEN LAYER
OUTPUT LAYER

$X_1$
$X_2$
$X_3$

$\hat{y}$

$Z = W_1 X_1 + W_2 X_2 + W_3 X_3 + b$
$a = $ ACTIVATION FUNCTION

UNIT

# WHY ACTIVATION FUNCTIONS?

EX. WITH NO ACTIVATION — $a = Z$

$a^{[1]} = Z^{[1]} = W^{[1]} X + b^{[1]}$   LAYER 1
$a^{[2]} = Z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$   LAYER 2

PLUG IN $a^{[1]}$

$a^{[2]} = W^{[2]} (W^{[1]} X + b^{[1]}) + b^{[2]}$
$= W^{[2]} W^{[1]} X + W^{[2]} b^{[1]} + b^{[2]}$

$W' X + b'$   ← LINEAR FUNCTION

WE COULD JUST AS WELL HAVE SKIPPED THE WHOLE NEURAL NET & USED LIN. REGR.
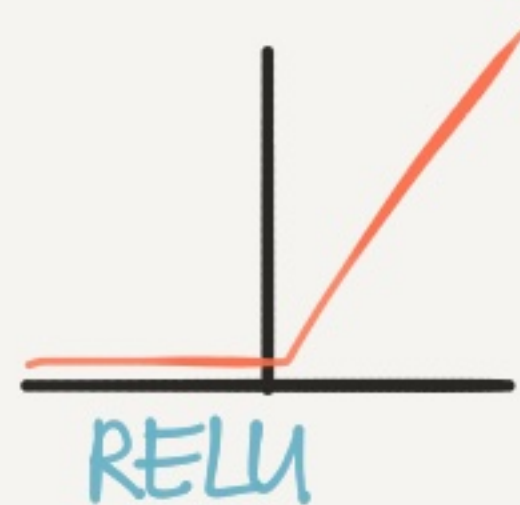
# ACTIVATION FUNCTIONS

SIGMOID $\sigma$
0.5
BINARY CLASSIFIER
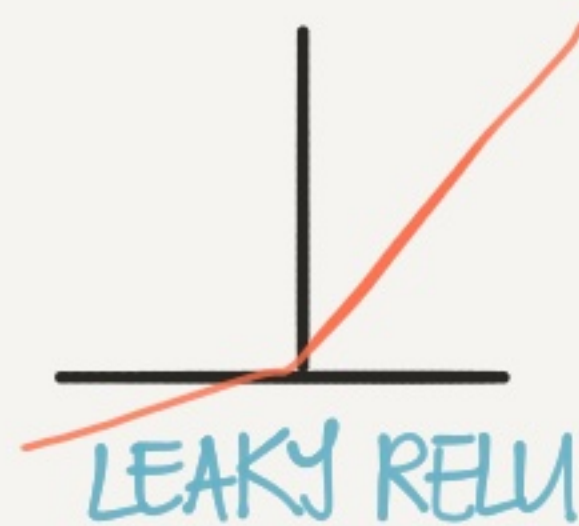- ONLY USED FOR OUTPUT LAYER

SLOW GRAD DESCENT SINCE SLOPE IS SMALL FOR LARGE/SMALL VAL

TANH
1
-1
NORMALIZED ⇒ GRADIENT DESCENT IS FASTER

RELU
DEFAULT CHOICE FOR ACTIVATION
SLOPE = 1/0

LEAKY RELU
AVOIDS UNDEF SLOPE AT $\emptyset$ BUT RARELY USED IN PRACTICE

# INITIALIZING $W + b$

WHAT IF: INIT TO $\emptyset$

THIS WILL CAUSE ALL THE UNITS TO BE THE SAME AND LEARN EXACTLY THE SAME FEATURES

SOLUTION: RANDOM INIT

BUT ALSO WANT THEM SMALL SO RAND * 0.01

← HYPERPARAM

@TessFerrandez

# DEEP
## NEURAL NETS
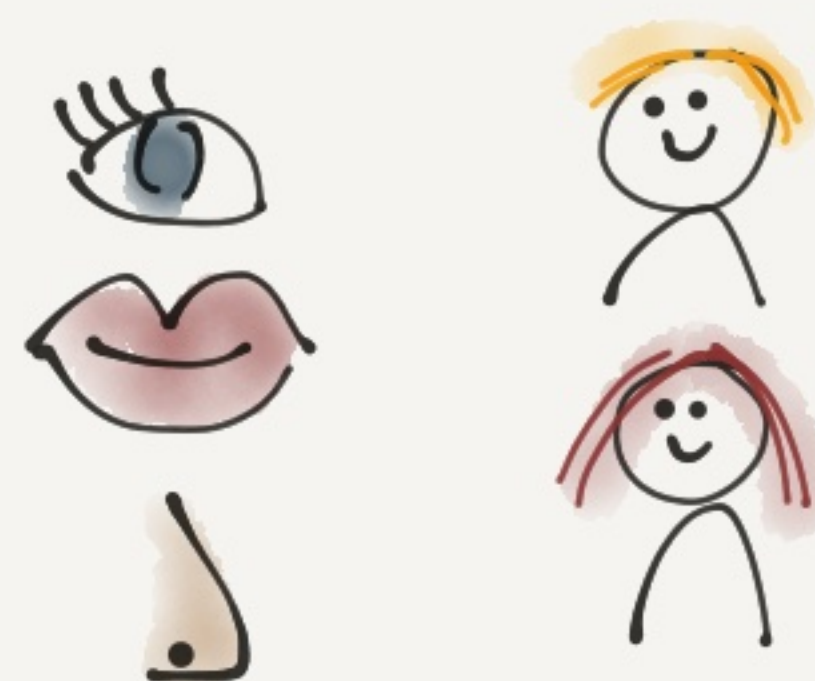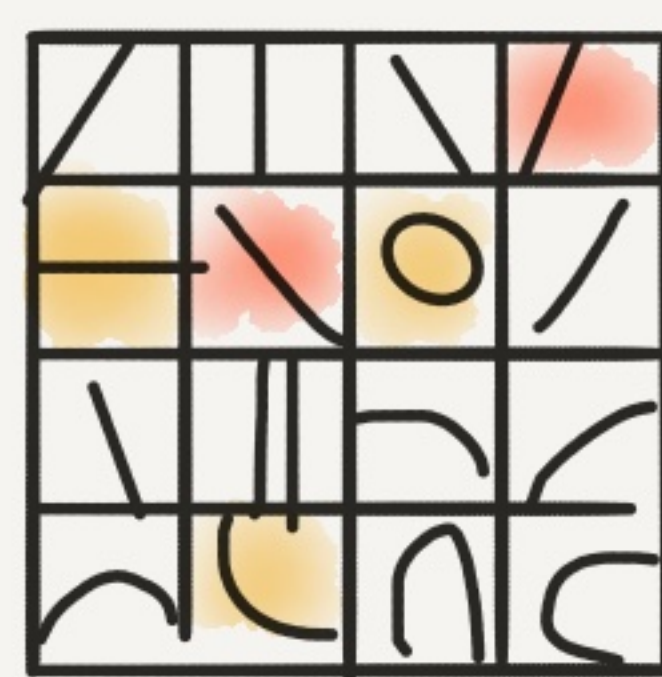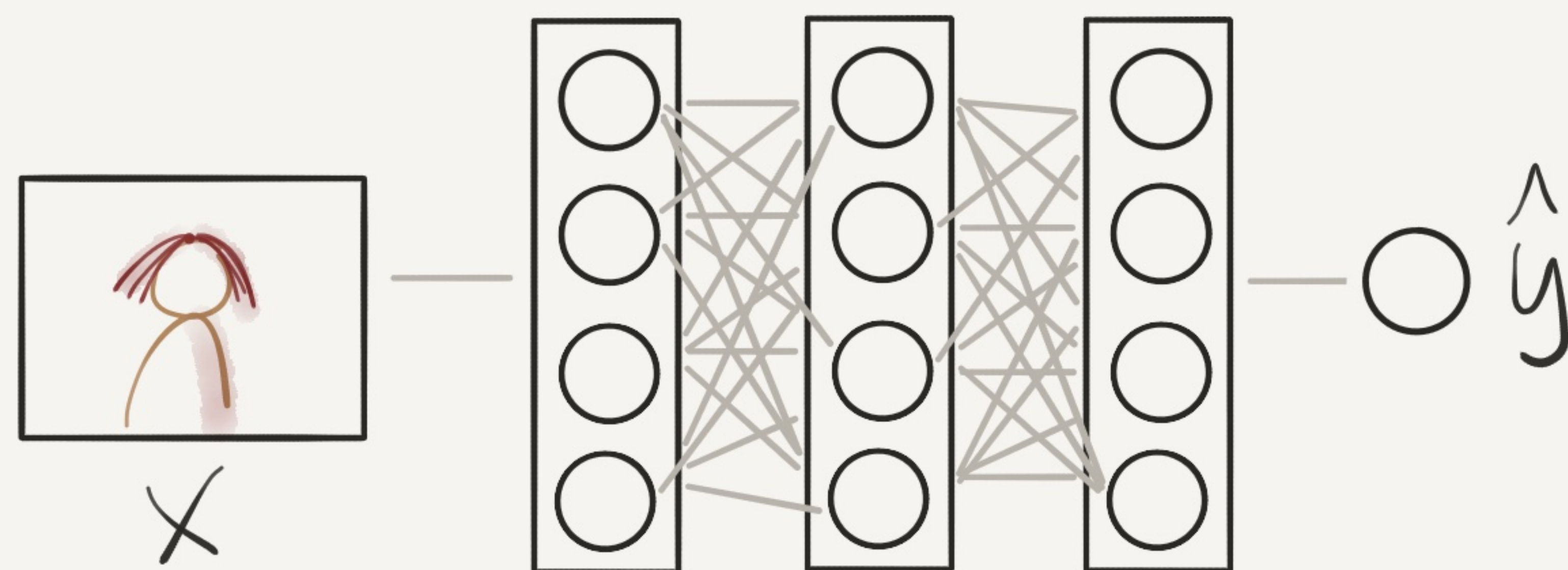
## WHY DEEP NEURAL NETS?

THERE ARE FUNCTIONS A SMALL DEEP NET CAN COMPUTE THAT SHALLOW NETS NEED EXP. MORE UNITS TO COMP.



$X$

$\hat{y}$

**VERY DATA HUNGRY**

**NEED LOTS OF COMPUTER POWER**

**ALWAYS VECTORIZE**
VECTOR MULT·CHEAPER THAN FOR LOOPS
**COMPUTE ON GPUs**

## LOTS OF HYPERPARAMS

LEARNING RATE $\alpha$     # HIDDEN UNITS
# ITERATIONS     CHOICE OF ACTIVATION
# HIDDEN LAYERS     MOMENTUM
    MINI·BATCH SIZE
    REGULARIZATION

LOW LEVEL AUDIO WAVE FEATURES ↗↘ PITCH — PHONEMES — WORDS — SENTENCES

CAT

@TessFerrandez