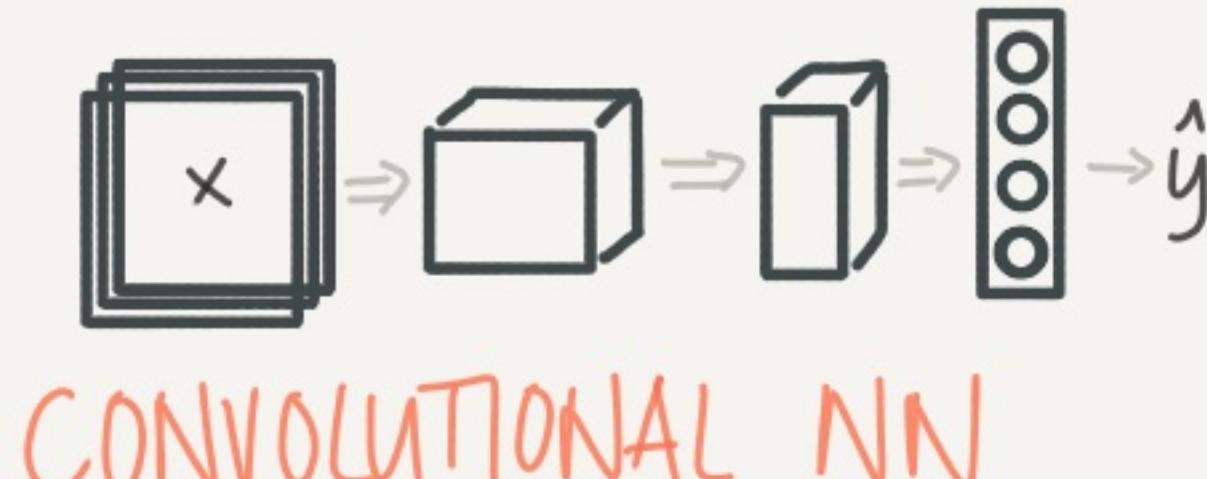
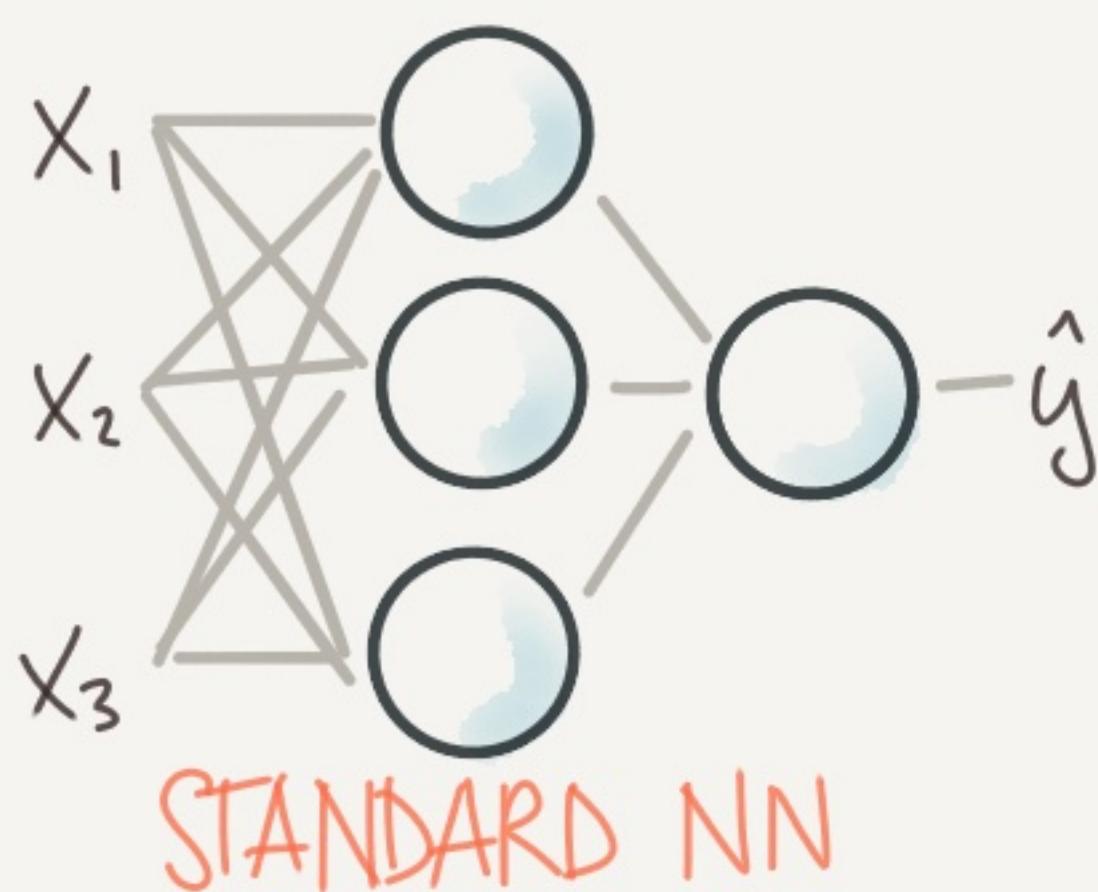


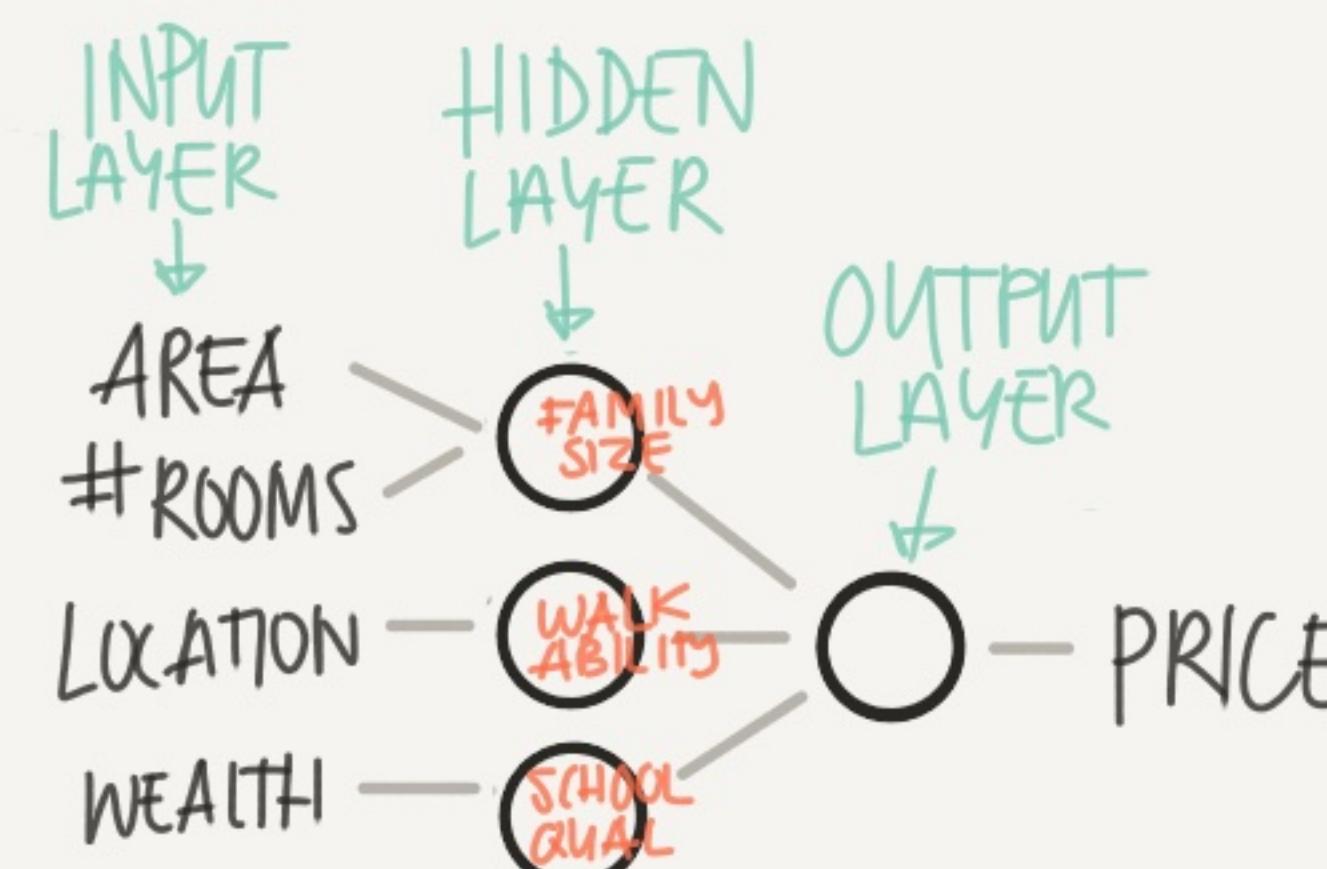
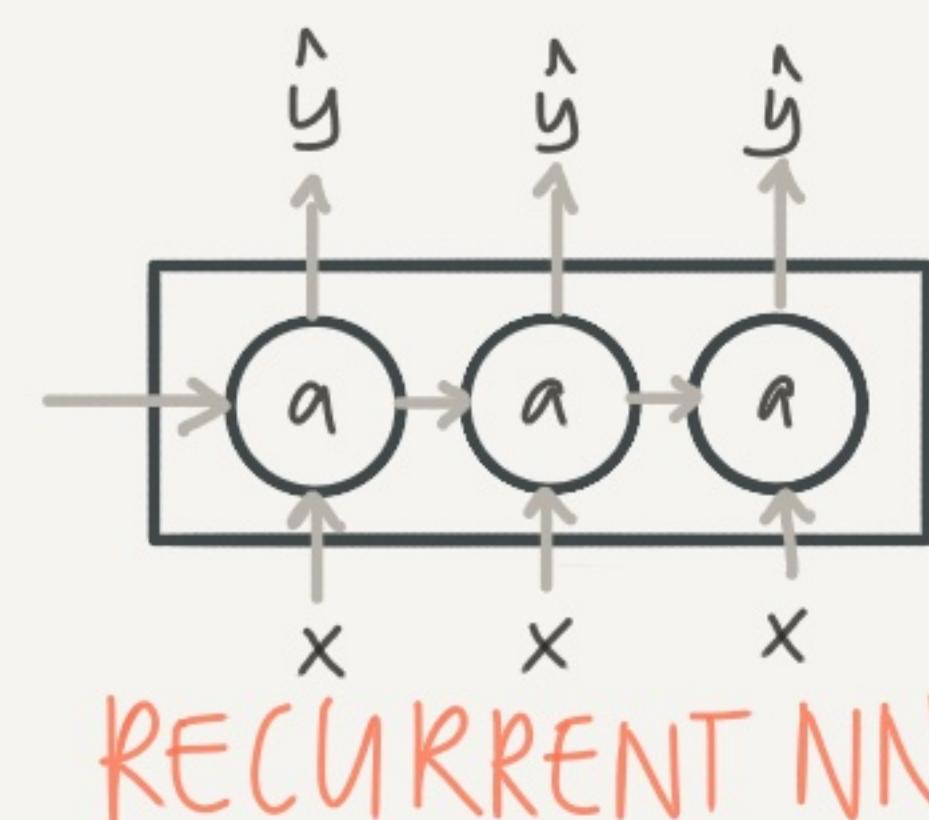
INTRO TO DEEP LEARNING

SUPERVISED LEARNING

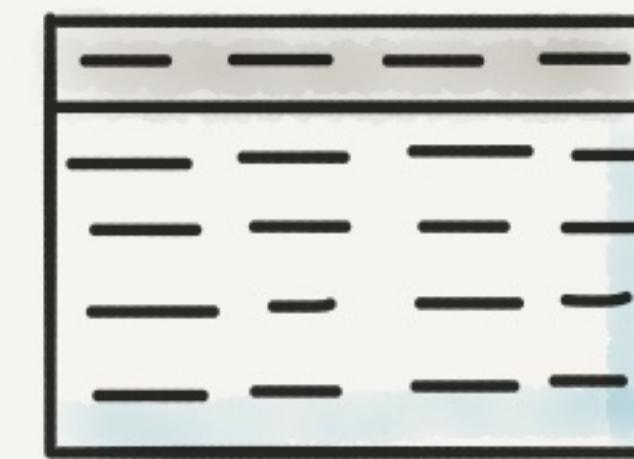
INPUT: X	OUTPUT: y	NN TYPE
HOME FEATURES	PRICE	STANDARD NN
AD+USER INFO	WILL CLICK ON AD (0/1)	
IMAGE	OBJECT (1...1000)	CONV. NN (CNN)
AUDIO	TEXT TRANSCRIPT	RECURRENT NN (RNN)
ENGLISH	CHINESE	
IMAGE/RADAR	POS OF OTHER CARS	CUSTOM/HYBRID



NETWORK ARCHITECTURES



NNs CAN DEAL WITH BOTH STRUCTURED & UNSTRUCTURED DATA



STRUCTURED

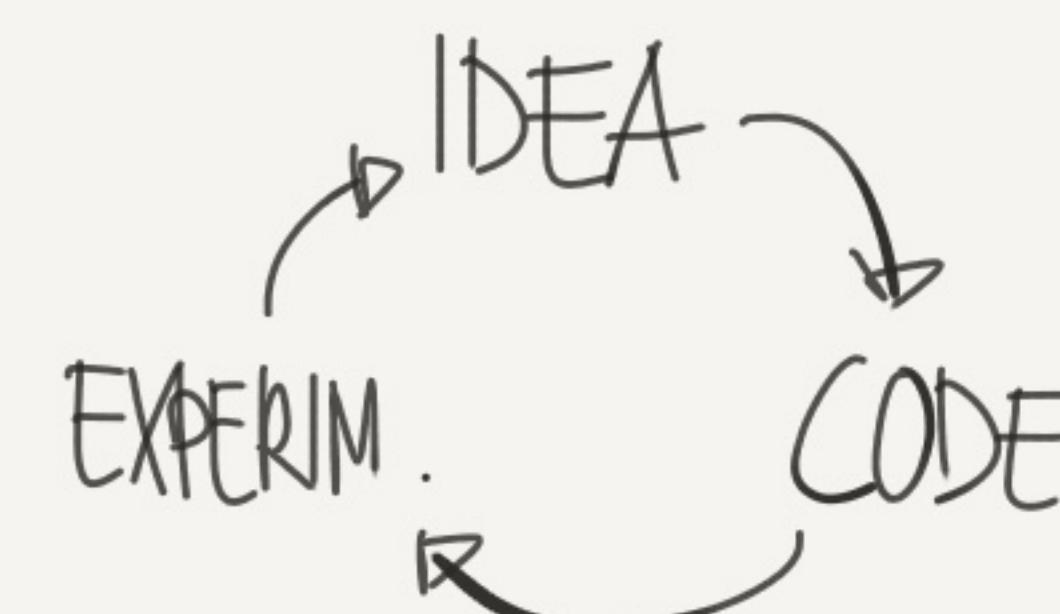
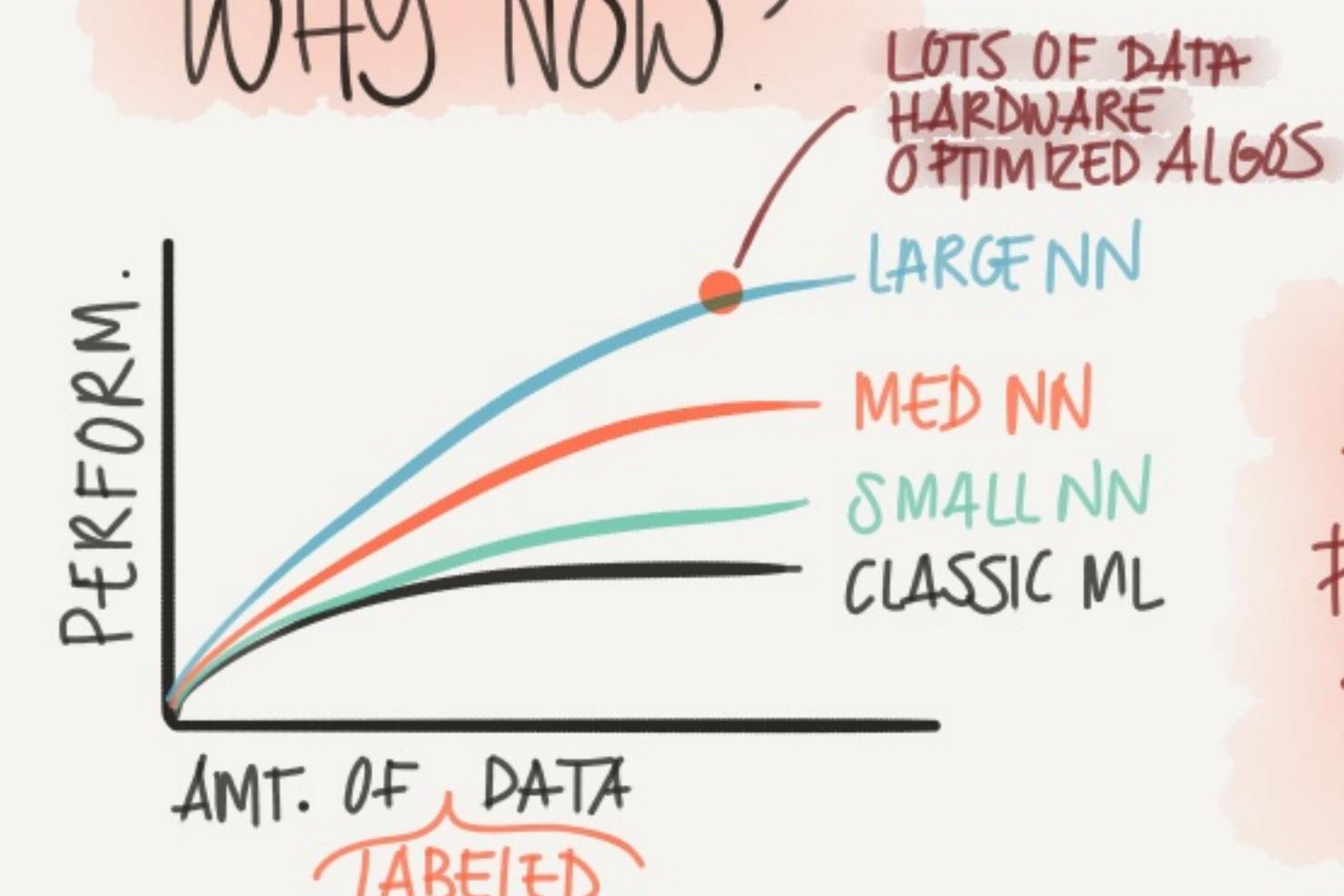


"THE QUICK BROWN FOX"

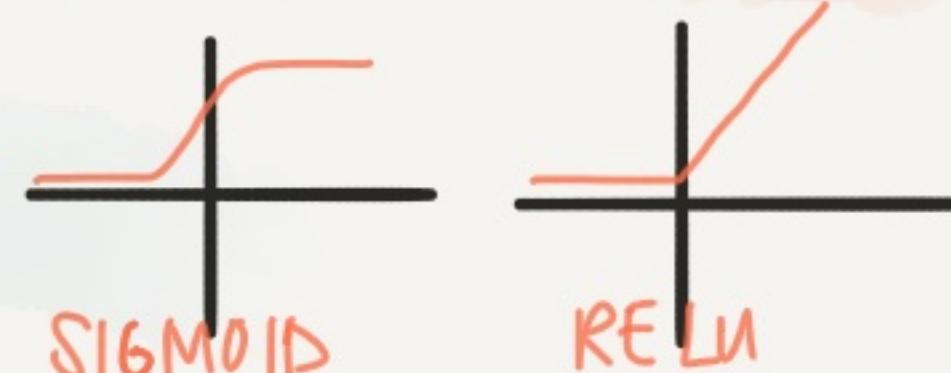
UNSTRUCTURED

HUMANS ARE GOOD
AT THIS

WHY NOW?

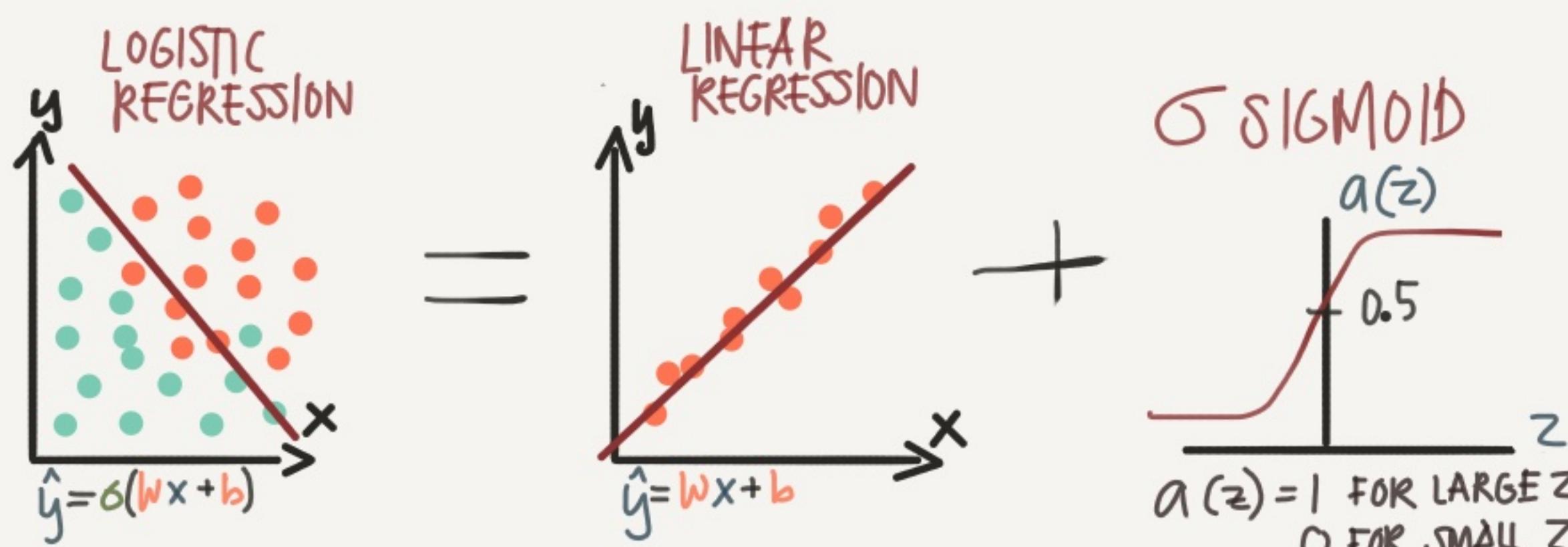
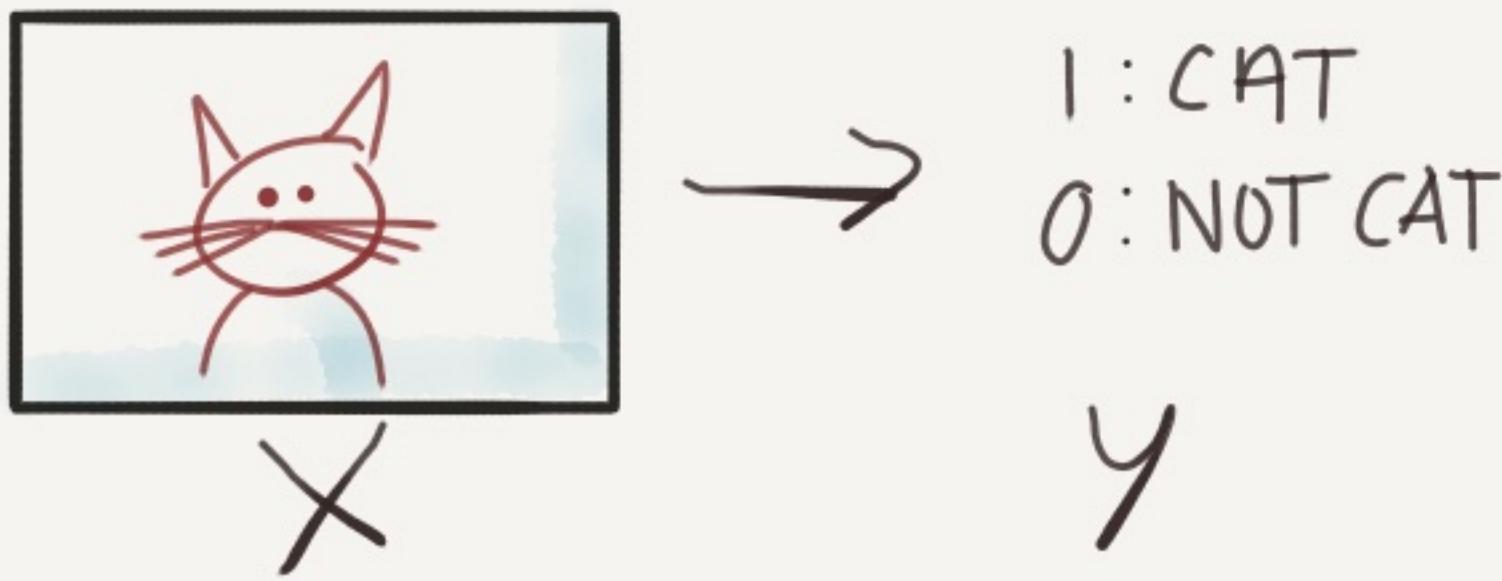


FASTER COMPUTATION
IS IMPORTANT TO SPEED UP
THE ITERATIVE PROCESS



© Tess Ferrandez

BINARY CLASSIFICATION



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)

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

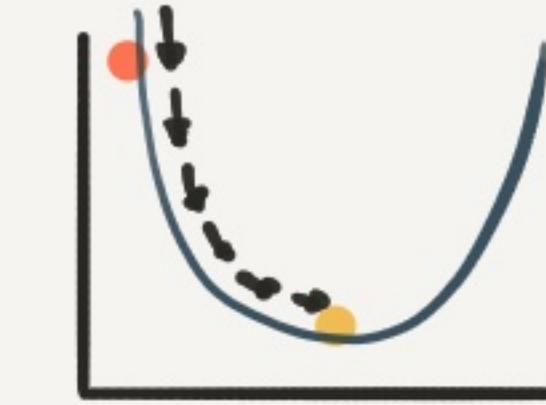
$$\text{COST} = J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

COST = LOSS FOR THE ENTIRE DATASET

LOGISTIC REGRESSION

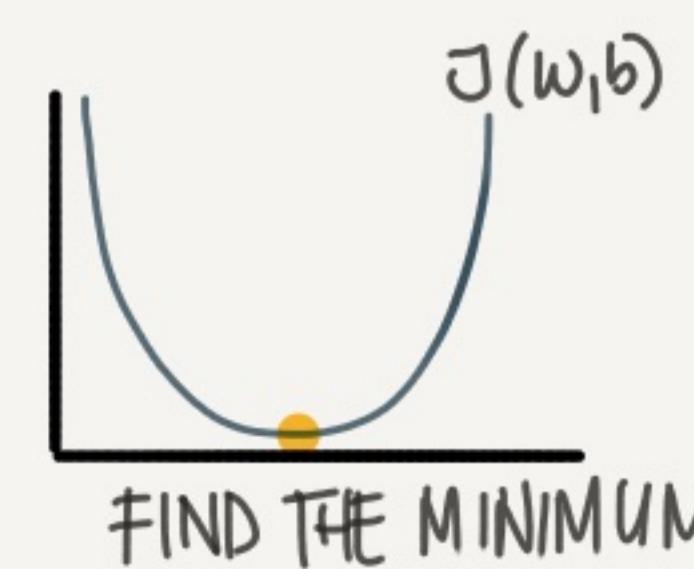
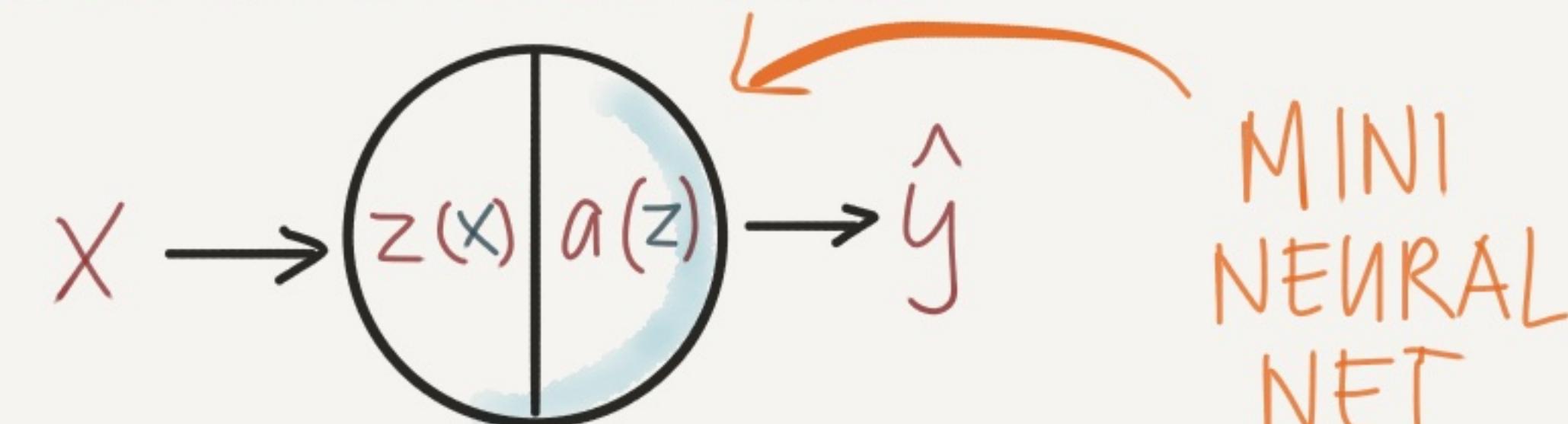
AS A NEURAL NET

FINDING THE MINIMUM WITH GRADIENT DESCENT



1. FIND THE DOWNSHILL DIRECTION (USING DERIVATIVES)
 2. WALK (UPDATE w & b) AT A α LEARNING RATE
- REPEAT UNTIL YOU REACH BOTTOM (CONVERGE)

PUTTING IT ALL TOGETHER



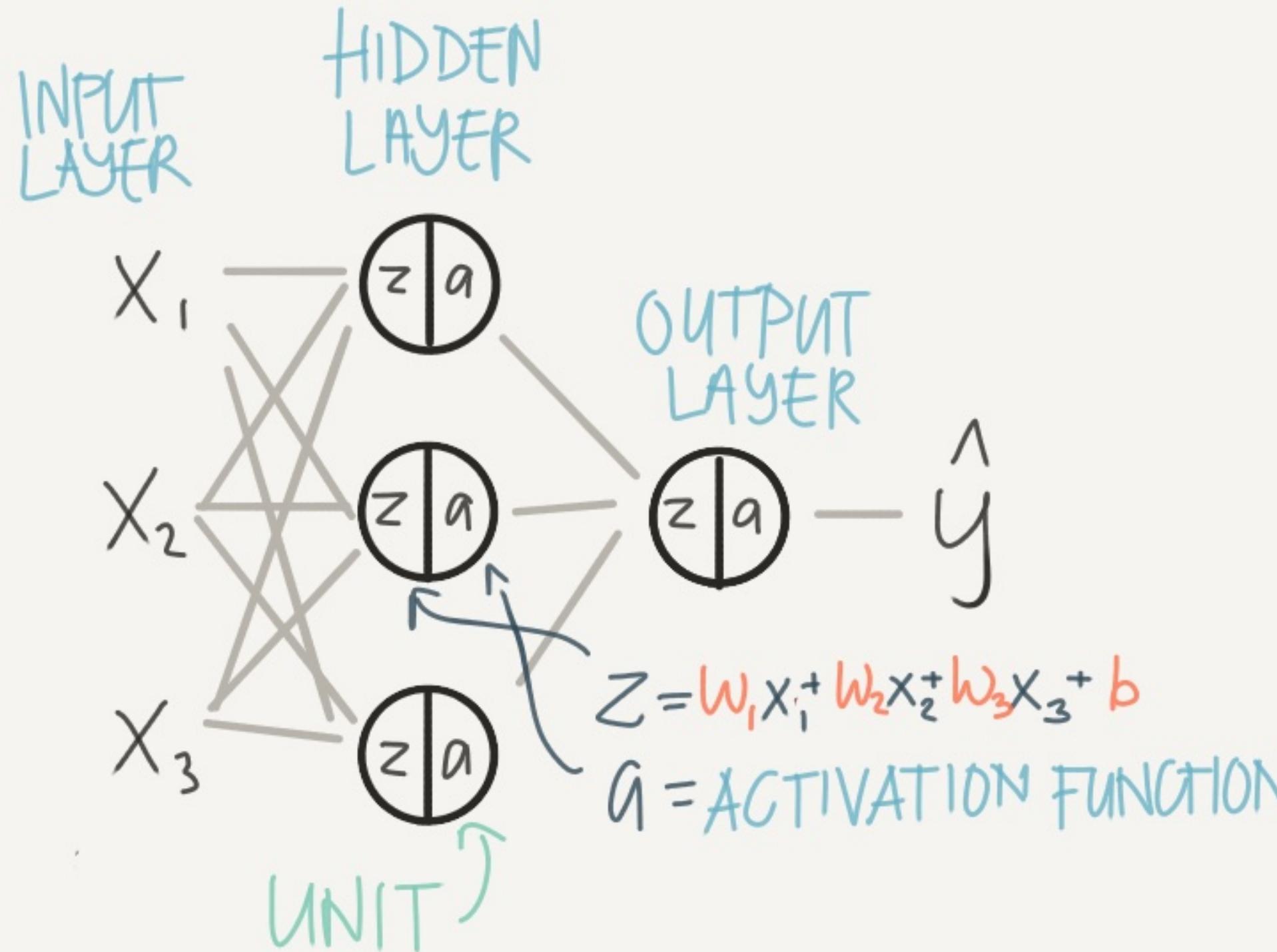
$$z(x) = w*x + b$$

$$\hat{y} = a(z) = \sigma \text{SIGMOID}(z)$$

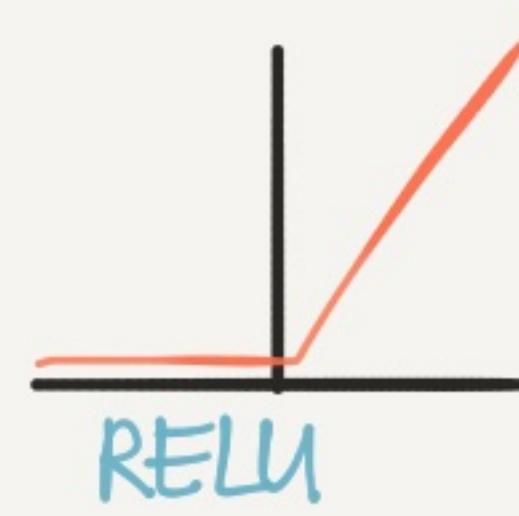
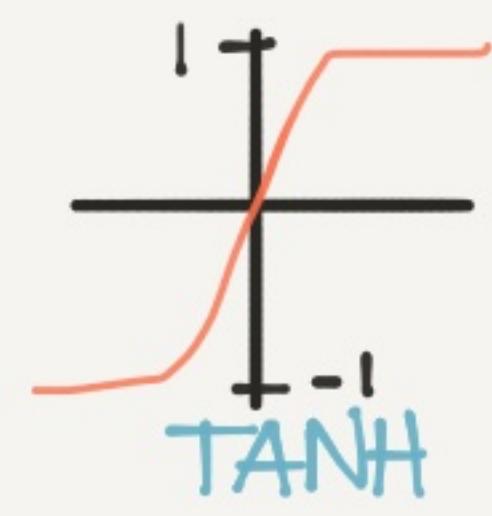
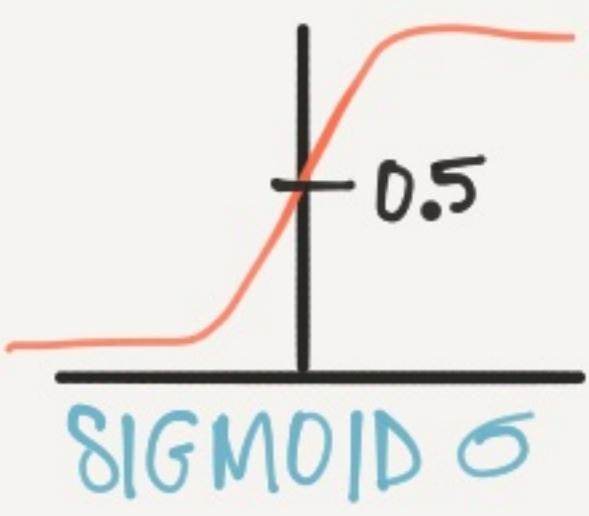
1. FORWARD PROPAGATION • CALCULATE \hat{y}
2. BACKWARD PROPAGATION • GRADIENT DESCENT + UPDATE w & b

REPEAT UNTIL IT CONVERGES

2 LAYER NEURAL NET



ACTIVATION FUNCTIONS

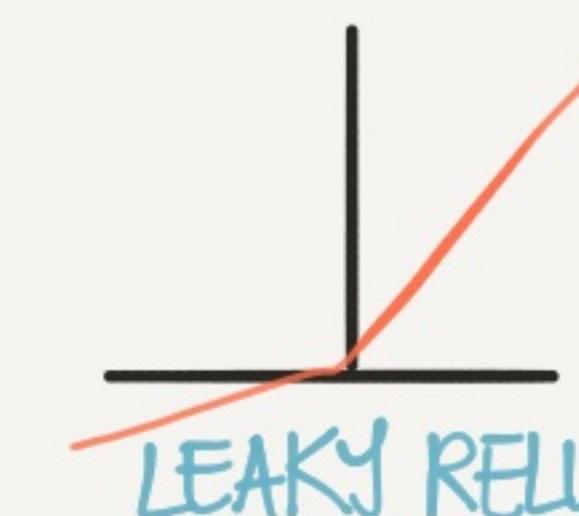


BINARY CLASSIFIER
ONLY USED FOR
OUTPUT LAYER

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

NORMALIZED
 \Rightarrow GRADIENT
DESCENT IS
FASTER

DEFAULT
CHOICE FOR
ACTIVATION
SLOPE = 1/0



AVoids UNDEF
SLOPE AT 0
BUT RARELY
USED IN PRACTICE

SHALLOW NEURAL NETS

WHY ACTIVATION FUNCTIONS?

EX. WITH NO ACTIVATION - $a = z$

$$\begin{aligned} a^{[1]} &= z^{[1]} = W^{[1]} X + b^{[1]} \\ a^{[2]} &= z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} \end{aligned}$$

LAYER 1
LAYER 2

PLUG IN $a^{[1]}$

$$\begin{aligned} a^{[2]} &= W^{[2]}(W^{[1]} X + b^{[1]}) + b^{[2]} \\ &= \underbrace{W^{[2]} W^{[1]} X}_{W' X} + \underbrace{W^{[2]} b^{[1]} + b^{[2]}}_{b'} \end{aligned}$$

LINEAR
FUNCTION

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 SD RAND $\neq 0.01$

HYPERPARAM

© Tess Ferrandez

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

SETTING UP YOUR ML APP

CLASSIC ML

100 - 1000 SAMPLES

TRAIN	DEV	TEST
60%	20%	20%

ALL FROM SAME PLACE
DISTRIBUTION

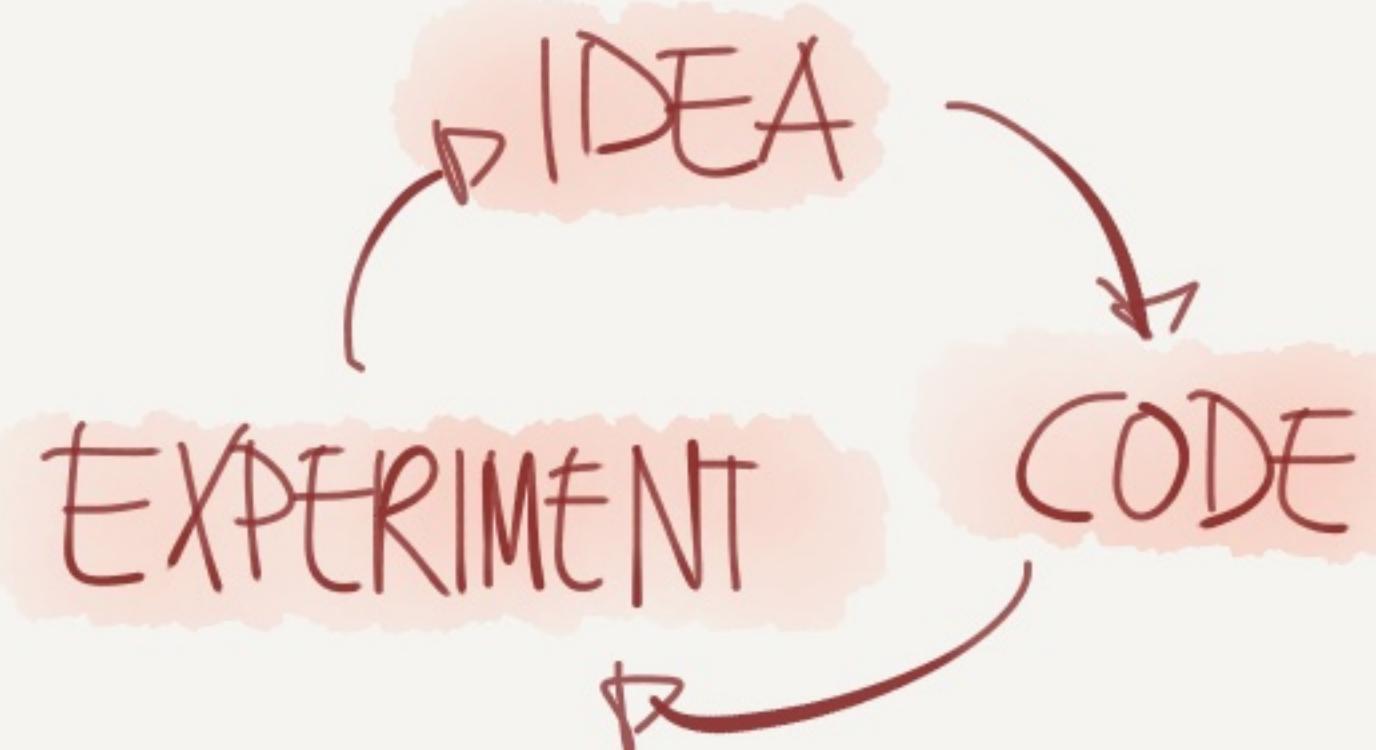
DEEP LEARNING

1M SAMPLES

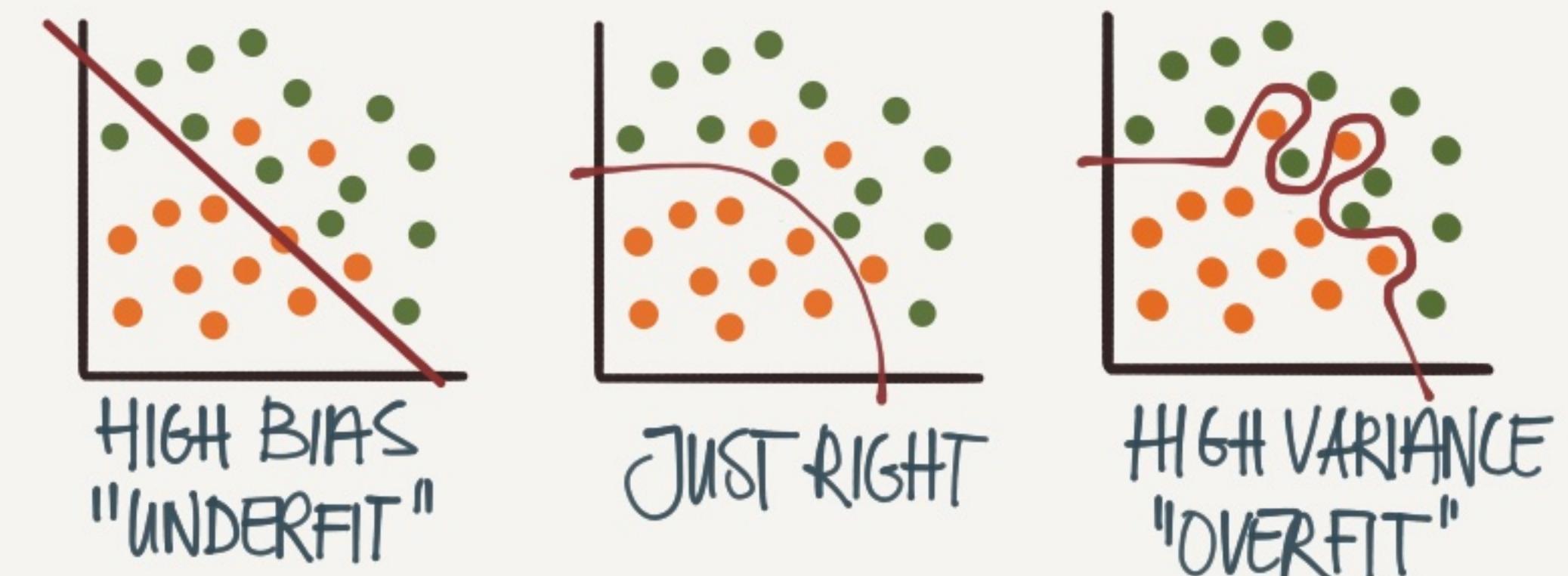
TRAIN	D	T
98%	1%	1%



 TIP
DEV & TEST SHOULD COME
FROM SAME DISTRIBUTION



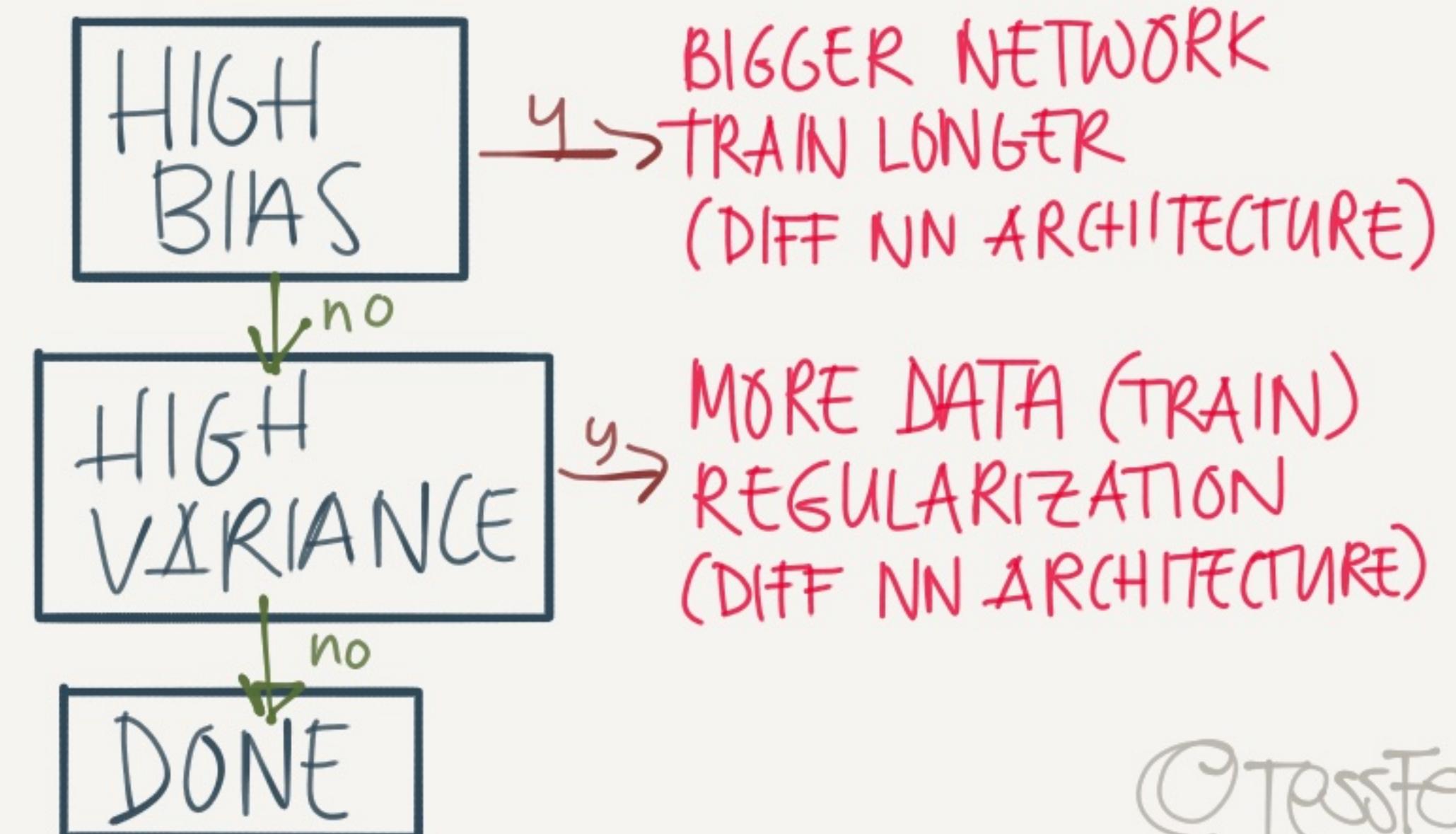
BIAS / VARIANCE



		ERROR				
		TRAIN	1%	15%	15%	0.5%
TRAIN	TEST	11%	16%	30%	1%	
		HIGH VARIANCE	HIGH BIAS	HIGH BIAS & VARIANCE	LOW BIAS & VARIANCE	

ASSUMING
HUMANS GET 0% ERROR

THE ML RECIPE



© Tess Fernandez

REGULARIZATION

PREVENTING OVERFITTING

L2 REGULARIZATION

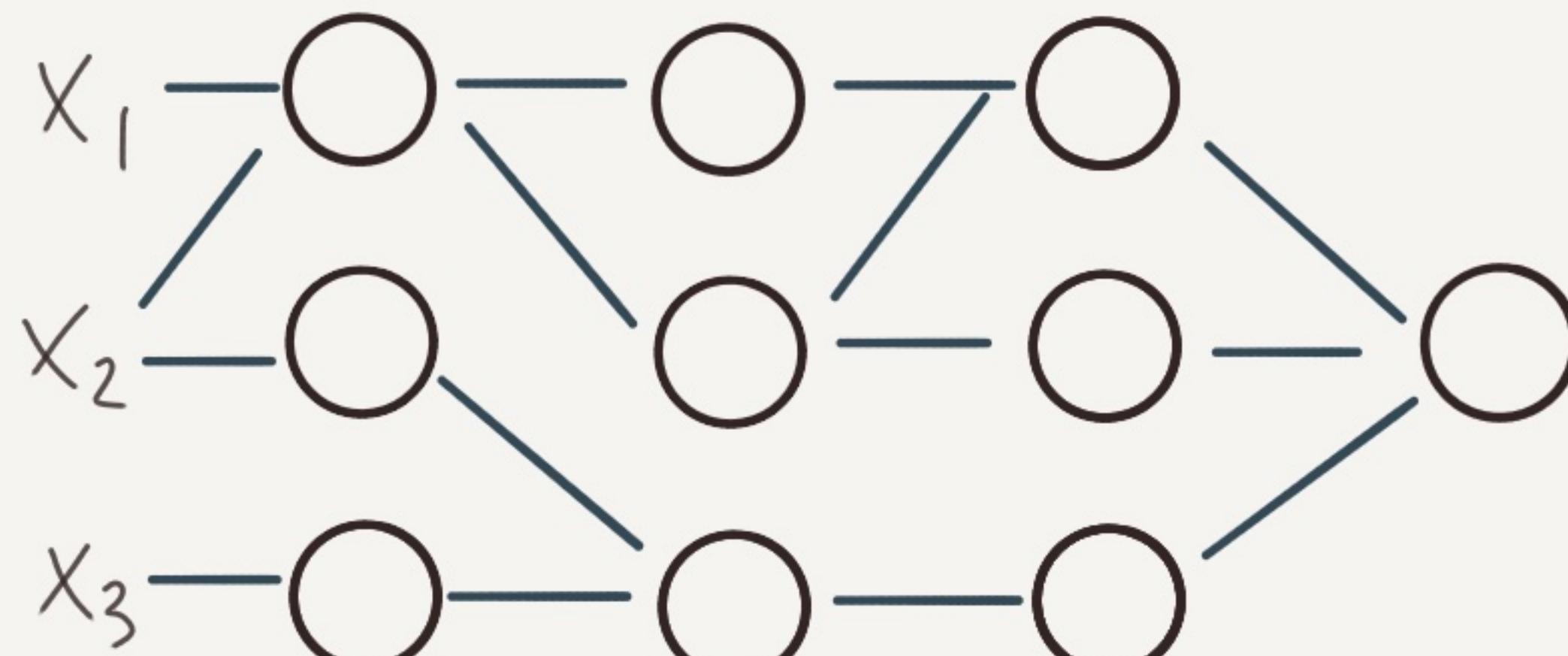
$$\text{COST: } J(w, b) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}_i, y_i) + \frac{\lambda}{2m} \|w\|_2^2$$

EUCLIDEAN NORM

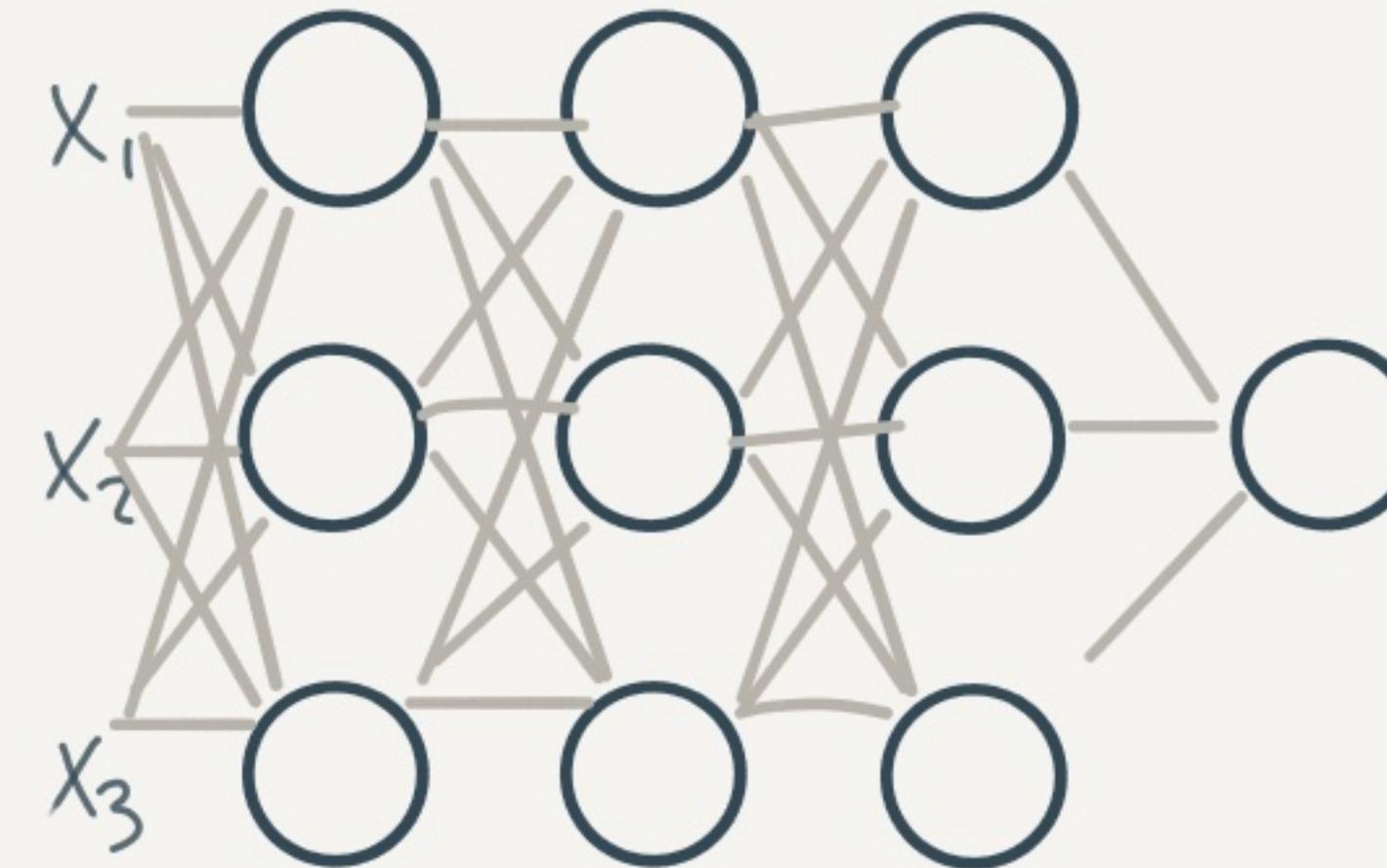
L1 REGULARIZATION

$$\text{COST: } J(w, b) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}_i, y_i) + \frac{\lambda}{m} \|w\|_1$$

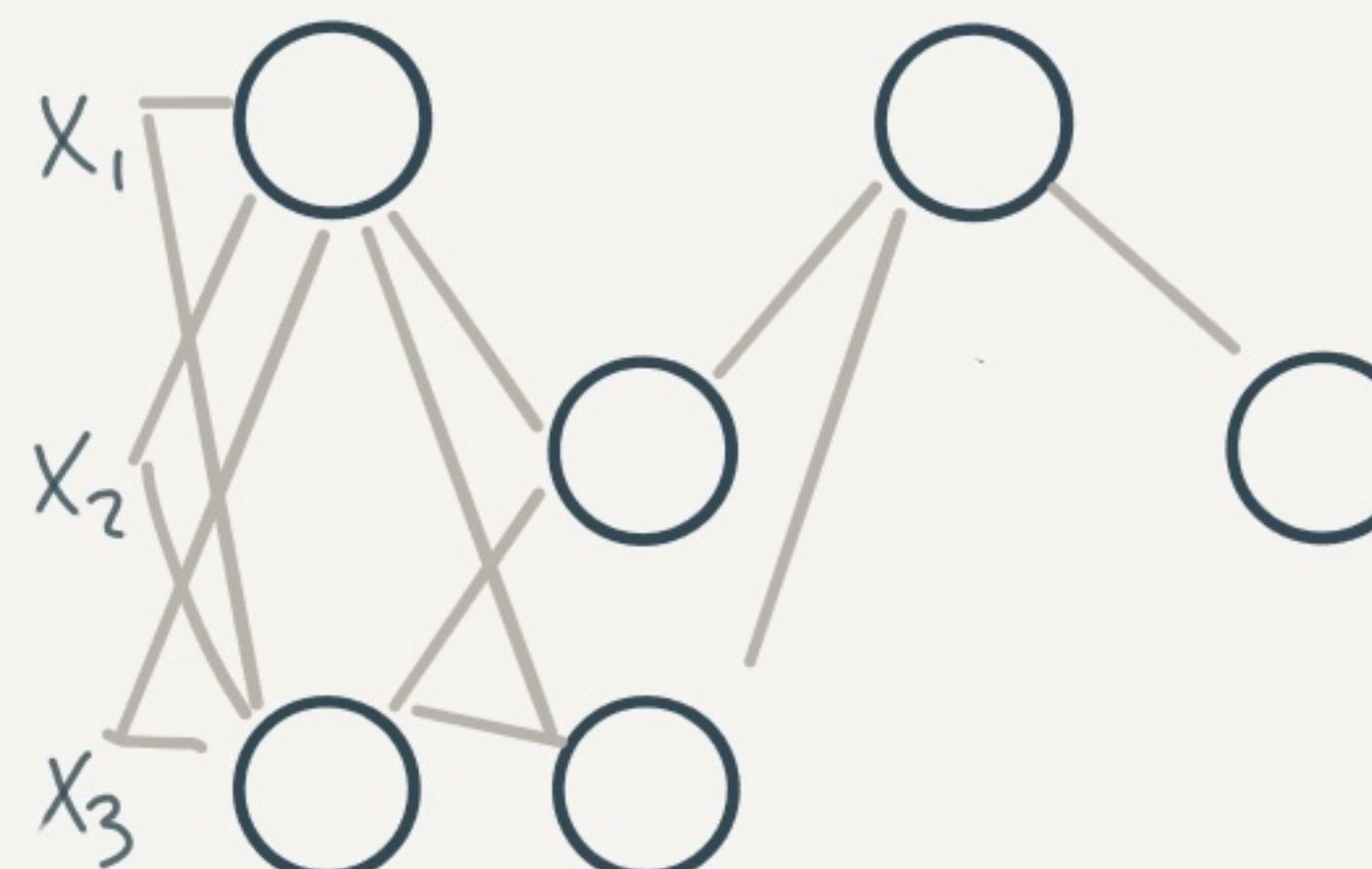
BOTH PENALIZE LARGE WEIGHTS \Rightarrow
 SOME WILL BE CLOSE TO $0 \Rightarrow$
 SIMPLER NETWORKS



DROPOUT



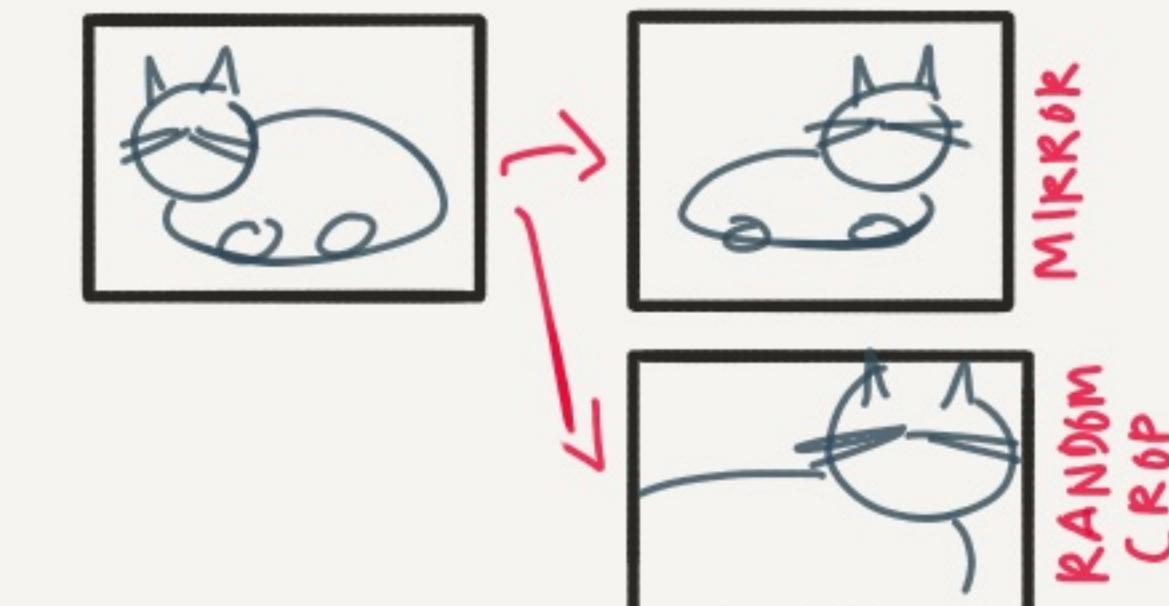
FOR EACH ITERATION i SAMPLE
 SOME NODES ARE RANDOMLY
 DROPPED (BASED ON KEEP-PROB)



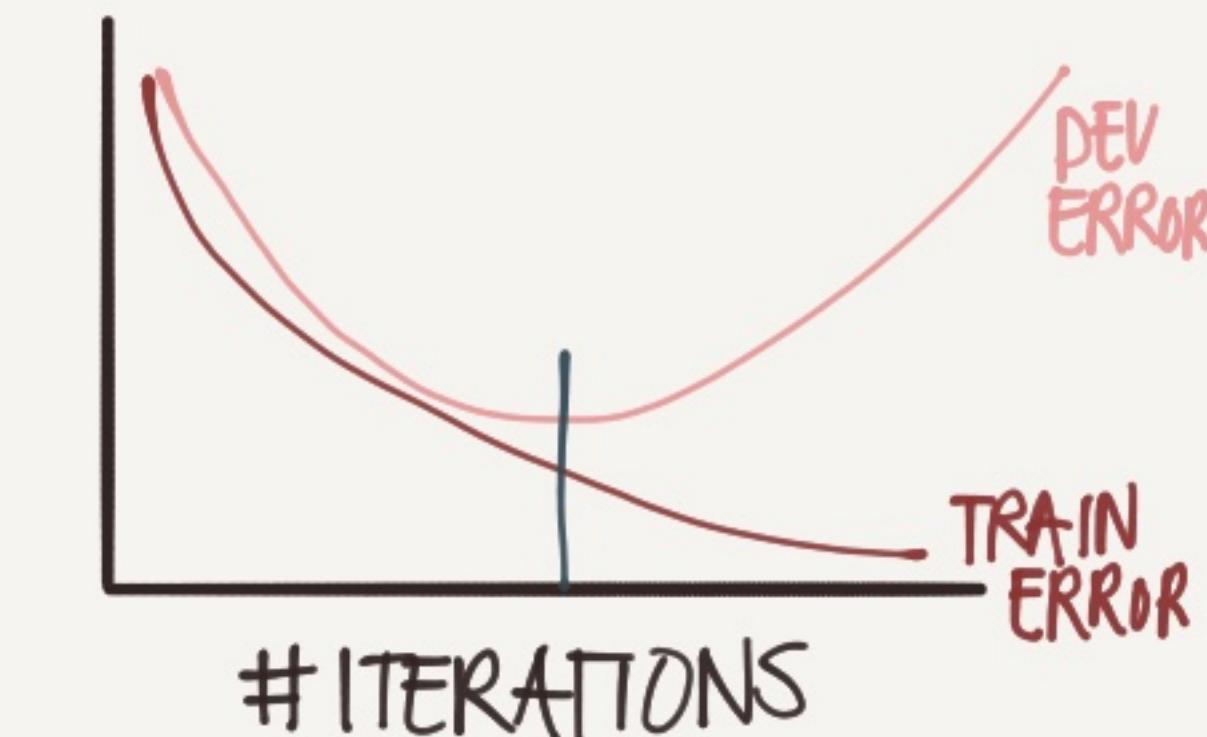
WE GET SIMPLER NWs
 & LESS CHANCE TO RELY ON
 SINGLE FEATURES

OTHER REGULARIZATION TECHNIQUES

DATA AUGMENTATION
 GENERATE NEW PICS FROM EXISTING



EARLY STOPPING

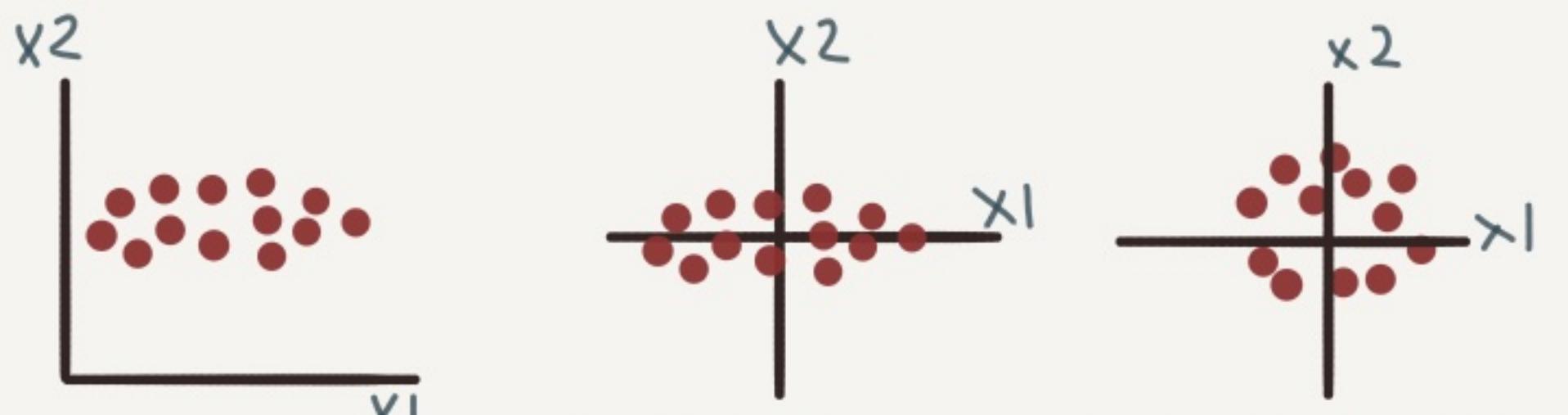


PROBLEM: AFFECTS BOTH
 BIAS & VARIANCE

OPTIMIZING

TRAINING

NORMALIZING INPUTS



STEP1: CENTER
AROUND 0,0

STEP2: SCALE
SO VARIANCE IS SAME
 $Cx -1 \rightarrow 1$

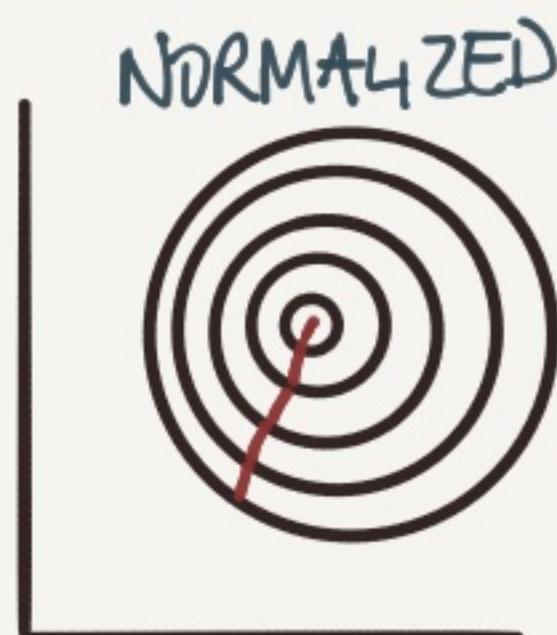


TIP
USE SAME AVG/VAR TO
NORMALIZE DEV/TEST

WHY DO WE DO THIS?



IF WE NORMALIZE, WE CAN USE A MUCH
LARGER LEARNING RATE α



DEALING WITH VANISHING/EXPLODING GRADIENTS

Ex: DEEP NW (L LAYERS)

$$\hat{y} = \underbrace{w^{[L-1]} w^{[L-2]} \cdots w^{[0]}}_{w} x + b$$

IF $w = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \Rightarrow 0.5^{L-1} \Rightarrow$ VANISHING

OR $w = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \Rightarrow 1.5^{L-1} \Rightarrow$ EXPLODING

IN BOTH CASES GRADIENT DESCENT
TAKES A VERY LONG TIME

PARTIAL SOLUTION: CHOOSE INITIAL
VALUES CAREFULLY

$$w^{[l]} = \text{rand} * \sqrt{\frac{2}{n^{l-1}}} \quad (\text{FOR RELU})$$

$$\# \text{inputs} \quad (\text{FOR TANH})$$

$$\sqrt{\frac{1}{n^{l-1}}} \quad (\text{FOR TANH})$$

SETS THE VARIANCE

GRADIENT CHECKING

IF YOUR COST DOES NOT
DECREASE ON EACH ITER
YOU MAY HAVE A
BACKPROP BUG.

GRADIENT CHECKING
APPROXIMATES THE
GRADIENTS SO YOU
CAN VERIFY CALC.

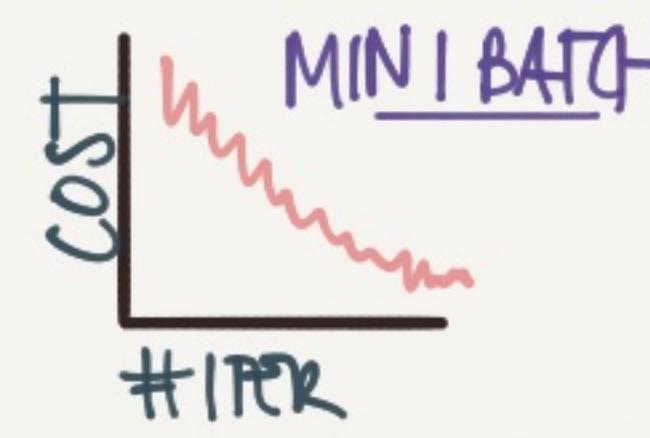
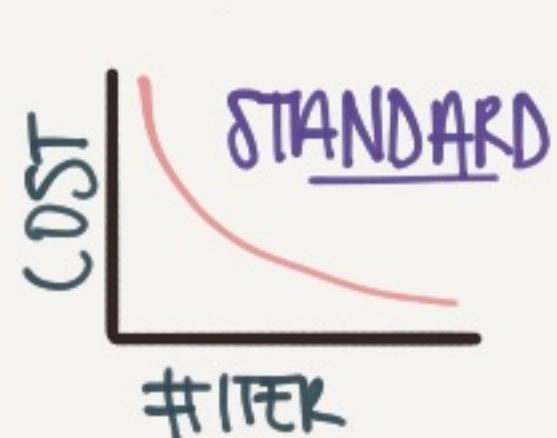
NOTE ONLY USE
WHEN DEBUGGING
SINCE IT'S SLOW

OPTIMIZATION ALGORITHMS

MINI-BATCH GRAD. DESCENT

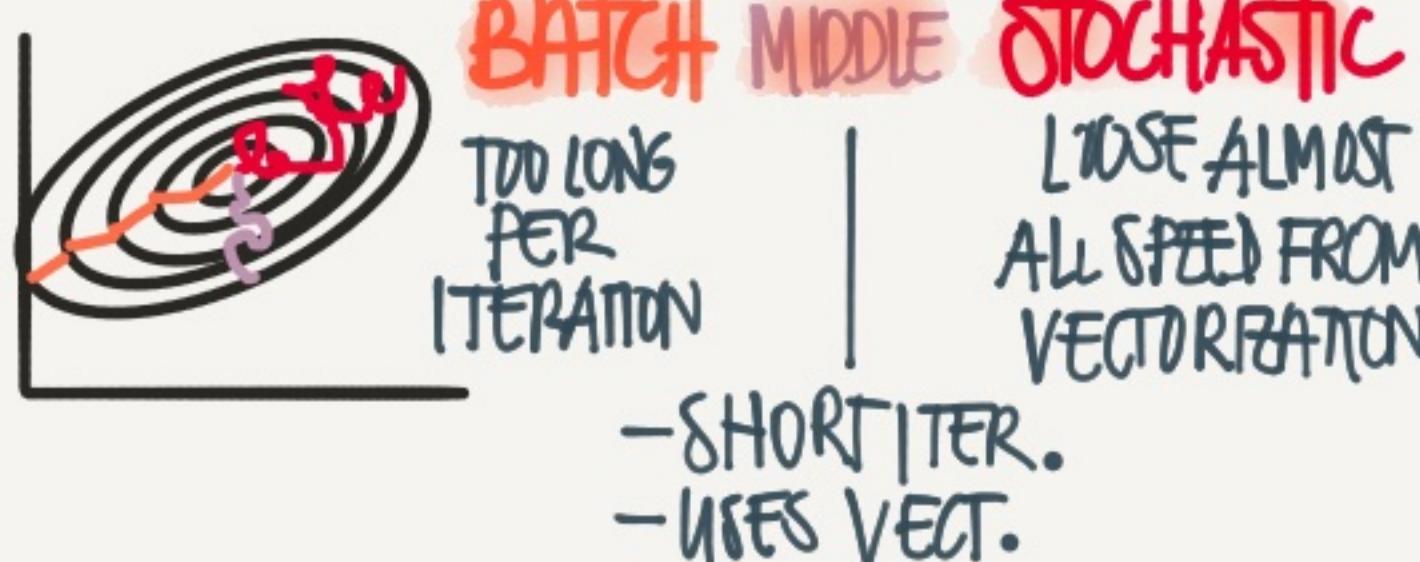


SPLIT YOUR DATA INTO MINI-BATCHES & DO GRAD DESCENT AFTER EACH BATCH THIS WAY YOU CAN PROGRESS AFTER JUST A SHORT WHILE



CHOOSING THE MINIBATCH SIZE

$\delta \text{SIZE} = m \rightarrow$ BATCH GRAD DESC.
 $\delta \text{SIZE} = 1 \rightarrow$ STOCHASTIC GRAD DESC



TIP
 IF YOU HAVE < 2000 SAMPLES
 USE $\text{SIZE}=2000$
 OTHERWISE, USE 64, 128, 256...
 SO X+Y FITS IN CPU/GPU CACHE

GRADIENT DESCENT W. MOMENTUM



WE WANT TO REDUCE OSCILLATION \updownarrow SO WE GET TO THE GOAL FASTER

SOLUTION: SMOOTH OUT THE CURVE BY TAKING AN EXPONENTIALLY WEIGHTED AVERAGE OF THE DERIVATIVES (i.e. LAST ONE HAS MORE IMPORTANCE)

RMSProp - ROOT MEAN SQUARED



NORMALIZE GRADIENT USING A MOVING AVG.

$$S_{dw} = \beta S_{dw} + (1-\beta) dw^2$$

$$S_{db} = \beta S_{db} + (1-\beta) db^2$$

$$w = w - \alpha \frac{dw}{\sqrt{S_{dw}}} \quad b = b - \alpha \frac{db}{\sqrt{S_{db}}}$$

ADAM OPTIMIZATION

COMBO OF GD w/ MOMENTUM & RMSProp

LEARNING RATE DECAY

IDEA: USE A LARGE α IN THE BEGINNING THEN DECREASE AS WE GET CLOSER TO GOAL

$$\text{OPTION 1: } \alpha = \frac{1}{1 + \text{DECAYRATE} \cdot \text{EPOCH}} \alpha_0$$

$$\text{EXponential: } \alpha = 0.95^{\text{EPOCH}} \alpha_0$$

$$\text{OPTION 3: } \alpha = \frac{k}{\sqrt{\text{EPOCH}}} \alpha_0$$

$$\text{OPTION 4: } \alpha = \frac{k}{\sqrt{t}} \alpha_0$$

$$\text{OPTION 5: DISCRETE STAIRCASE}$$

$$\text{OPTION 6: MANUAL}$$

EPOCH = 1 PASS THROUGH THE DATA

HYPERPARAM TUNING

WHICH HYPERPARAMS ARE MOST IMPORTANT?

α LEARNING RATE

HIDDEN UNITS

MINIBATCH SIZE

β MOMENTUM, TURN = 0.9

LAYERS

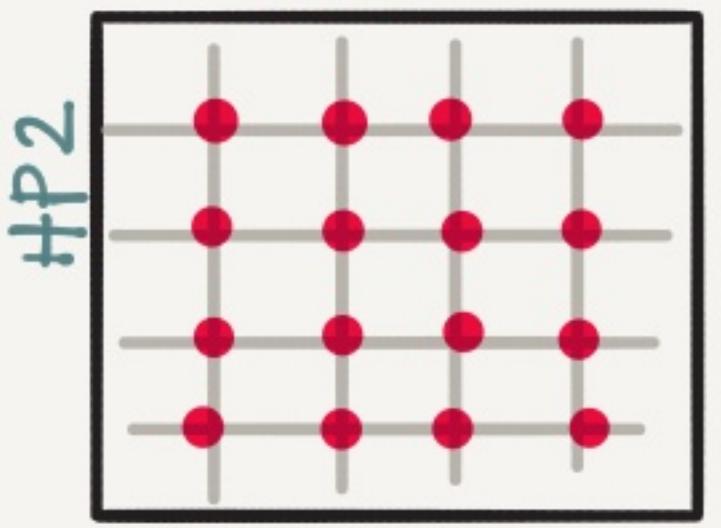
LEARNING RATE DECAY

$\beta_1 = 0.9 \quad \beta_2 = 0.999 \quad \epsilon = 10^{-8}$ (ADAM)

TESTING VALUES

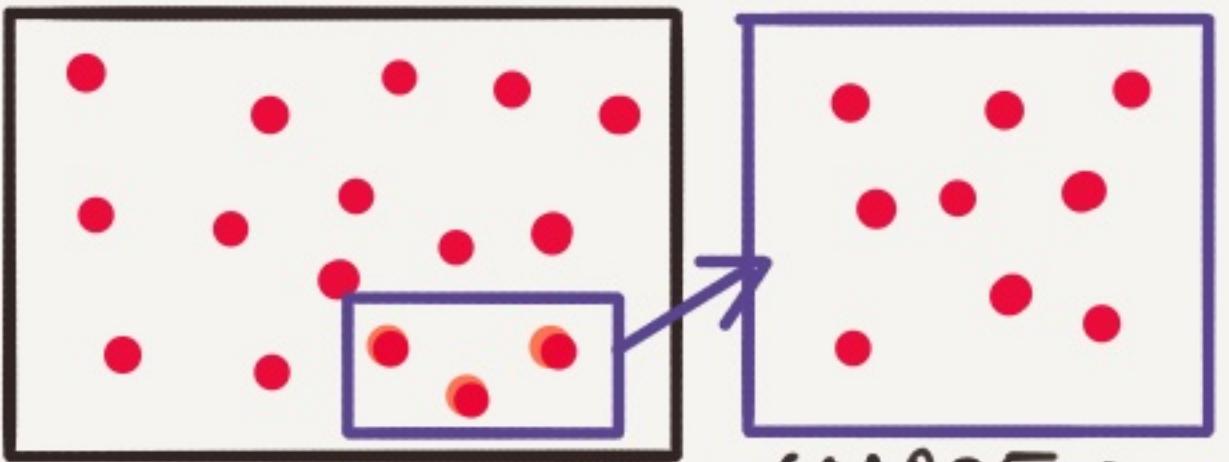
CLASSIC ML

HP1



GRID SEARCH

SOLUTION



RANDOM SEARCH + COARSE \rightarrow DENSE

PROBLEM: ONE ITERATION TAKES A LONG TIME & IN 16 GO'S WE HAVE ONLY TRIED 4 α - BUT 4 DIFF ϵ

NOT AS IMPORTANT

MY PANDA IS ACTUALLY A MIS-CATEGORIZED CAT BECAUSE I CAN'T DRAW PANDAS



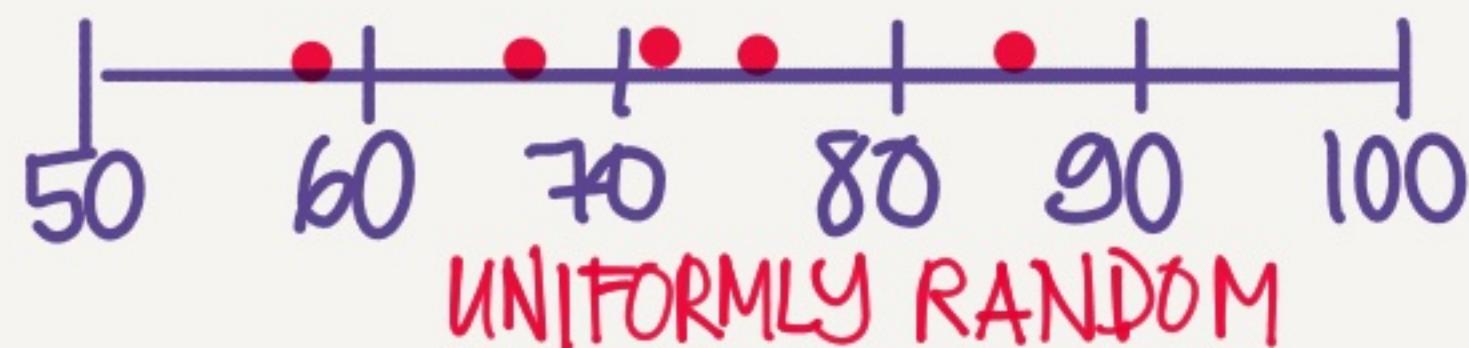
BABYSIT ONE MODEL & TUNE



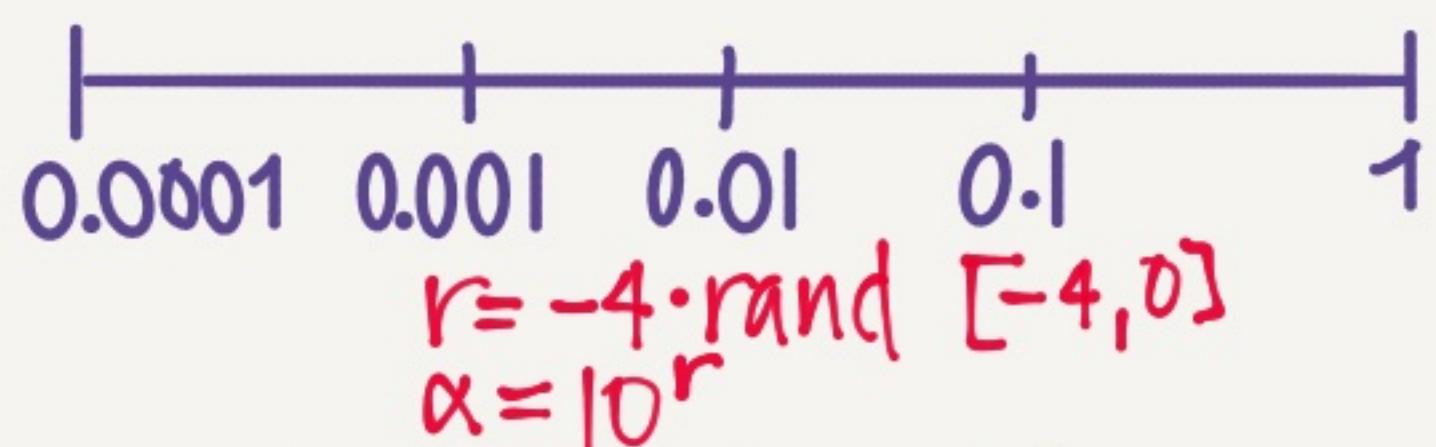
SPAWN LOTS OF MODELS W DIFF HP
GOOD IF YOU HAVE LOTS OF SHARE COMP POWER

USE AN APPROPRIATE SCALE

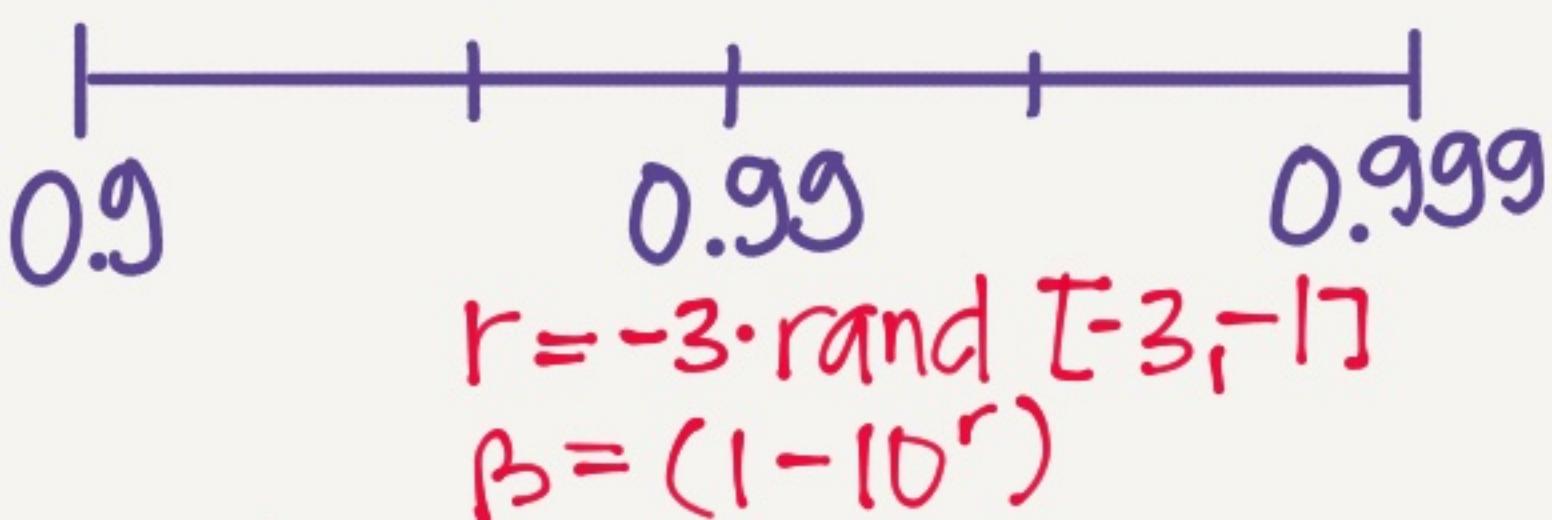
HIDDEN UNITS



α LEARNING RATE



β EXP WEIGHT AVE



TIP
RE-EVALUATE YOUR HYP. PARAMS EVERY FEW MONTHS

PANDA VS CAVIAR

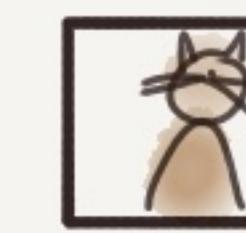
MISC. EXTRAS

BATCH NORMALIZATION

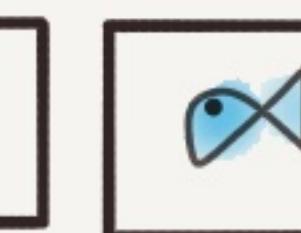
NORMALIZE LAYER OUTPUT

- SPEEDS UP TRAINING
- MAKES WEIGHTS DEEPER IN NW MORE ROBUST (COVARIATE SHIFT)
- SIGHT REGULARIZING EFFECT

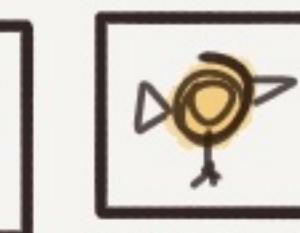
MULTICLASS CLASSIFIC.



CAT



#FISH



BABY CHICK



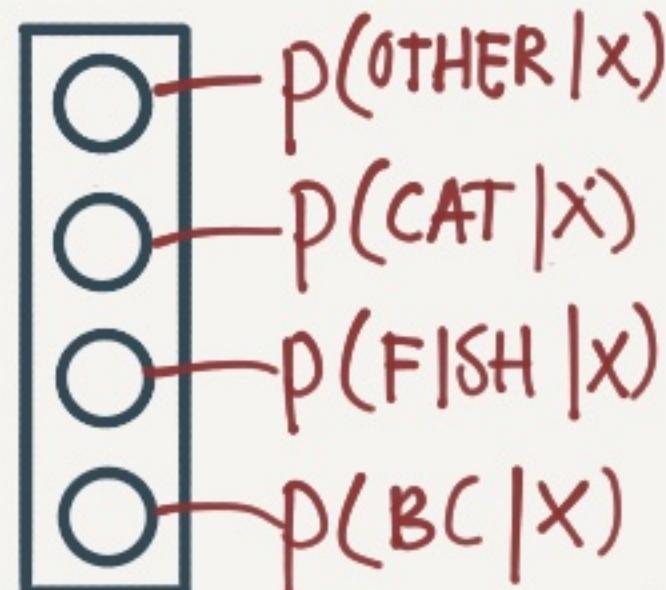
OTHER

$$C = \# \text{CLASSES} = 4$$

SOFTMAX ACTIVATION

$$t = e^{(z^{[i]})}$$

$$a^{[i]} = \frac{t}{\sum t_i}$$



SUM: 1

$$\text{EX: } z^{[i]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ -3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^{-3} \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 2.1 \end{bmatrix}$$

$$\Rightarrow a^{[i]} = \frac{t}{176.3} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.02 \\ 0.114 \end{bmatrix} = 11.4\% \text{ PROB IT'S A BABY CHICK}$$

© TessFerrandez

STRUCTURING YOUR ML PROJECTS

SETTING YOUR GOAL

* GOAL SHOULD BE A SINGLE #

	PRECISION	RECALL	
A	95%	90%	IS A OR B BEST?
B	98%	85%	

	PRECISION	RECALL	F1
A	95%	90%	92.4%
B	98%	85%	91%

F1 = HARMONIC MEAN BETW.
RECALL & PRECISION

* DEFINE OPTIMIZING VS
SATISFYING METRICS

	ACCURACY	RUNTIME
A	90%	80ms
B	92%	95ms
C	95%	1500ms

MAXIMIZE ACC.
GIVEN TIME < 100ms

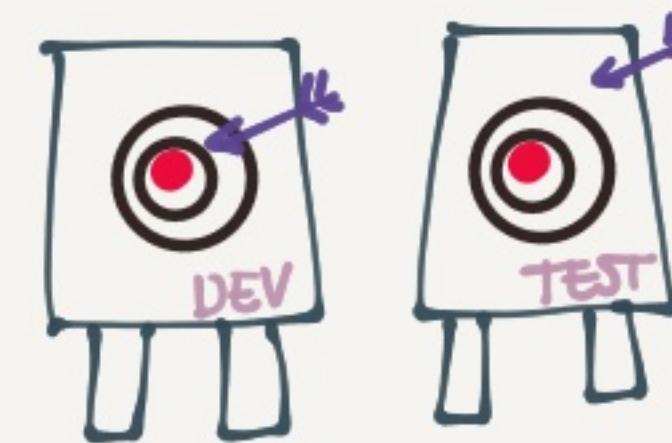
ACCURACY =
OPTIMIZING
RUNTIME =
SATISFYING

SELECTING YOUR DEV/TEST SETS

DATA

US
UK
EUROPE
S.AM
INDIA
CHINA
AUST.

OPTION 1:
DEV = UK, US, EUR
TEST = REST



IF DEV & TEST ARE DIFF
& WE OPTIMIZE FOR DEV
WE WILL MISS THE TEST TARGET

HUMAN LEVEL PERF



WHY DOES ACC
SLOW DOWN WHEN
WE SURPASS HUMAN
LEVEL PERF?

TYPICAL HUMAN 3%
TYPICAL DOCTOR 1%
EXPERIENCED DR. 0.7%
TEAM OF EXP DRs. 0.5%

↑ HUMAN LEV PERF
(PROXY FOR BAYES)

- OFTEN CLOSE TO BAYES
- A HUMAN CAN NO LONGER HELP IMPROVE (INSIGHTS)
- DIFFICULT TO ANALYSE BIAS/VARIANCE

CAT CLASSIFICATION

	A	B	BLURRY
HUMAN	1%	7.5%	
TRAIN ERR	8%	8%	AVOIDABLE BIAS
DEV ERR	10%	10%	VARIANCE

FOCUS ON BIAS FOCUS ON VARIANCE

HUMAN TRAIN BIGGER NETW.
| AVOIDABLE BIAS } TRAIN LONGER/BETTER OPT. (RMSprop, ADAM)
TRAIN | ALSO CHANGE NN ARCH OR HYPERPARAMS
| VARIANCE } MORE DATA (TRAIN)
DEV } REGULARIZATION NN ARCHITECTURE

	A	B	
HUMAN	0.5	0.5	AVOIDABLE BIAS
TRAIN ERR	0.6	0.3	VARIANCE
DEV ERR	0.8	0.4	
AVOID. BIAS	0.1	?	DON'T KNOW IF WE OVERFIT OR IF WE'RE CLOSE TO BAYES

OPTIONS TO PROCEED ARE UNCLEAR

ERROR ANALYSIS

YOU HAVE 10% ERRORS, SOME ARE DOGS MIS-CLASSIFIED AS CATS. SHOULD YOU TRAIN ON MORE DOG PICS?

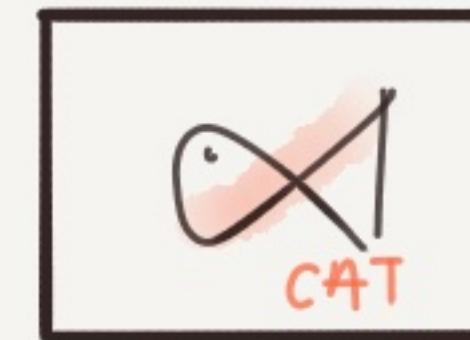
1. PICK 100 MIS-LABLED
2. COUNT ERROR REASONS

Dog	Blurry	Insta Filter	Big Cat	...
1	1		1	
2			1	
3		1		
...				
100		1		
5	...			

5% OF ALL ERRORS

FOCUSING ON DOGS. THE BEST WE CAN HOPE FOR IS 9.5% ERROR

YOU FIND SOME INCORR. LABELED DATA IN THE DEV SET. SHOULD YOU FIX IT?



DL ALGORITHMS ARE PRETTY ROBUST TO RANDOM ERRORS. BUT NOT TO SYSTEMATIC ERR.
(EX. ALL WHITE CATS INCORR LABLED AS MICE)

ADD EXTRA COL. IN ERROR ANALYSIS AND USE SAME CRITERIA

NOTE IF YOU FIX DEV YOU SHOULD FIX TEST AS WELL.

FOR NEW PROJ. ·
BUILD 1ST SYSTEM QUICK & ITERATE

EX: SPEECH RECOGNITION



WHAT SHOULD YOU FOCUS ON?

NOISE
ACCENTS
FAR FROM MIKE

1. START QUICKLY DEV/TEST METRICS
2. GET TRAIN-SET
3. TRAIN
4. BIAS/VARIANCE ANAL
5. ERROR ANALYSIS
6. PRIORITIZE NEXT STEP