# INTRO TO DEEP LEARNING

## Neural Network Diagram

INPUT LAYER → HIDDEN LAYER → OUTPUT LAYER

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

## SUPERVISED LEARNING

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

## NETWORK ARCHITECTURES

$X_1$
$X_2$ → $\hat{y}$
$X_3$

STANDARD NN

CONVOLUTIONAL NN

$\hat{y}$ $\hat{y}$ $\hat{y}$
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)

IDEA → CODE → EXPERIM. → (cycle)

FASTER COMPUTATION IS IMPORTANT TO SPEED UP THE ITERATIVE PROCESS

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

SIGMOID    RELU

@TessFerrandez

# TRAIN vs DEV/TEST MISMATCH

## AVAILABLE DATA

200k PRO CAT PICS FROM INTERNET

10k BLURRY CAT PICS FROM APP — WHAT WE CARE ABT

### HOW DO WE SPLIT → TRAIN/DEV/TEST?

**OPTION 1: SHUFFLE ALL**

| 205k (TRAIN) | D | T |
|---|---|---|

2.5k

**PROBLEM:** DEV/TEST IS NOW MOSTLY WEB IMG (NOT REPRES. OF END SCENARIO)

**SOLUTION:** LET DEV/TEST COME FROM APP · THEN SHUFFLE 5k OF APP PICS w̄ WEB FOR TRAIN

| 205k | 2.5 | 2.5 |
|---|---|---|

WEB + APP — APP APP

## BIAS & VARIANCE w MISMATCHED TRAIN/DEV

| HUMANS | ~0% |
|---|---|
| TRAIN | 1% |
| DEV ERR | 10% |

IS THIS DIFF DUE TO THE MODEL NOT GENERALIZING OR IS DEV DATA MUCH HARDER
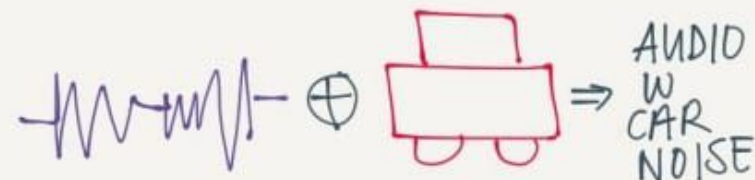
**A: CREATE A TRAIN·DEV SET THAT WE DON'T TRAIN ON**

| TRAIN | FD | D | T |
|---|---|---|---|

|  | A | B | C | D |
|---|---|---|---|---|
| TRAIN | 1% | 1% | 10% | 10% |
| TRAIN·DEV | 9% | 15% | 11% | 11% |
| DEV | 10% | 10% | 12% | 20% |
|  | VARIANCE | TRAIN/DEV MISMATCH | BIAS | BIAS + DATA MISMATCH |

## ADDRESSING DATA MISMATCH

EX. CAR GPS · TRAINING DATA IS 10.000H OF GENERAL SPEECH DATA

1. CARRY OUT MANUAL ERROR ANALYSIS TO UNDERSTAND THE DIFFERENCE (EX NOISE, STREET NUMBERS)

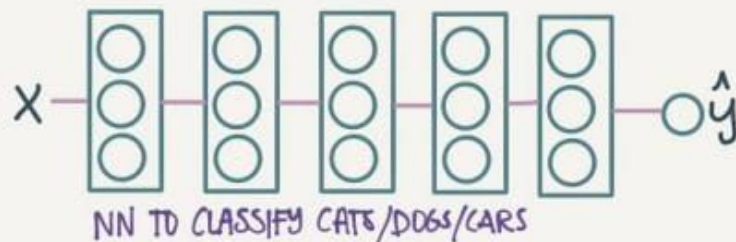2. TRY TO MAKE TRAIN MORE SIMILAR TO DEV OR GATHER MORE DEV·LIKE TRAIN·DATA

⊕ ⟹ AUDIO w CAR NOISE

**NOTE** BE CAREFUL · IF YOU ONLY HAVE 1 HR OF CAR NOISE & APPLY IT TO 10K HR SPEECH YOU MAY OVERFIT TO THE CAR NOISE

@TessFerrandez

# EXTENDED LEARNING

## TRANSFER LEARNING

PROBLEM: YOU WANT TO CLASSIFY SOME MEDICAL IMG. YOU HAVE AN NN THAT CLASSIFIES CATS

$X$ ———[ooo][ooo][ooo][ooo][ooo]——— $o\,\hat{y}$

NN TO CLASSIFY CATS/DOGS/CARS

**OPTION 1:** YOU ONLY HAVE A FEW RADIOLOGY IMAGES

SOLUTION: INIT W. WEIGHTS FROM CAT NN ONLY RETRAIN LAST LAYER(S) ON RADIOLOGY IMAGES

**OPTION 2** YOU HAVE LOTS OF RADIOLOGY IMG.

SOLUTION: INIT WITH WEIGHTS FROM CAT NN RETRAIN ALL LAYERS

TESS NOTES

THIS IS MICROSOFT CUSTOM VISION

## MULTI TASK LEARNING

TRAINING ON MULT. TASKS AT ONCE

DETECT
CAR
STOP SIGN
PEDESTR.
TRAFFIC LIGHT

$\hat{y}$ →
1
1
0

STOP

UNLIKE SOFMAX · MANY THINGS CAN BE TRUE

$$COST: J(\omega, b) = \frac{1}{m} \sum_{i=1}^{m} \sum_{j=1}^{m} \mathcal{L}(\hat{y}_j^{(i)}, y_j^{(i)})$$
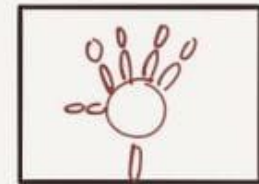
SUMMING OVER ALL OUTP OPTIONS

WE COULD HAVE JUST TRAINED 4 NNS INSTEAD BUT... MT LEARNING MAKES SENSE WHEN

A. THE LEARNING DATA YOU HAVE FOR THE DIFF TASKS IS QUITE SIMILAR — & THE AMOUNTS (EG. 1K CARS, 1K STOP SIGNS)

B. THE SUM OF THE DATA ALLOWS YOU TO TRAIN A BIG ENOUGH NN TO DO WELL ON ALL TASKS

IN REALITY TRANSFER LEARNING IS USED MORE OFTEN

## END-TO-END LEARNING

FROM X-RAY OF CHILDS HAND TELL ME THE AGE OF THE CHILD

TYPICAL SOLN:

1. LOCATE BONES TO FIND LENGTHS USING ML

2. TRAIN MODEL TO PREDICT AGE BASED ON BONE LENGTH

### END-TO-END

RADIOLOGY IMG ——→ CHILD AGE

PROS:
- LET'S THE DATA SPEAK (MAYBE IT FINDS RELATIONS WE'RE UNAWARE OF)
- LESS HAND-DESIGNING OF COMPONENTS NEEDED

CONS:
- NEEDS LARGE AMTS OF LABLED DATA $(x \to y)$
- EXCLUDES POTENTIALLY USEFUL HAND-MADE COMPONENTS

@Tessferrandez

# 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

---

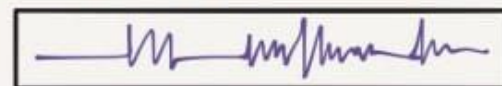YOU FIND SOME INCORR. LABLED DATA IN THE DEV SET. SHOULD YOU FIX IT?

🐟 CAT

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
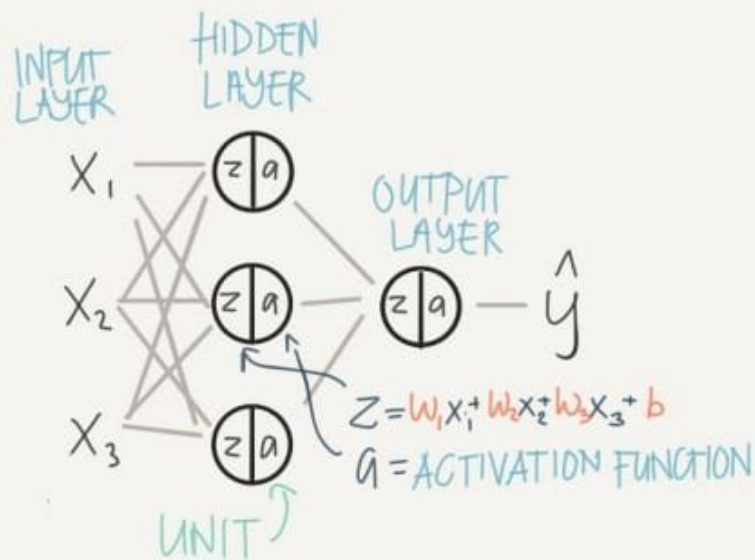
EX: SPEECH RECOGNITION

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

## 2 LAYER NEURAL NET

INPUT LAYER    HIDDEN LAYER

$x_1$

OUTPUT LAYER

$x_2$ — $\hat{y}$

$x_3$

$z = W_1 x_1 + W_2 x_2 + W_3 x_3 + b$

$a = $ ACTIVATION FUNCTION

UNIT

## WHY ACTIVATION FUNCTIONS?

EX. WITH NO ACTIVATION — $a = z$

$$a^{[1]} = z^{[1]} = W^{[1]} x + b^{[1]} \quad \text{LAYER 1}$$
$$a^{[2]} = z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} \quad \text{LAYER 2}$$

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

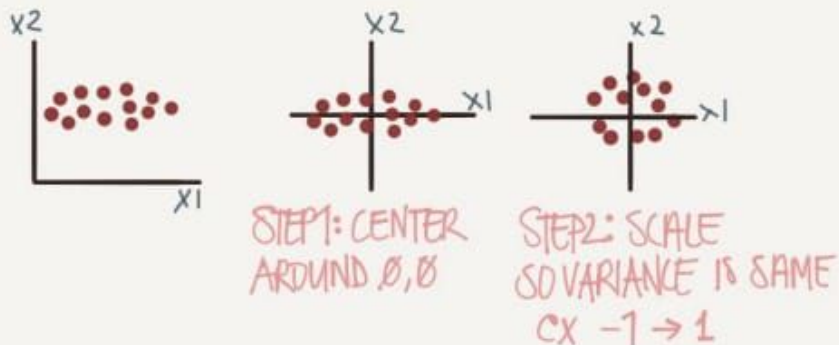← HYPERPARAM

@TessFerrandez

# OPTIMIZING
## TRAINING

## NORMALIZING INPUTS



STEP 1: CENTER AROUND 0,0

STEP 2: SCALE SO VARIANCE IS SAME
$cx \; -1 \rightarrow 1$

**TIP**
USE SAME AVG/VAR TO NORMALIZE DEV/TEST

### WHY DO WE DO THIS?



UNNORMALIZED          NORMALIZED

IF WE NORMALIZE, WE CAN USE A MUCH LARGER LEARNING RATE $\alpha$

---

## DEALING WITH VANISHING/EXPLODING GRADIENTS

Ex: DEEP NW (L LAYERS)

$$\hat{y} = W^{[L-1]} W^{[L-2]} \cdots W^{[1]} x + b$$

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

$$W^{[1]} = \text{rand} * \sqrt{\frac{2}{n^{[L-1]}}} \quad \left(\begin{array}{c}\text{FOR}\\\text{RELU}\end{array}\right)$$

#inputs

XAVIER $\sqrt{\frac{1}{n^{[L-1]}}}$ $\left(\begin{array}{c}\text{FOR}\\\text{TANH}\end{array}\right)$

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

@TessFerrandez

# STRUCTURING
## YOUR ML PROJECTS

## SETTING YOUR GOAL

### ★ A GOAL SHOULD BE A SINGLE #

PRECISION, RECALL

|   | | |
|---|---|---|
| A | 95% | 90% |
| B | 98% | 85% |

IS A OR B BEST?

PRECISION, RECALL, F1

|   | | | |
|---|---|---|---|
| A | 95% | 90% | 92.4% |
| B | 98% | 85% | 91% |

A IS BEST

F1 = HARMONIC MEAN BETW. RECALL & PRECISION

### ★ DEFINE OPTIMIZING VS SATISFICING METRICS

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

MAXIMIZE ACC. GIVEN TIME ≤ 100ms

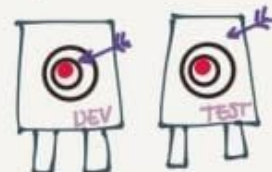ACCURACY = OPTIMIZING
RUNTIME = SATISFICING

## SELECTING YOUR DEV/TEST SETS

### DATA

US
UK
EUROPE
S.AM
INDIA
CHINA
AUST.

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

DEV        TEST

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

### HUMAN LEVEL PERF

BAYES OPTIMAL ERROR
HUMAN LEVEL PERF

MEDICAL IMG CLASS
TYPICAL HUMAN    3%
TYPICAL DOCTOR   1%
EXPERIENCED DR.  0.7%
TEAM OF EXP DRS. 0.5%
↑ HUMAN LEV PERF (PROXY FOR BAYES)

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

1. OFTEN CLOSE TO BAYES
2. A HUMAN CAN NO LONGER HELP IMPROVE (INSIGHTS)
3. DIFFICULT TO ANALYSE BIAS/VARIANCE

## CAT CLASSIFICATION

IS BLURRY

|   | A | B |
|---|---|---|
| HUMAN | 1% | 7.5% |
| TRAIN ERR | 8% | 8% |
| DEV ERR | 10% | 10% |

← AVOIDABLE BIAS
← VARIANCE

FOCUS ON BIAS | FOCUS ON VARIANCE

HUMAN } AVOIDABLE BIAS
TRAIN } VARIANCE
DEV

AVOIDABLE BIAS:
- TRAIN BIGGER NETW.
- TRAIN LONGER/BETTER OPT. ALGOS (RMSProp, ADAM)
- CHANGE NN ARCH OR HYPERPARAMS

VARIANCE:
- MORE DATA (TRAIN)
- REGULARIZATION
- NN ARCHITECTURE

|   | A | B |
|---|---|---|
| HUMAN | 0.5 | 0.5 |
| TRAIN ERR | 0.6 | 0.3 |
| DEV ERR | 0.8 | 0.4 |
| AVOID. BIAS | 0.1 | ? |

← AVOIDABLE BIAS
← VARIANCE

DON'T KNOW IF WE OVERFIT OR IF WE'RE CLOSE TO BAYES

OPTIONS TO PROCEED ARE UNCLEAR

@TessFerrandez

# HYPERPARAM
## TUNING

### WHICH HYPERPARAMS ARE MOST IMPORTANT?

$\alpha$ LEARNING RATE
\# HIDDEN UNITS
MINIBATCH SIZE
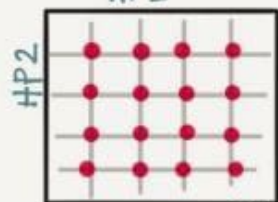$\beta$ MOMENTUM. TURN = 0.9
\# LAYERS
LEARNING RATE DECAY

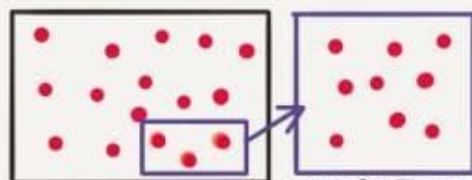$\beta_1 = 0.9$  $\beta_2 = 0.999$  $\varepsilon = 10^{-8}$ (ADAM)

### TESTING VALUES

#### CLASSIC ML

HP1
HP2

GRID SEARCH

PROBLEM: ONE ITERATION TAKES A LONG TIME & IN 16 GO'S WE HAVE ONLY TRIED 4 $\alpha$ - BUT 4 DIFF $\varepsilon$
NOT AS IMPORTANT

#### SOLUTION

RANDOM SEARCH + COARSE → DENSE

---

## USE AN APPROPRIATE SCALE

### \# HIDDEN UNITS

```
50   60   70   80   90   100
```
UNIFORMLY RANDOM

### $\alpha$ LEARNING RATE

```
0.0001  0.001  0.01   0.1    1
```
$r = -4 \cdot rand\ [-4, 0]$
$\alpha = 10^r$

### $\beta$ EXP WEIGHT AVE

```
0.9        0.99        0.999
```
$r = -3 \cdot rand\ [-3, -1]$
$\beta = (1 - 10^r)$

### TIP
RE-EVALUATE YOUR HYP-PARAMS EVERY FEW MONTHS

### PANDA VS CAVIAR

MY PANDA IS ACTUALLY A MISCLASSIFIED 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

---

## MISC. EXTRAS

### BATCH NORMALIZATION

NORMALIZE LAYER OUTPUT
- SPEEDS UP TRAINING
- MAKES WEIGHTS DEEPER IN NW MORE ROBUST (COVARIATE) SHIFT
- SLIGHT REGULARIZING EFFECT

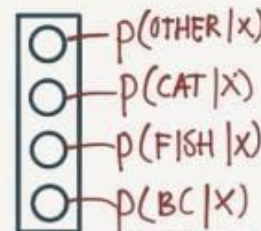### MULTICLASS CLASSIFIC.

CAT   FISH   BABY CHICK   OTHER
$C = \#\ CLASSES = 4$

SOFTMAX ACTIVATION
$t = e^{(z^{[L]})}$
$a^{[L]} = \dfrac{t}{\sum t_i}$

p(OTHER|x)
p(CAT|x)
p(FISH|x)
p(BC|x)

SUM: 1

EX: $z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$  $t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$

$\overline{176.3}$

$\Rightarrow a^{[L]} = \dfrac{t}{176.3} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.02 \\ 0.114 \end{bmatrix}$

11.4% PROB IT'S A BABY CHICK

# 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

STANDARD — COST / #ITER

MINI BATCH — COST / #ITER

### CHOOSING THE MINIBATCH SIZE

SIZE = $m$ → BATCH GRAD DESC.
SIZE = 1 → STOCHASTIC GRAD DESC

**BATCH** MIDDLE **STOCHASTIC**

BATCH: TOO LONG PER ITERATION
STOCHASTIC: LOOSE ALMOST ALL SPEED FROM VECTORIZATION
— SHORT ITER.
— USES VECT.

**TIP**
IF YOU HAVE < 2000 SAMPLES
USE SIZE = 2000
OTHERWISE, USE 64, 128, 256...
SO X+y FITS IN CPU/GPU CACHE

---

## GRADIENT DESCENT W. MOMENTUM

SLOWER LEARNING
FASTER LEARNING

WE WANT TO REDUCE OSCILLATION ↕ 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.

$$S_{dw} = \beta S_{dw} + (1-\beta) dw^2$$
$$S_{db} = \beta S_{db} + (1-\beta) db^2$$
$$w = w - \alpha \frac{dw}{\sqrt{S_{dw}}} \qquad b = b - \alpha \frac{db}{\sqrt{S_{db}}}$$

---

## ADAM OPTIMIZATION
COMBO OF GD W MOMENTUM & RMSProp

## LEARNING RATE DECAY

IDEA: USE A LARGE $\alpha$ IN THE BEGINNING. THEN DECREASE AS WE GET CLOSER TO GOAL

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

EXPONENTIAL: $\alpha = 0.95^{EPOCH} \alpha_0$

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

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

OPTION 5: DISCRETE STAIRCASE

OPTION 6: MANUAL

EPOCH = 1 PASS THROUGH THE DATA

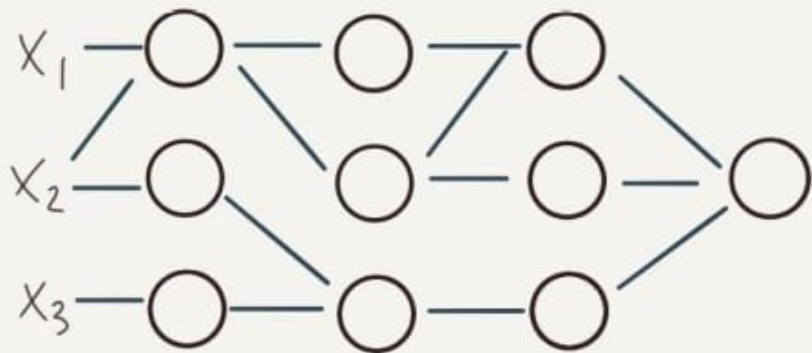# REGULARIZATION
## PREVENTING OVERFITTING

## L2 REGULARIZATION

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

← EUCLIDEAN NORM

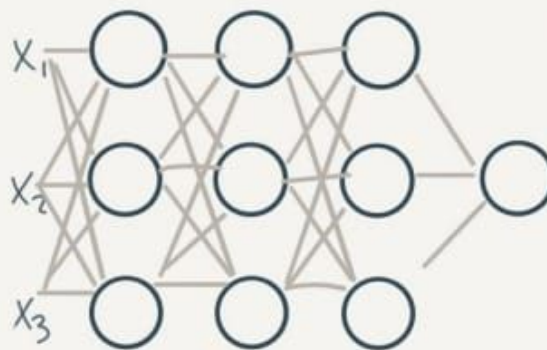## L1 REGULARIZATION

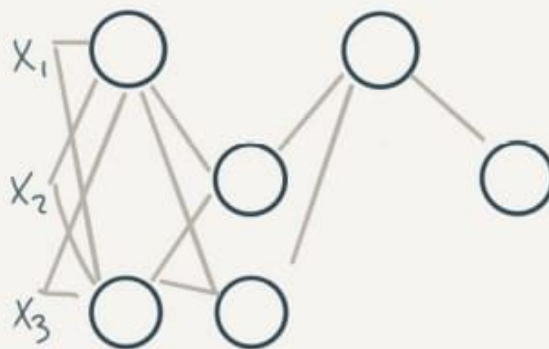$$COST: J(w,b) = \frac{1}{m}\sum_{i=1}^{m} \mathcal{L}(\hat{y},y) + \frac{\lambda}{m} \|w\|_1$$

BOTH PENALIZE LARGE WEIGHTS ⇒ SOME WILL BE CLOSE TO ∅ ⇒ SIMPLER NETWORKS

## DROPOUT

$X_1$ $X_2$ $X_3$

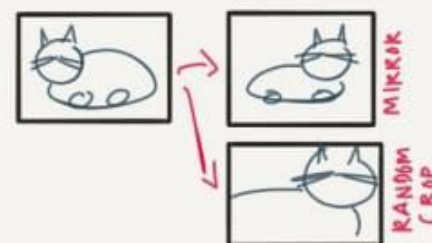FOR EACH ITERATION & SAMPLE SOME NODES ARE RANDOMLY DROPPED (BASED ON KEEP-PROB)
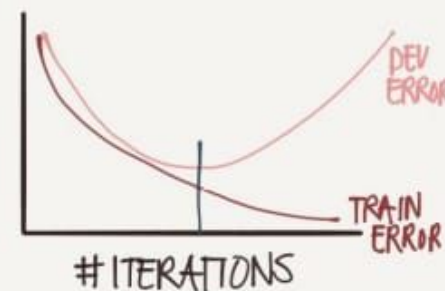
$X_1$ $X_2$ $X_3$

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

## OTHER REGULARIZATION TECHNIQUES

### DATA AUGMENTATION
GENERATE NEW PICS FROM EXISTING

MIRROR

RANDOM CROP

### EARLY STOPPING

DEV ERROR

TRAIN ERROR

#ITERATIONS

PROBLEM: AFFECTS BOTH BIAS & VARIANCE

@TessFerrandez

# SETTING UP
## YOUR ML APP

## CLASSIC ML
100 - 10000 SAMPLES

| TRAIN | DEV | TEST |
|-------|-----|------|
| 60% | 20% | 20% |

ALL FROM SAME PLACE
DISTRIBUTION

## DEEP LEARNING
1M SAMPLES

| TRAIN | | D | T |
|-------|--|---|---|
| 98% | | 1% | 1% |

EX: TRAIN — PRO CAT PICS FROM INTERNET

DEV/TEST — BLURRY CAT PICS FROM APP

**TIP** — DEV & TEST SHOULD COME FROM SAME DISTRIBUTION

IDEA → CODE → EXPERIMENT → IDEA

## BIAS/VARIANCE

| HIGH BIAS "UNDERFIT" | JUST RIGHT | HIGH VARIANCE "OVERFIT" |
|---|---|---|

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

ASSUMING HUMANS GET 0% ERROR

## THE ML RECIPE

HIGH BIAS —y→ BIGGER NETWORK, TRAIN LONGER (DIFF NN ARCHITECTURE)

↓ no

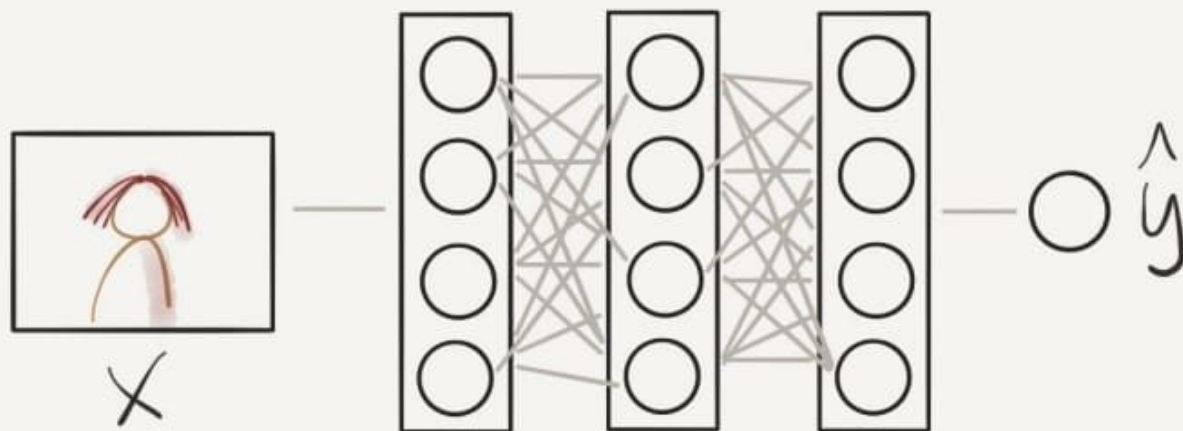HIGH VARIANCE —y→ MORE DATA (TRAIN), REGULARIZATION (DIFF NN ARCHITECTURE)

↓ no

DONE

# 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.

$x$

$\hat{y}$

LOW LEVEL AUDIO WAVE FEATURES
↗ ↘ PITCH

PHONEMES — WORDS — SENTENCES

CAT

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

@TessFerrandez