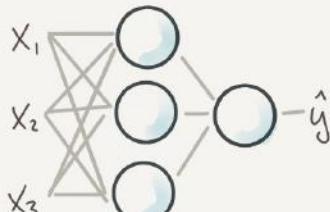


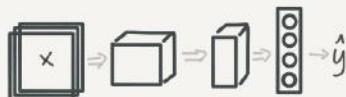
# INTRO TO DEEP LEARNING

## SUPERVISED LEARNING

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

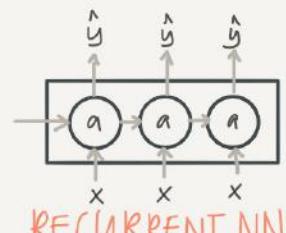


STANDARD NN

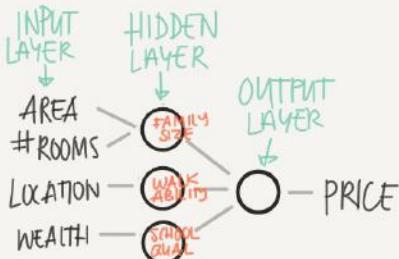


CONVOLUTIONAL NN

## NETWORK ARCHITECTURES



RECURRENT NN



NNs CAN DEAL WITH BOTH STRUCTURED & UNSTRUCTURED DATA

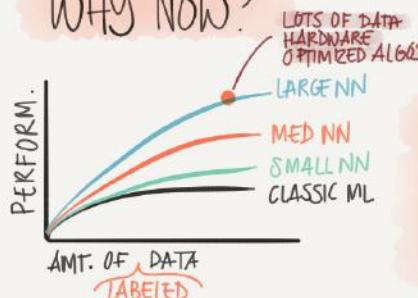


STRUCTURED



"THE QUICK BROWN FOX"  
UNSTRUCTURED  
HUMANS ARE GOOD  
AT THIS

## WHY NOW?



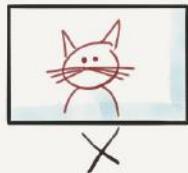
ONE OF THE BIG BREAKTHROUGHS HAS BEEN MOVING FROM SIGMOID TO RELU FOR FASTER GRADIENT DESCENT



IDEA → EXPERIM. → CODE

FASTER COMPUTATION IS IMPORTANT TO SPEED UP THE ITERATIVE PROCESS

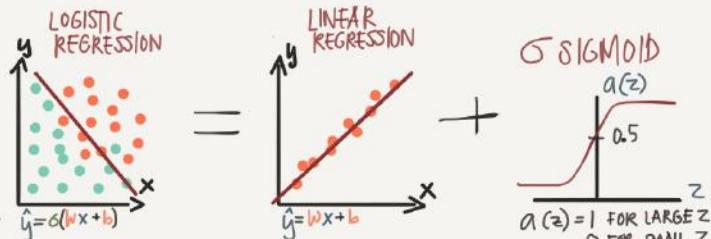
## BINARY CLASSIFICATION



1: CAT  
0: NOT CAT

$y$

Calculate current loss (forward propagation)  
Calculate current gradient (backward propagation)  
Update parameters (gradient descent)



THE TASK IS TO LEARN  $w \in b$  BUT HOW?

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

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

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

COST = LOSS FOR THE ENTIRE DATASET

The lower the better

$$w := w - \alpha \frac{\Delta J(w)}{\Delta w}$$

$$b := b - \alpha \frac{\Delta J(b)}{\Delta b}$$

} to change the COST function

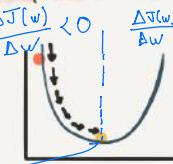
# LOGISTIC REGRESSION

## AS A NEURAL NET

Slope = Derivatives

## FINDING THE MINIMUM WITH GRADIENT DESCENT

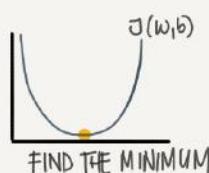
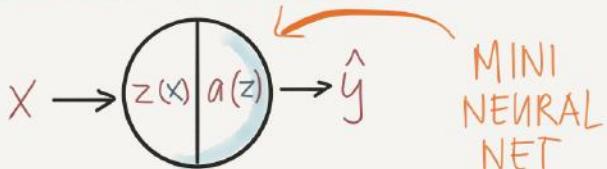
$$w := w - \alpha \frac{\Delta J(w)}{\Delta w}$$



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

REPEAT UNTIL YOU REACH BOTTOM (CONVERGE)

## PUTTING IT ALL TOGETHER



$$z(x) = wx + b$$

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

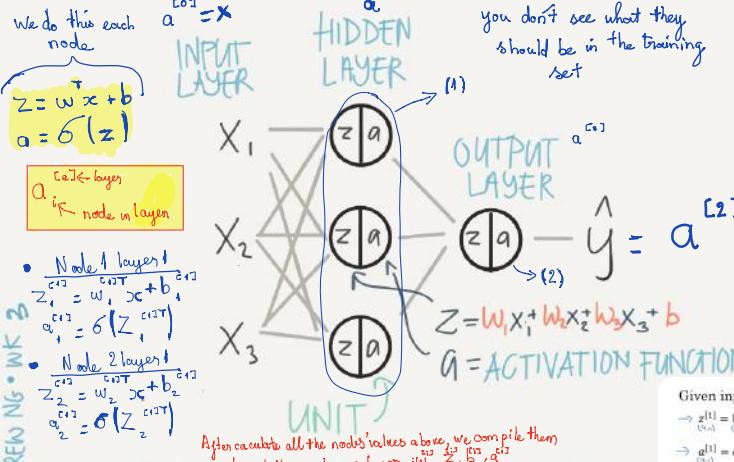
1. FORWARD PROPAGATION • CALCULATE  $\hat{y}$
2. BACKWARD PROPAGATION • GRADIENT DESCENT + UPDATE  $w \in b$

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

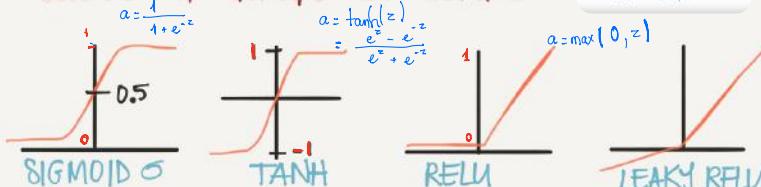
REPEAT UNTIL IT CONVERGES

# SHALLOW NEURAL NETS

## 2 LAYER NEURAL NET



## ACTIVATION FUNCTIONS



### BINARY CLASSIFIER

- ONLY USED FOR OUTPUT LAYER  
→ GRADIENT DESCENT IS FASTER

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

Better than Sigmoid for sure!

### DEFAULT CHOICE FOR ACTIVATION

SLOPE = 1/0

AVOIDS UNDEF SLOPE AT 0 BUT RARELY USED IN PRACTICE

Your output	Output layer
Binary Classification (0, 1)	Sigmoid
All other types	ReLU
⇒ Hidden layer: ReLU	

## WHY ACTIVATION FUNCTIONS?

EX. WITH NO ACTIVATION -  $a = z$

$$a^{[1]} = z^{[1]} = w^{[1]} x + b^{[1]}$$

$$a^{[2]} = z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

LAYER 1  
LAYER 2

PLUG IN  $a^{[1]}$

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

LINEAR FUNCTION

Logistic Regression looks like a sigmoid with 0 but should not for neural network.

## 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

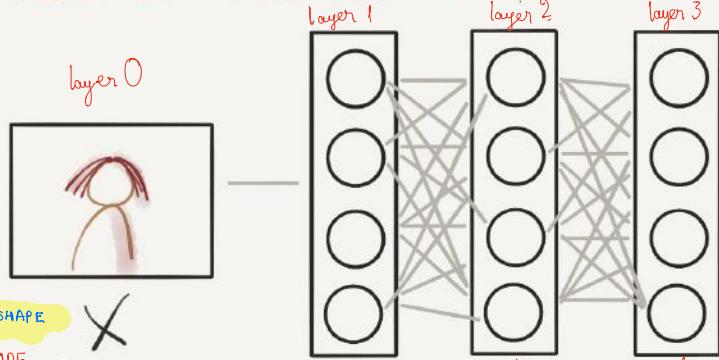
Activation Function is the foundation of neural network

© Tess Ferrandez

# DEEP NEURAL NETS

- $l = \# \text{ of layers}$
- $l = 3$
- $n^{(l)} = \# \text{ units in layer } l$
- $n^{(1)} = n^{(2)} = n^{(3)} = n = 4$
- $n^{(1)} = n_x = 1$
- $a^{(l)} = \text{activation in layer } l$
- $a^{(l)} = g(z^{(l)})$
- $w^{(l)}$ : weights for  $z^{(l)}$

## WHY DEEP NEURAL NETS?



VECTOR SHAPE

1 NODE	
$z^{(l)}$	$(n^{(l)}, 1)$
$w^{(l)}$	$(n^{(l)}, n^{(l-1)})$
$a^{(l-1)}$	$(n^{(l-1)}, 1)$
$b^{(l)}$	$(n^{(l)}, 1)$



1 LAYER	
$z^{(l)}$	$(n^{(l)}, m)$
$W^{(l)}$	$(n^{(l)}, n^{(l-1)})$
$A^{(l-1)}$	$(n^{(l-1)}, m)$
$b^{(l)}$	$(n^{(l)}, 1) \rightarrow (n^{(l)}, m)$



LOW LEVEL  
AUDIO WAVE  
FEATURES  
 $\rightarrow$  PITCH

— PHONEMES — WORDS — SENTENCES

CAT

Layer neural net work

[LINEAR  $\rightarrow$  tanh] ( $L-1$  times)  $\rightarrow$  LINEAR  $\rightarrow$  SIGMOID

© Tess Ferrandez

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

Layer  $l$ :  $w^{(l)}$ ,  $b^{(l)}$ , Output  $a^{(l)}$ , Cache  $z^{(l)}$   
 Forward: Input  $a^{(l-1)}$ , Output  $a^{(l)}$ , Cache  $z^{(l)}$   

$$z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)} \quad | \quad a^{(l)} = g(z^{(l)})$$
  
 Backward: Input  $\Delta a^{(l)}$ , output  $\Delta a^{(l-1)}$ , Cache  $\{z^{(l)}\}$ ,  $\Delta w^{(l)}$ ,  $\Delta b^{(l)}$

VERY DATA HUNGRY

NEED LOTS OF COMPUTER POWER

ALWAYS VECTORIZE  
VECTOR MULT. CHEAPER THAN FOR LOOPS  
COMPUTE ON GPUs

Influence the parameters  
 $(w, b)$

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

2

# SETTING UP YOUR ML APP

## CLASSIC ML

100 - 10000 SAMPLES

TRAIN	DEV	TEST
60%	20%	20%

ALL FROM SAME PLACE  
DISTRIBUTION

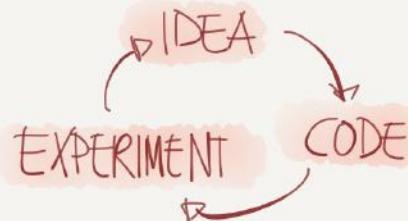
## DEEP LEARNING

1M SAMPLES

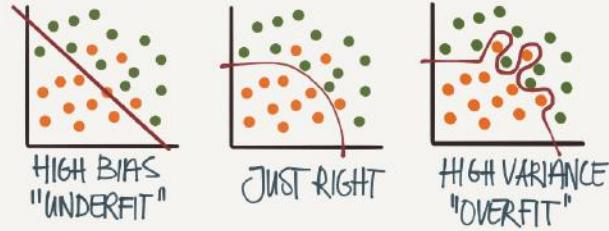
TRAIN	DEV	TEST
98%	1%	1%



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



## BIAΣ/VARIANCE



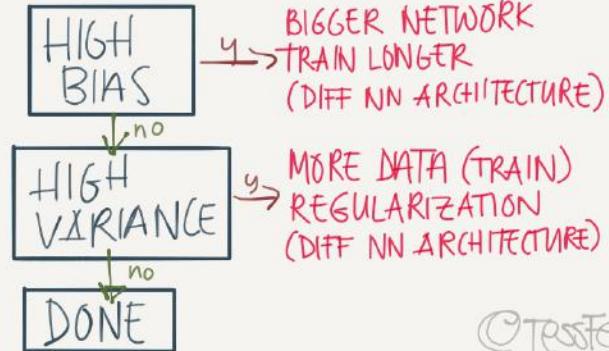
	ERROR			
	TRAIN	15%	15%	0.5%
TRAIN	1%	15%	15%	0.5%
DEV / TEST	11%	16%	30%	1%
DEV / TEST	11%	16%	30%	1%
	HIGH <sup>Overfit</sup> BIAS	HIGH <sup>Underfit</sup> BIAS	HIGH BIAS & VARIANCE	LOW BIAS & VARIANCE
	HIGH VARIANCE			

This is a really bad scenario that you should avoid.

you don't really evaluate well the model so the difference in error % is too large.

ASSUMING HUMANS GET 0% ERROR

## THE ML RECIPE



Cases of no bias variance trade off can happen in deep learning as we can have both happen at the same time

# REGULARIZATION

## PREVENTING OVERTFITTING

"Weight decay"

### L2 REGULARIZATION

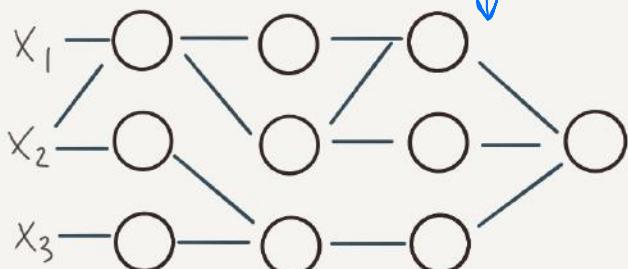
$$\text{COST: } J(w, b) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}_i, y_i) + \frac{\lambda}{2m} \|w\|_2^2$$

$\lambda$  - regularization parameter  
 $\|w\|_2^2 = \sum_{j=1}^{n_p} w_j^2 = w^T w$   
 regularization function  
 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



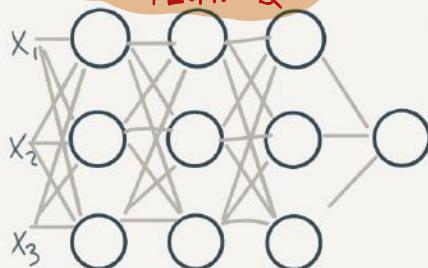
$\lambda$  - regularization parameter

$$\|w\|_2^2 = \sum_{j=1}^{n_p} w_j^2 = w^T w$$

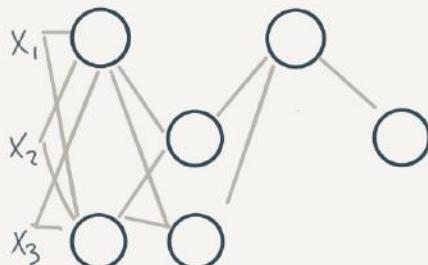
regularization function

EUCLIDEAN NORM

## DROPOUT TECHNIQUE



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

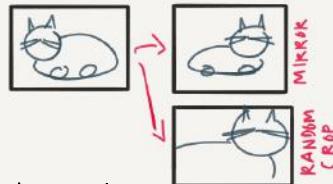


WE GET SIMPLER NETS  
 & LESS CHANCE TO RELY ON  
 SINGLE FEATURES

## OTHER REGULARIZATION TECHNIQUES

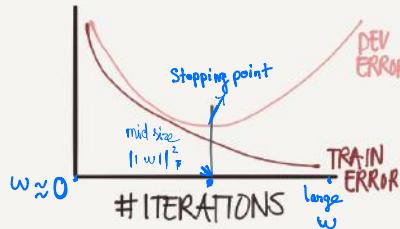
### DATA AUGMENTATION

GENERATE NEW PICS FROM EXISTING



4  $\rightarrow$  4  $\rightarrow$  distortion

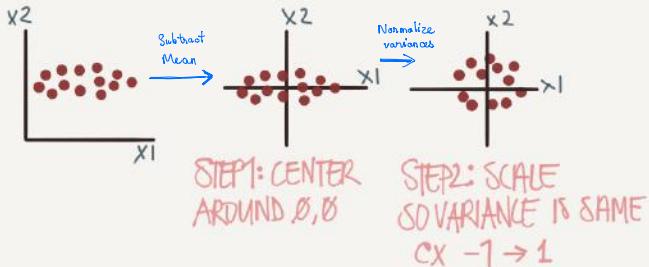
### EARLY STOPPING



PROBLEM: AFFECTS BOTH  
 BIAS & VARIANCE

# OPTIMIZING TRAINING

## NORMALIZING INPUTS



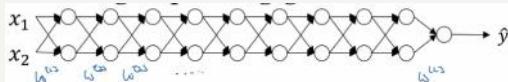
**TIP**  
USE SAME AVG/VAR TO NORMALIZE DEV/TEST

## WHY DO WE DO THIS?



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

we don't need that much time to reach to the center like the left



This makes training really difficult because the more layers of your deep NN, the bigger or smaller your activation

$$\begin{aligned} a^{(L)} &= g(z^{(L)}) \\ &= g(w^{(L)}x^{(L-1)} + b^{(L)}) \end{aligned}$$

## DEALING WITH VANISHING/EXPLODING GRADIENTS

Ex: DEEP NW (L LAYERS)

$$\hat{y} = w^{(L-1)}w^{(L-2)} \dots w^{(1)}x + b$$

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

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

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

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

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

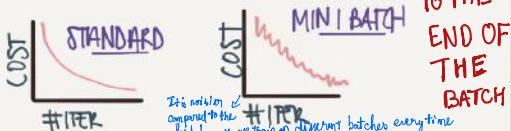
# OPTIMIZATION

## 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 WITHOUT WAITING



### CHOOSING THE MINIBATCH SIZE

$\text{SIZE} = m \rightarrow \text{BATCH GRAD DESC. } (x^{(1)}, y^{(1)}) = (x, y)$   
 $\text{SIZE} = 1 \rightarrow \text{STOCHASTIC GRAD DESC. }$  every example is its own mini batch



**TIP:** IF YOU HAVE < 2000 SAMPLES USE SIZE=2000

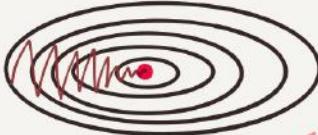
OTHERWISE, USE 64, 128, 256... SO X+Y FITS IN CPU/GPU CACHE

$$\mathbf{x} = [x^{(1)}, x^{(2)}, \dots, x^{(m)}] \quad x^{(1)}, x^{(2)}, \dots, x^{(m)}$$

$m = 5 \text{ millions}$   
each minibatch has 1,000 each

### GRADIENT DESCENT W. MOMENTUM

→ Speed up gradient descent



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.

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

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

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

### 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

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

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

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

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

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

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

EPOCH = 1 PASS THROUGH THE DATA



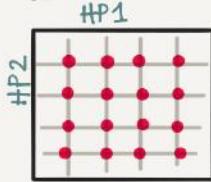
# HYPERTPARAM TUNING

WHICH HYPERPARAMS ARE MOST IMPORTANT?

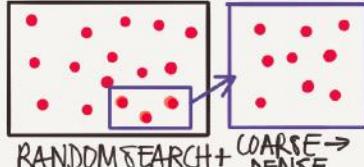
- $\alpha$  LEARNING RATE
- # HIDDEN UNITS
- MINIBATCH SIZE
- $\beta$  MOMENTUM, TURN = 0.9 }  
# LAYERS  
LEARNING RATE DECAY }
- $\beta_1 = 0.9 \quad \beta_2 = 0.999 \quad \epsilon = 10^{-8}$  (ADAM) } very rare

## TESTING VALUES

### CLASSIC ML



GRID SEARCH  
SOLUTION



PROBLEM: ONE ITERATION TAKES A LONG TIME & IN 16 GO'S WE HAVE ONLY TRIED 4 $\alpha$  - BUT 4 DIFF  $\epsilon$

NOT AS IMPORTANT

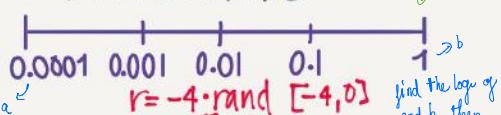
MY PANDA IS ACTUALLY A MISCLASSIFIED CAT BECAUSE I CAN'T DRAW PANDAS  
BABYSIT ONE MODEL & TUNE  
PANDA VS CAVIAR  
SPAWN LOTS OF MODELS W DIFF HP  
GOOD IF YOU HAVE LOTS OF SHARP COMP POWER

USE AN APPROPRIATE SCALE

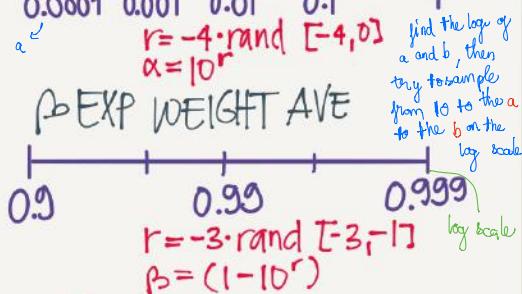
### # HIDDEN UNITS



### $\alpha$ LEARNING RATE



### $\beta$ EXP WEIGHT AVE



TIP  
RE-EVALUATE YOUR HYP. PARAMS EVERY FEW MONTHS

## MISC. EXTRAS

### BATCH NORMALIZATION

NORMALIZE LAYER OUTPUT

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

### MULTICLASS CLASSIFIC.

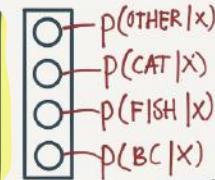


C = # CLASSES = 4 (0, ..., 3)

### SOFTMAX ACTIVATION

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

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



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

$$\Rightarrow a^{[i]} = \frac{t}{\sum t} = \frac{\begin{bmatrix} 0.842 \\ 0.042 \\ 0.02 \\ 0.14 \end{bmatrix}}{176.3} = 0.842$$

$\Sigma t = 176.3$   
11.4% PROB IT'S A BABY CHICK

© TessFernandez

3

# STRUCTURING YOUR ML PROJECTS

## SETTING YOUR GOAL

\* GOAL SHOULD BE A SINGLE #

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

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

F1 = HARMONIC MEAN BETW.  
RECALL & PRECISION

This is the best one

\* DEFINE OPTIMIZING VS SATISFYING METRICS

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

MAXIMIZE ACC.  
GIVEN TIME  $\leq$  100ms

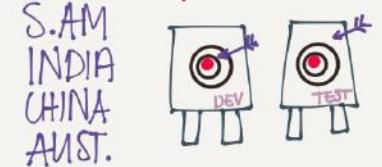
ACCURACY = OPTIMIZING  
RUNTIME = SATISFYING

If satisfying below a threshold then we choose the optimizing

## SELECTING YOUR DEV/TEST SETS

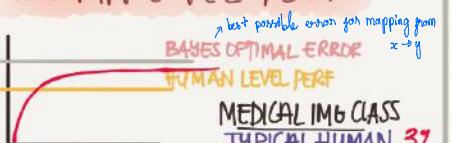
### DATA

OPTION 1:  
 US → DEV = UK, US, EUR  
 UK → TEST = REST  
 EUROPE → S.AM, INDIA, CHINA, AUST.



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?

MEDICAL IMAGING CLASS  
 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

Available bias: the bias is as closed to human as possible so there is no need to actually do more about it \*

## CAT CLASSIFICATION

Priority for Bayes error (0.5%)	Emon	
	A	B
HUMAN	1%	7.5%
TRAIN ERR	8%	8%
DEV ERR	10%	10%

FOCUS ON BIAS      FOCUS ON VARIANCE

### HUMAN

TRAIN BIGGER NETS.

TRAIN LONGER/BETTER OPT. (RMSPROP ALGO'S)

CHANGE NN ARCH OR HYPERPARAMS

### TRAIN

MORE DATA (TRAIN)

REGULARIZATION

NN ARCHITECTURE

### DEV

	A	B
HUMAN	0.5	0.5

	A	B
TRAIN ERR	0.6	0.3

	A	B
DEV ERR	0.8	0.4

	A	B
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

# ERROR ANALYSIS

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

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

	DOG	BLURRY	INSTA FILTER	BIG EAT	...
1	1		1		
2				1	
3		1			
...					
100					
5	43%	61%	12%	1	5% OF ALL ERRORS

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

Ways to determine if needed to go through mislabel

- Overall Dev set error
- Errors due incorrect labels
- Errors due to other causes

Correcting incorrect dev/test set examples

- Apply same process to your dev + test
- Consider examining examples algorithm got right/wrong
- Train and dev/test data may now come from slightly different distributions

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 LABELED AS MICE)

ADD EXTRA COL. IN ERROR ANALYSIS AND USE SAME CRITERIA

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

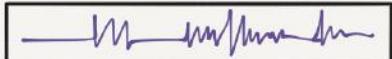
Ideas to improve cat-detection:

- Fix pictures of dogs being recognized as cats
- Fix great cats (lion, panthers,...) being misrecognized
- Improve performances on blurry pics

Based on this we can decide which is the worthiest of our time to do

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

# TRAIN vs DEV/TEST MISMATCH

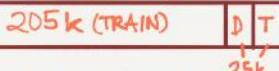
## AVAILABLE DATA

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



PROBLEM: DEV/TEST IS NOW  
MOSTLY WEB IMGS (NOT REPR. OF END SCENARIO)

SOLUTION: LET DEV/TEST COME  
FROM APP. THEN SHUFFLE 5K  
OF APP PICS IN WEB FOR TRAIN



## BIAS & VARIANCE IN MISMATCHED TRAIN/DEV

HUMANS	~0%
TRAIN	1%
DEV ERR	10%

IS THIS DIFF  
DUETO THE MODEL  
NOT GENERALIZING  
OR IS DEV DATA  
MUCH HARDER

A: CREATE A TRAIN-DEV SET  
THAT WE DON'T TRAIN ON

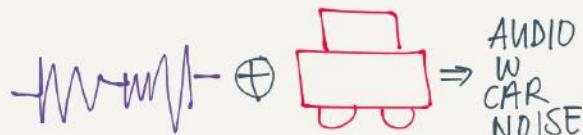
TRAIN		FD	D	T
A	B	C	D	
Human error: 0%				This is too large.
TRAIN	1%	1%	10%	10%
TRAIN-DEV	9%	15%	11%	11%
DEV	10%	10%	12%	20%
VARIANCE	TRAINY	BIAS	BIAS + DATA	MISMATCH
Training:	variance	evaluate on the training distribution		
Train - dev:	data mismatch			
Dev:	degree of overfitting	evaluate on the dev-test set distribution		
Test:	to tie dev set			

Human level: 4%  
↓ avoidable bias  
Training: 7%  
↑ variance  
Train - dev: 10%  
Dev: 12%  
Test: 12%  
↓ evaluate on the  
training distribution  
↓ evaluate on the dev-test  
set distribution  
↓ to tie dev set

## 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 (ex: artificial data synthesis,



## 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.

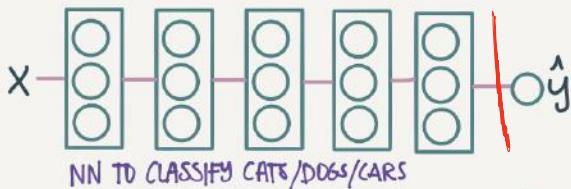
# EXTENDED LEARNING

## TRANSFER LEARNING

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

✓ pre-train the weights and parameters first on a different dataset before applying on the one we need

↙ pretraining / fine tuning



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

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

**OPTION 2:** YOU HAVE LOTS OF RADIOLOGY IMB.

SOLUTION: INIT WITH WEIGHTS FROM CAT NN  
RETRAIN ALL LAYERS

THIS IS MICROSOFT CUSTOM VISION

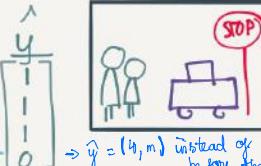
It is useful for the system to learn the initial weights so that they can perform better on harder datasets regardless of new training size.

\* Never Train in reverse | initial size < target size, Initial low level feature could be helpful for target learning

## MULTITASK LEARNING

### TRAINING ON MULT. TASKS AT ONCE

DETECT CAR STOP SIGN PEDESTRIAN TRAFFIC LIGHT



$\hat{y} = (1, m)$  instead of  $(1, m)$  before that

UNLIKE SOFTMAX. MANY THINGS CAN BE TRUE

usual logistic loss

$$\text{COST: } J(w, b) = \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^4 \ell(y_i, \hat{y}_j)$$

SUMMING OVER ALL OUTP OPTIONS

WE COULD HAVE JUST TRAINED 4 NN'S INSTEAD BUT.. MT LEARNING MAKES SENSE WHEN

A. THE LEARNING DATA YOU HAVE FOR THE DIFF TASKS IS QUITE SIMILAR - E.g. THE AMOUNTS (E.G. 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

C. Usually: amount of data you have for each task is quite similar

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 SPLIT:

1. LOCATE BONES TO FIND LENGTHS USING ML
2. TRAIN MODEL TO PREDICT AGE BASED ON BONELLENGTH

## END-TO-END

RADIOLOGY  $\longrightarrow$  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 DATA ( $X \rightarrow Y$ )
- EXCLUDES POTENTIALLY USEFUL HAND-MADE COMPONENTS

LABELLED  
© Tessierandez

# 4 CONVOLUTION

## FUNDAMENTALS

### COMPUTER VISION

IMAGE  
CLASSIFICATION



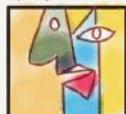
CAT OR  
NOT CAT

OBJECT  
DETECTION



WHERE IS  
THE CAR?

NEURAL  
STYLE  
TRANSFER



PAINT ME  
LIKE PICASSO

PROBLEM: IMAGES CAN BE BIG  
 $1000 \times 1000 \times 3$  (RGB) = 3M

WITH 1000 HIDDEN UNITS WE  
NEED  $3M \times 1000 = 3B$  PARAMS

SOLUTION: USE CONVOLUTIONS  
IT'S LIKE SCANNING OVER YOUR  
IMG WITH A MAGNIFYING GLASS  
OR FILTER

ALSO SOLVES THE PROBLEM  
THAT THE CAT IS NOT  
AWAYS IN THE SAME  
LOCATION IN THE IMG

### CONVOLUTION

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

INPUT 6x6 IMAGE

10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0

INPUT 6x6 IMAGE



THIS IS LIKE ADDING  
AN 'INSTA' FILTER THAT  
JUST SHOWS OUTLINES

WE COULD HARD-CODE FILTERS. JUST LIKE WE  
CAN HARD-CODE HEURISTIC RULES... BUT.... A MUCH BETTER  
WAY IS TO TREAT THE FILTER # AS PARAMS  
TO BE LEARNED

$$(3 \times 1) \quad (1 \times 1) \quad (2 \times 1) \quad (0 \times 0) \quad (5 \times 0)$$

(3x1)

R

\*

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

# CONVENTIONAL NEURAL NETS • COURSERA

$$(6 \times 6)_{n \times n} \xrightarrow{f \times f} (3 \times 3) = (4 \times 4)_{[(n-f+1) \times (n-f+1)]}$$

## PADDING

PROBLEM: IMAGES SHRINK

$$6 \times 6 \rightarrow 3 \times 3 \rightarrow 4 \times 4$$

PROBLEM: EDGES GET LESS 'LOVE'

SOLUTION: PAD W. A BORDER OF 0's BEFORE CONVOLVING

0	0	0	6	0	0	0	0
0	3	0	1	2	7	4	0
0	1	5	8	9	3	1	0
0	2	7	2	5	1	3	0
0	0	1	3	1	7	8	0
0	4	8	1	6	2	8	0
0	2	4	5	2	3	9	0
0	0	0	0	0	0	0	0

TWO COMMONLY USED  
PADDING OPTIONS  
(HOW MUCH TO PAD)

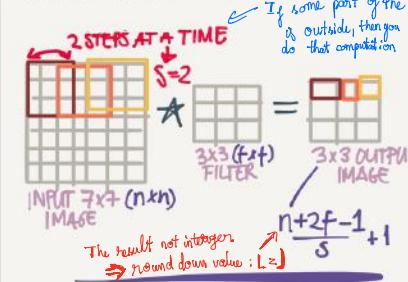
'VALID'  $\Rightarrow P=0$  NO PADDING  
'SAME'  $\Rightarrow P=\frac{f-1}{2}$  OUTPUT SIZE = INPUT SIZE  
FILTER SIZE

$$(n-f+1) \times (n-f+1)  
(n+2P-f+1) \times (n+2P-f+1)$$

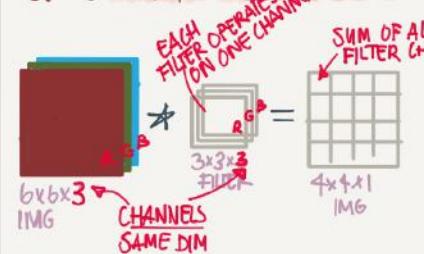
NOTE: ALL CONVOLUTION IDEAS CAN BE APPLIED TO 1D AS WELL LIKE EKG SIGNALS. AND 3D LIKE CT-SCANS

## STRIDE

WHAT PACE YOU SCAN WITH

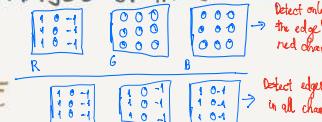


## CONVOLUTIONS OVER VOLUMES



THIS ALLOWS US TO DETECT FEATURES IN COLOR IMAGES FOR EXAMPLE

MAYBE WE WANT TO FIND ALL EDGES OR MAYBE ORANGE BLOBS

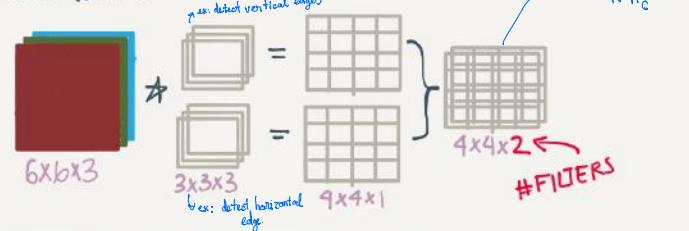


$$\text{Parameters: } ((m_i \times n_j \times d_i) + 1) \times k$$

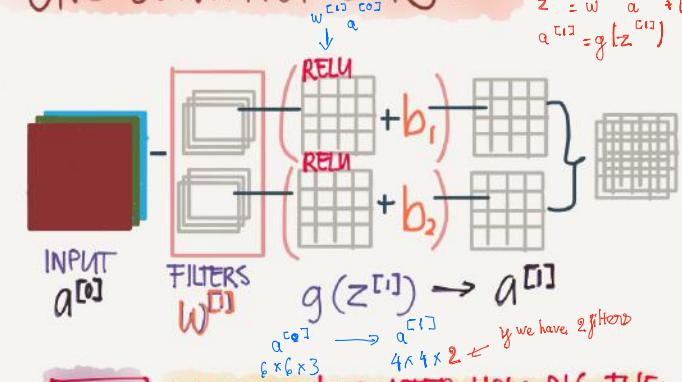
Most deep learning engineers don't need to bother with the details of the backward pass

## MULTIPLE FILTERS

DETECTING MULTIPLE FEATURES AT A TIME



## ONE CONV. NET LAYER



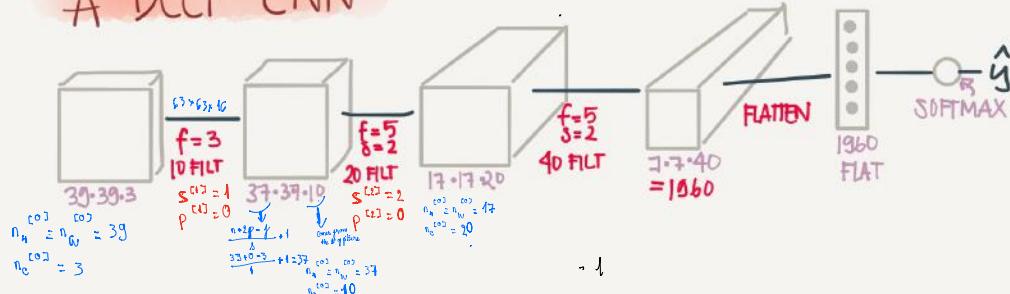
**NOTE** IT DOESN'T MATTER HOW BIG THE INPUT IS - THE LEARNABLE PARAMS  $W$  &  $b$  ONLY DEPEND ON THE # OF FILTERS AND THEIR SIZES.

$$W = 3 \cdot 3 \cdot 3 \cdot 2 = 54 \quad \left\{ \begin{array}{l} \text{56 PARAMS} \\ \text{TO LEARN} \end{array} \right.$$

$$b = 2$$

© TessFernandez

# A DEEP CNN



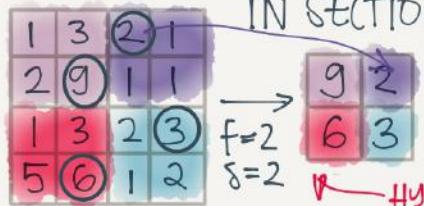
A LOT OF THE WORK IS FIGURING OUT HYPERPARAMS  
 $= \# \text{FILTERS}, \text{STRIDE}, \text{PADDING} \text{ ETC}$

TYPICALLY SIZE  $\rightarrow$  TREND DOWN  
 $\# \text{FILTERS} \rightarrow$  TREND UP

## TYPICAL CONV.NET LAYERS

- CONVOLUTION
- POOLING
- FULLY CONNECTED

### POOLING (MAX)

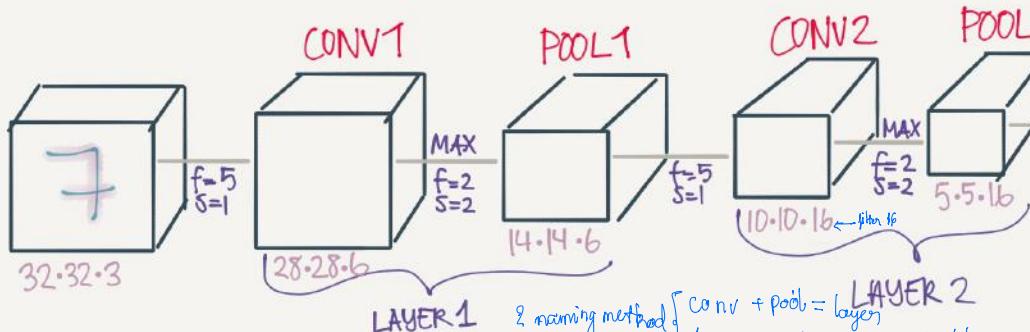


FIND MAX VAL IN SECTION

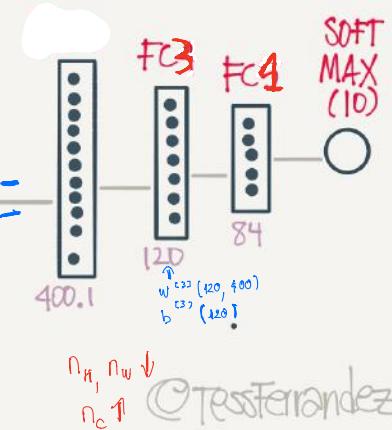
- \* REDUCES SIZE OF REPRES.
- \* SPEEDS UP COMPUTATION
- \* MAKES SOME OF THE DETECTED FEAT. MORE ROBUST

## CONV NET EXAMPLE BASED ON LeNet-5

### DETECTING HANDWRITTEN DIGITS



2 training method:  
 $\left\{ \begin{array}{l} \text{conv} + \text{pool} = \text{layer} \\ 1 \text{ conv} + 1 \text{ layer, } 1 \text{ pool} = 1 \text{ layer} \end{array} \right.$

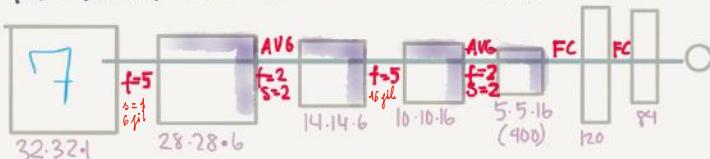


© TessFernandez

# CLASSIC CONV. NETS

**LeNet-5**: to recognize hand written digits

DOCUMENT CLASSIFICATION



TRENDS: HEIGHT/WIDTH GO DOWN  $n_h, n_w \downarrow$   
CHANNELS GO UP  $n_c \uparrow$

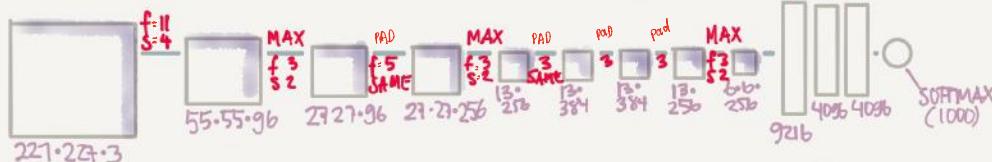
COMMON PATTERN: A COUPLE OF CONV(1<sup>st</sup>)/POOL LAYERS FOLLOWED BY A FEW FC

OLD STUFF: USED AVG POOLING INST. OF MAX  
PADDING WAS NOT VERY COMMON  
IT USED SIGMOID/TANH INST. OF RELU

# AlexNet

IMAGE CLASSIFICATION

$\approx 60M$  PARAMETERS



- SIMILAR TO LeNet BUT MUCH BIGGER
- USES RELU
- THE NN THAT GOT RESEARCHERS INTERESTED IN VISION AGAIN

# VGG-16

ALL CONV. LAYERS HAVE SAME PARAMS  
 $f=3 \times 3, s=1, p=same$   
AND POOLING LAYER  $2 \times 2, s=2$



- VERY DEEP
- EASY ARCHITECTURE
- # FILTERS DOUBLE 64, 128, 256, 512

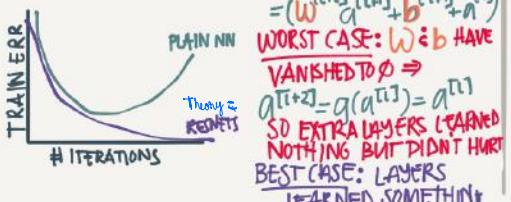
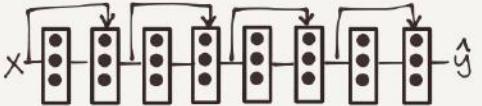
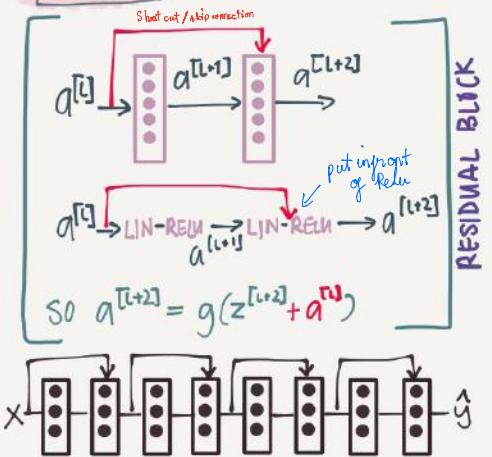
# CONVOLUTIONAL NEURAL NETS • COURSERA

## SPECIAL NETWORKS

**ResNets**: an experiment when training error is too high (add noise to the end of neural network)

**PROBLEM: DEEP NN OFTEN SUFFER PROBLEMS W VANISHING OR EXPLODING GRADIENTS**

### SOLUTION) RESIDUAL NETS



## NETWORK IN NETWORK (1x1 CONVOLUTION)

6	5	3	2
4	1	0	5
5	8	2	4
0	3	6	1

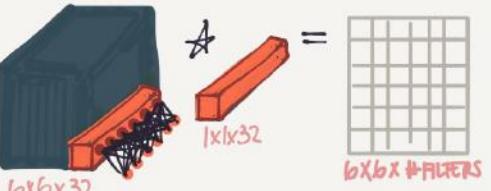
 $\star 2 =$ 

12	10	6	4
8	2	18	10
10	16	4	8
0	6	12	2

### 1x1 CONVOLUTION

IT SEEMS PRETTY USELESS, BUT IT ACTUALLY SERVES 2 PURPOSES

### 1. NETWORK IN A NETWORK



LEARN COMPLEX, NON-LINEAR RELATIONSHIPS ABOUT A SLICE OF A VOLUME

### 2. REDUCING # CHANNELS

28x28x192
-----------

 $\star$ 

1x1x92
--------

 $=$ 

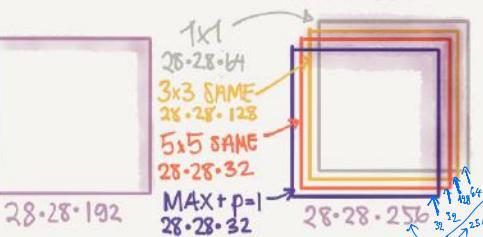
28x28x32
----------

32 FILT

## INCEPTION NETWORKS

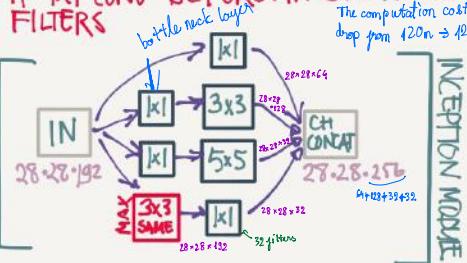
Called "GooleNet" by the Google employees to pay homage to the LeNet we talked above

INSTEAD OF CHOOSING A 1x1, 3x3, 5x5 OR A POOLING LAYER - CHOOSE ALL



PROBLEM: VERY EXPENSIVE TO COMPUTE

SOLUTION: SHRINK THE # CHANNELS W/ A 1x1 CONV BEFORE APPLYING ALL THE FILTERS



TO BUILD AN INCEPTION NETWORK YOU MAINLY STACK A BUNCH OF INCEPTION MODULES



INCEPTION THE MOVIE

© TessFerrandez

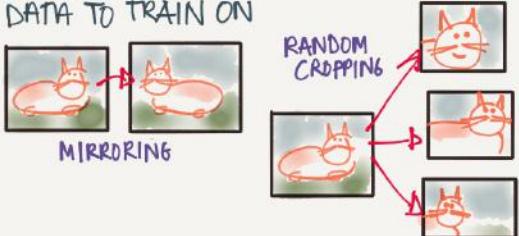
# PRACTICAL ADVICE

## USE OPEN SOURCE IMPLEMENTATIONS

SOME OF THE PAPERS ARE HARD TO IMPLEMENT FROM SCRATCH - USING OS YOU CAN REUSE OTHER PPLS WORK  
DON'T FORGET TO CONTRIBUTE

## DATA AUGMENTATION

WE ALMOST ALWAYS NEED MORE DATA TO TRAIN ON

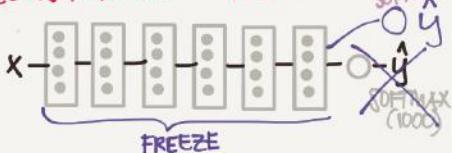


## TRANSFER LEARNING



WANT TO TRAIN A CLASSIFIER FOR YOUR CATS BUT DON'T HAVE ENOUGH PICTURES

**SOLUTION** DOWNLOAD SOMEONE ELSE'S PRETRAINED NET & WEIGHTS



FREEZE THE PARAMS, AND JUST REPLACE THE SOFTMAX LAYER WITH YOUR OWN & TRAIN

IF YOU HAVE MORE PICS • RETRAIN A FEW OF THE LATER LAYERS (MAYBE INITIALIZING WITH THE PRETRAINED WEIGHTS)

## STATE OF COMPUTER VISION

WE HAVE LOTS OF DATA  
SIMPLER ALGORITHMS  
LESS HAND ENGINEERING

- SPEECH RECD.

- IMAGE RECOGNITION

- OBJECT DETECTION  
IMGS IN LABELED BOXES

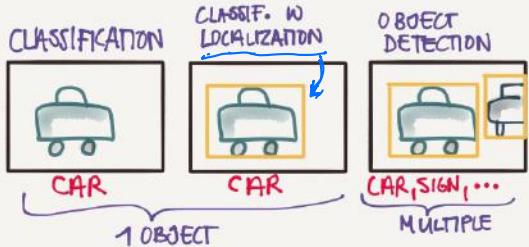
WE HAVE LITTLE LABELED DATA  
MORE HAND ENGINEERING

TIPS FOR DOING WELL ON BENCHMARKS/COMPETITIONS

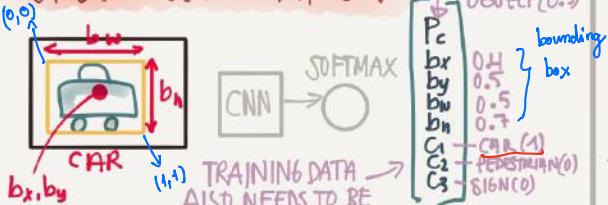
- \* ENSEMBLING. AVG OUTPUTS FROM MULT NN
- \* MULTI-CROP AT TEST TIME AVG OUTPUTS FROM MULTIPLE CROPS OF THE IMAGE

IN PRACTICE THEY ARE NOT USED IN PRODUCTION BECAUSE THEY ARE COMPUTE & MEM EXPENSIVE

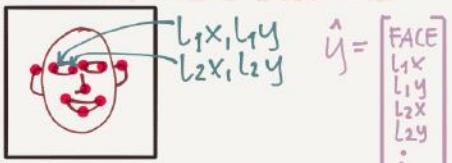
# DETECTION ALGORITHMS



## OBJECT LOCALIZATION



## LANDMARK DETECTION



TO DETECT LANDMARKS IN THE FACE (CORNER OF MOUTH ETC) LABEL THE X, Y COORDS OF THE LANDMARK

USED FOR SENTIMENT ANALYSIS & FOR EFFECTS LIKE PLACING CROWN ON HEAD ETC.

## SLIDING WINDOWS DETECTION



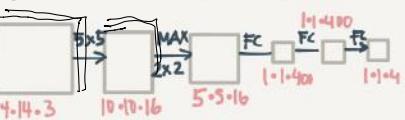
1. CREATE TINYLY CROPPED IMG OF CARS (LOTS)

2. SLIDE A WINDOW OVER THE IMG. & CLASSIFY THIS WINDOW CAR(VO) AGAINST YOUR OTHER CARS

3. REPEAT WITH SLIGHTLY LARGER WINDOW SIZE

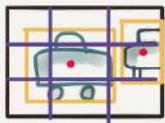
PROBLEM: VERY EXPENSIVE (TO COMPUTE)

SINCE ALL WINDOWS SHARE A LOT OF THE COMPUTATIONS WE CAN DO THIS MUCH CHEAPER W CONVOLUTIONS



NOW WE JUST PASS THROUGH ONCE AND CALC ALL AT THE SAME TIME  
EACH OF THE 4 VALS ARE RESULTS FOR EACH OF THE 4 WINDOWS

## YOLO - You Only Look Once



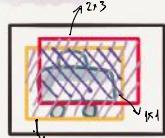
1. SPLIT IMG INTO X(g) GRID CELLS
2. FOR EACH CELL, SAY IF IT CONTAINS CAR + BOUNDING BOX (IF CELL CONTAINS THE MID POINT)

IN PRACTICE WE MIGHT HAVE A 10x10 GRID

$$X \times g = 3 \times 3 \times 8$$

HOW DO YOU KNOW HOW GOOD IT IS?

HOW GOOD IS THE RED SQUARE?



$$IOU = \frac{2 \times 2}{2 \times 2 + 2 \times 3} = \frac{4}{8} = \frac{1}{2}$$

## INTERSECTION OVER UNION

GENERALLY IF IOU  $\geq 0.5$  IT IS REGARDED AS CORRECT  $\rightarrow$  if lower than 0.5, we keep both box

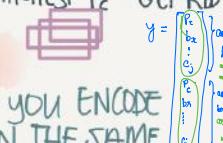
WHAT IF MULTIPLE SQUARES CLAIM THE SAME CAR?

## NON-MAX SUPPRESSION

IF TWO BOUNDING BOXES HAVE A HIGH IOU - PICK THE ONE W HIGHEST P<sub>c</sub> - GET RID OF THE REST.

## ANCHOR BOXES

ANCHOR BOXES LET YOU ENCODE MULTIPLE OBJECTS IN THE SAME SQUARE



# FACE RECOGNITION

FACE  
VERIFICATION



IS THIS PETE?  
99% ACC  
⇒ PRETTY GOOD

FACE  
RECOGNITION



WHO IS THIS?  
(OUT OF K PERSONS)  
IF K = 100 NEED  
MUCH HIGHER THAN  
99% (99.9%, something)

## ONE SHOT LEARNING

NEED TO BE ABLE TO RECOGNIZE  
A PERSON EVEN THOUGH YOU ONLY  
HAVE ONE SAMPLE IN YOUR DB.  
YOU CAN'T TRAIN A CNN WITH  
A SOFTMAX (EACH PERSON) BECAUSE

- (A) YOU DON'T HAVE ENOUGH SAMPLES
- (B) IF A NEW PERSON JOINS YOU  
NEED TO RETRAIN THE NETWORK

**[SOLUTION]** LEARN A SIMILARITY  
FUNCTION

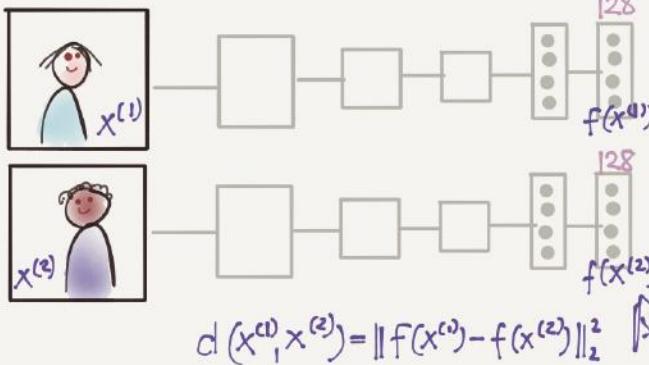
$$d(\text{img}1, \text{img}2) = \text{degree of difference}$$

BUT HOW DO YOU LEARN THIS?

- Input: image, name/ID (1 : 1)
  - Output: whether the input image is that of the claimed person.
  - "Is this the claimed person?"
- Recognition
  - Has a database of K persons
  - Get an input image
  - Output ID if the image is any of the K persons (or not recognized)
  - "Who is this person?"

## SIAMESE NETWORK

## DeepFace



$$d(x^{(1)}, x^{(2)}) = \|f(x^{(1)}) - f(x^{(2)})\|_2^2$$

## TRIPLET LOSS

FaceNet



$$\text{WANT } \|f(A) - f(P)\|^2 \leq \|f(A) - f(N)\|^2 \Rightarrow d(A, P) - d(A, N) \leq 0$$

BUT WE WANT A GOOD MARGIN, SO...

$$d(A, P) - d(A, N) + \alpha \leq 0$$

↑ push the d(A,P) and d(A,N) further away from each other

HOW DO WE CHOOSE TRIPLETS (A,P,N)  
TO TRAIN ON?

- IF A/P ARE VERY SIMILAR, & A/N ARE VERY DIFFERENT  
TRAINING IS VERY EASY.

SELECT A/N THAT ARE PRETTY SIMILAR TO TRAIN A GOOD NET

Training set: 10k pictures of 1k persons

"Joe small" if same people  
"big" if different people

LEARN THE PARAMS OF  
THE NN SUCH THAT

- IF  $x^{(i)}, x^{(j)}$  ARE THE SAME PERSON  $\cdot d(x^{(i)}, x^{(j)}) \Rightarrow \text{SMALL}$
- IF  $x^{(i)}, x^{(j)}$  ARE DIFFERENT PEOPLE  $\cdot d(x^{(i)}, x^{(j)}) \Rightarrow \text{LARGE}$

WE CAN ACCOMPLISH  
THIS WITH THE TRIPLET  
LOSS FUNCTION



**TIP** PRECOMPUTE ENCODINGS  
FOR PPL IN YOUR DB, SO YOU  
DON'T HAVE TO SAVE IMAGES  
& COMPUTE ENCODINGS AT RUN-  
TIME

SOME BIG COMPANIES  
HAVE ALREADY TRAINED  
NETWORKS ON LARGE  
AMTS OF PHOTOS SO  
YOU MAY JUST  
WANT TO REUSE  
THEIR WEIGHTS

© TessFerrandez

# NEURAL STYLE TRANSFER



WE CAN VISUALIZE WHAT A NETWORK LEARNS BY LOOKING AT WHAT IMAGES (PARTS) ACTIVATED EACH UNIT MOST



BUT HOW DOES THIS HELP US GENERATE AN IMAGE IN THE STYLE OF ANOTHER?

## IDEA:

1. GENERATE A RANDOM IMAGE
  2. OPTIMIZE THE COST FUNCTION
- $$J(G) = \alpha J_{\text{CONTENT}}(C, G) + \beta J_{\text{STYLE}}(S, G)$$
- find  $G$  that minimizes  $J(G)$
- if  $C \sim G$  then they are "similar"  
if  $S \sim G$  then they are "similar"
3. UPDATE EACH PIXEL

## CONTENT COST FUNCTION

- USE A PRE-TRAINED CONVNET (ex VGG)
- SELECT A HIDDEN LAYER SOMEWHERE IN THE MIDDLE (not too deep or shallow)  
LATER → COPIES LARGER FEATURES
- LET  $a^{(l)(c)}$  &  $a^{(l')(c')}$  BE THE ACTIVATIONS
- IF  $a^{(l)(c)} \approx a^{(l')(c')}$  ARE SIMILAR THEY HAVE SIMILAR CONTENT  
BECAUSE THEY BOTH TRIGGER THE SAME HIDDEN UNITS

HOW DO WE TELL IF THEY ARE SIMILAR?

$$J_{\text{CONTENT}}(C, G) = \frac{1}{2} \| a^{(l)(c)} - a^{(l')(c')} \|_F^2$$

## CAPTURING THE STYLE

USING THE STYLE IMG AND THE ACTIVATIONS IN A LAYER. LOOK THROUGH THE ACTIVATIONS IN THE DIFFERENT CHANNELS TO SEE HOW CORRELATED THEY ARE

WHEN WE SEE PATTERNS LIKE THIS DO WE USUALLY SEE IT WITH PATCHES LIKE THESE?



⇒ if "yes" then they are "correlated"  
if "no" then they are "unrelated"

$a^{(l)(i,j,k)} = \text{activation at } (i,j,k)$

## STYLE MATRIX

CREATE A MATRIX OF HOW CORRELATED THE ACTIVATIONS ARE, FOR EACH POS  $(x,y)$  & CHANNEL PAIR  $(k,k')$  FOR THE STYLE IMG & GENERATED

$$G_{kk'} = \sum_{i=1}^{n_h} \sum_{j=1}^{n_w} a^{(l)(i,j,k)} \cdot a^{(l)(i,j,k')}$$

## THE STYLE COST FUNCTION

$$J(S, G) = \| G^{(S)} - G^{(G)} \|_F^2$$

FROBENIUS NORM

TO GET MORE VISUALLY PLEASING IMAGES IF YOU CALC  $J(S, G)$  OVER MULTIPLE LAYERS



$$J_{\text{style}}(S, G) = \sum_l \lambda_l J_{\text{style}}^{(l)}(S, G)$$

## 5

# RECURRENT NEURAL NETWORKS

## SEQUENCE MODELS • COURSERA

Because  
of this  $x$  and  $y$   
right  $\leftarrow$   
different  
length

## SEQUENCE PROBLEMS

IN	OUT	PURPOSE
THE QUICK BROWN FOX JUMPED...	THE QUICK BROWN FOX JUMPED...	SPEECH RECOGNITION
		MUSIC GENERATION
THERE IS NOTHING TO LIKE IN THIS MOVIE	★ ★ ★	SENTIMENT CLASSIFICATION
AGCCCTTGAGGAACTAG	AGCCCTTGAGGAACTAG	DNA SEQUENCE ANALYSIS
Vouslez-vous chanter avec moi?	DO YOU WANT TO SING WITH ME?	MACHINE TRANSLATION
	RUNNING	VIDEO ACTIVITY RECOGNITION
Yesterday Harry Potter met Hermoine Granger	Yesterday Harry Potter met Hermoine Granger	NAME ENTITY RECOGNITION

## NAME ENTITY RECOGNITION

output per input words

$$(x^{(1)}, x^{(2)}, \dots, x^{(T)}) \rightarrow y = (y^{(1)}, y^{(2)}, \dots, y^{(T)})$$

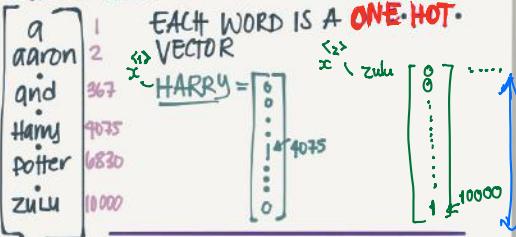
GRANGER INVENTED A NEW SPELL

$$y = \begin{matrix} 1 & 1 & 0 & \dots & 1 \end{matrix} \quad T_y = T_x$$

EXAMPLE OF A PROBLEM WHERE EVERY  $x^{(i)}$  HAS AN OUTPUT  $y^{(i)}$

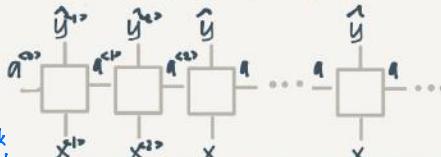
## HOW DO WE REPRESENT WORDS?

CREATE A VOCABULARY (EG 10K MOST COMMON WORDS IN YOUR TEXTS • OR DOWNLOAD EXISTING)



## DIFFERENT TYPES OF RNN

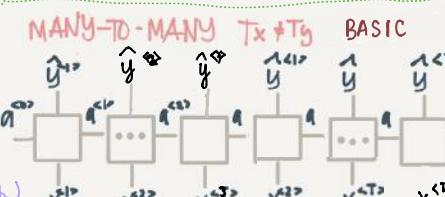
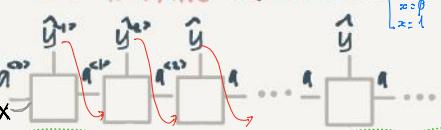
### MANY-TO-MANY $T_x = T_y$



### MANY-TO-ONE $T_x \neq T_y$

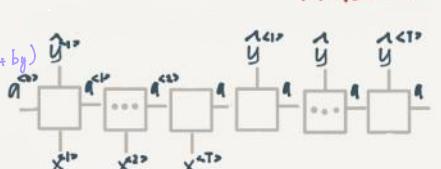


### ONE-TO-MANY $T_x = T_y$



OK

### TRANSLATION



LessTerrible

## RECURRENT NEURAL NET (RNN)



PREVIOUS RESULTS ARE PASSED IN AS INPUTS SO WE GET CONTEXT:

$$\begin{aligned} a^{(t+1)} &= g(W_{aa} a^{(t)} + W_{ay} y^{(t)} + b_a) \\ y^{(t+1)} &= g(W_{ya} a^{(t+1)} + b_y) \end{aligned}$$

Sigmoid

THE SAME  $W$  &  $b$  ARE USED IN ALL TIME STEPS

THE LOSS WE OPTIMIZE IS THE SUM OF  $\mathcal{L}(\hat{y}, y)$  FROM 1-T

## MORE ON RNNs

## LANGUAGE MODELLING

HOW DO YOU KNOW IF SOMEONE IS AID



- THE APPLE AND PAIR SALAD OR P  
 THE APPLE AND PEAR SALAD? PE

## THE PURPOSE OF A LIANG. MODEL CALCULATE THE PROBABILITIES

feed the data, it will calculate what is P for first word to be from the corpus you created

SO GIVEN: CATS AVERAGE 15 WHAT IS THE PROB.  
THE NEXT WORD IS HOURS?

## SAMPLING SENTENCES

1. TRAIN ON ALL HARRY POTTER BOOKS.
  2. RANDOMLY SELECT A WORD (ON OF THE TOP WORDS)  
*(EX. THE)*
  3. PASS THIS INTO THE NEXT TIMESTAMP  
AND SAMPLE A NEW WORD
  4. REPEAT UNTIL X WORDS OR YOU  
REACHED <EOS>

*CAN DO AT CHARACTER LEVEL AS WELL*



THE CAT WHO ALREADY ATE APPLES AND ORANGES  
AND A FEW MORE THINGS BUT IT WAS FULL  
THE CATS, WHO ALREADY ATE ...  
... WERE FULL

NEED TO REMEMBER  
SING/PLURAL FOR A LONG  
TIME

SINCE LONG SENTENCE  $\Rightarrow$  DEEP RNN  
WE GET THE VANISHING GRADIENTS PROB. WE  
HAVE IN STANDARD NNs - I.E. THE GRADIENTS  
FOR CAT/CATS HAVE LITTLE OR NO EFFECT  
ON WAS/WERE.  $\rightarrow$  then we need to figure  
out how to fix it - that gradient

**NOTE**) SOMETIMES YOU SEE EXPLODING GRAD  
(AS OVERFLOW NAN) BUT THIS IS EASILY FIXED  
WITH GRADIENT CLIPPING

## GATED RECURRENT UNIT

HELPS RECALL IF CAT WAS SING.  
OR PLURAL (go to d1 2 note)



THE GRU ACTS AS A MEMORY  
— AT EVERY TIMESTEP IT  
CALCULATES A NEW  $\tilde{c}$  TO STORE  
AND A GATE  $\Gamma_b$  DECIDES TO  
UPDATE  $c$  TO  $\tilde{c}$  OR NOT

you can reject  
any word that come out as **<UNK>**

YAY! YOU ARE NOW  
YOUR OWN J.K. ROWLING

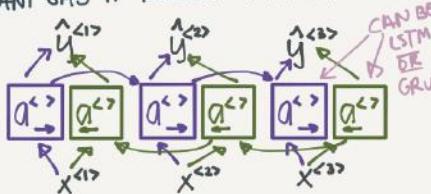
## LONG SHORT TERM MEMORY (LSTM) ( $a_t$ to $f_t, g_t$ , $c_t$ )

THE LSTM IS A VARIATION ON  
THE SAME THEME AS GRU  
BUT WITH AN ADDITIONAL F  
**FORGET GATE**

→ usually BRNN + LSTM is the best

HE SAID, 'TEDDY BEARS ARE ON SALE'  
HE SAID, 'TEDDY ROOSEVELT WAS A  
GREAT PRESIDENT'

PROBLEM: WITHOUT LOOKING FORWARD WE  
CAN'T SAY IF TEDDY IS A TOY OR A NAME



ONE DISADVANTAGE IS THAT YOU NEED THE FULL SENTENCE BEFORE YOU BEGIN.  
SO NOT SUITABLE FOR LIVE SPEECH REC.

The diagram illustrates a Deep Recurrent Neural Network (RNN) across four time steps. Each time step  $t$  consists of an input  $x^{(t)}$ , a hidden state  $y_t^{(t)}$ , and a cell state  $c_t^{(t)}$ . The hidden states are connected sequentially via arrows labeled  $W_h$  and  $b_h$ . The initial hidden state  $y_0^{(0)}$  and cell state  $c_0^{(0)}$  are also shown. A blue arrow points from the text "SINCE THEY ARE ALREADY MEMOROUS DEEP THEY USUALLY DON'T HAVE 4 LOT OF LAYERS" to the first layer of hidden states.

# SEQUENCE MODELS • COURSERA

## NLP & WORD EMBEDDINGS

MAN IS TO WOMAN AS  
KING IS TO QUEEN

PROBLEM: THE ONE-HOT REPR OF AN APPLE HAS NO INFO ABOUT ITS RELATIONSHIP TO OTHER FRUIT

I WANT A GLASS OF ORANGE —  
I WANT A GLASS OF APPLE —

SOLUTION: CREATE A MATRIX OF FEATURES TO DESCRIBE THE WORDS

## WORD EMBEDDINGS

	MAN	WOMAN	KING	QUEEN	APPLE	ORANGE
5391	9853	4914	3157	456	6257	
GENDER	-1	1	-0.95	0.97	0.00	0.01
ROYAL	0.01	0.02	0.93	0.95	-0.01	0.00
AGE	0.03	0.02	0.7	0.69	0.03	-0.02
FOOD	0.04	0.01	0.02	0.01	0.95	0.97
:						
e <sub>5391</sub>						

IN REALITY • THE FEATURES ARE LEARNED & NOT AS STRAIGHTFWD AS GENDER/AGE



t-SNE makes this happen  
VISUAL REPRESENTATION OF 300D WORD EMBEDDINGS (300P → 2D)

FIND(W): high degree of similarity between ARG. MAX SIM( $e_{W_1}, e_{KING} + e_{MAN} + e_{WOMAN}$ )  
 $SIM(U, V) = \frac{U^T V}{\|U\|_2 \|V\|_2}$  COSINE SIMILARITY  
 After the t-SNE mapping, the link won't apply to them

## USING WORD EMBEDDINGS

EX. NAME/ENTITY RECOGN



WITH WORD EMBEDDINGS WE UNDERSTAND THAT AN ORANGE FARMER IS A PERSON ⇒ SALLY SMITH = NAME

- APPLE ~ ORANGE ⇒ APPLES FARMER = PERSON
- USING WORD EMBEDDINGS TRAINED ON LOTS OF TEXT WE ALSO GET EMB FOR MORE UNCOMMON WORDS (DURIAN, CULTIVATOR)

EX. MAN IS TO WOMAN AS KING IS TO ?

EMBEDDING MATRIX

Analogies

e <sub>man</sub>	—	MAN	WOMAN	KING	QUEEN
5391	9853	4914	3157	456	6257
GENDER	-1	1	-0.95	0.97	0.00
ROYAL	0.01	0.02	0.93	0.95	-0.01
AGE	0.03	0.02	0.7	0.69	0.03
FOOD	0.04	0.01	0.02	0.01	0.95
:					
e <sub>5391</sub>					

e<sub>man</sub>

e<sub>woman</sub>

e<sub>king</sub>

e<sub>queen</sub>

e<sub>man</sub> - e<sub>woman</sub> ≈ e<sub>king</sub> - e<sub>queen</sub>

EVERY SIMILAR

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

try to find the word that has e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

e<sub>woman</sub> - e<sub>man</sub> ≈ e<sub>queen</sub> - e<sub>king</sub>

word

that has

- Let's use our previous example: "I want a glass of orange juice to go along with my cereal."
- We will choose a context and a target from the choices we have mentioned in the previous sections.
- Then we will calculate this for every pair:  $x_{ij} = \#$  times  $i$  appears in context of  $j$ .
- $x_{ij} = 0$ , if we choose a window pair, but they will not equal if we choose the previous words for example, in GloVe they use a window which means they are equal.

## IDEA: GLOVE WORD VECTORS

$c, t$  (connect with each other)

**TARGET**  
 $x_{ij}^t = \#$  TIMES WORD  $j$  APPEARS IN THE CONTEXT OF  $i^t$   
 (HOW RELATED THEY ARE)

$$\text{MINIMIZE } \sum_{i=1}^{10k} \sum_{j=1}^{10k} f(x_{ij}) (\theta_i^T e_j + b_i + b_j - \log x_{ij})^2$$

if no context  
 (also helps weighing very freq words (the, of, etc) & very infrequent (durian))

EVERYTHING LED UP TO THIS VERY SIMPLE ALGORITHM

- Definition of word embeddings
- If it is good for you, it is good for everyone else
- If you have enough data, you can try to implement one of the available experiments
- Because word embeddings are very computationally expensive to train, most practitioners will load a pre-trained set of embeddings
- A final note that you can't guarantee that the axis used to represent the features will be well aligned with what people think feature words are like, so don't overfit

## SENTIMENT CLASSIFICATION

X	y
THE DESSERT IS EXCELLENT	★★★▲
SERVICE WAS QUITE SLOW	★★
GOOD FOR A QUICK MEAL BUT NOTHING SPECIAL	★★★
COMPLETELY LACKING IN GOOD TASTE, GOOD SERVICE AND GOOD AMBIENCE	★

PROBLEM: YOU MAY NOT HAVE A LARGE DATASET  
 BUT YOU CAN USE AN EMBEDDING MATRIX  $E$  THAT IS ALREADY PRE-TRAINED

## IDEA: SIMPLE CLASSIFICATION

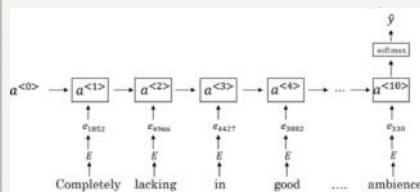


WORKS WELL FOR SHORT SENTENCES  
 BUT DOESN'T TAKE ORDER INTO ACCOUNT

"COMPLETELY LACKING IN GOOD TASTE,  
 GOOD SERVICE AND GOOD AMBIENCE"

THIS MAY BE SEEN AS A +++ REVIEW

## IDEA: USE AN RNN



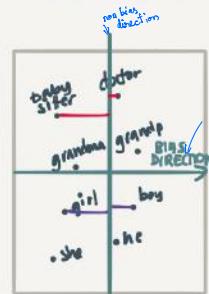
THIS CAN NOW TAKE INTO  
 ACCOUNT THAT COMPLETELY LACKING  
 NEGATES THE WORD GOOD

## ELIMINATING BIAS IN WORD EMBEDDINGS

MAN IS TO COMPUTER PROGRAMMER AS  
 WOMAN IS TO HOME MAKER X

SOMETIMES THE TEXT CONTAINS SENSES WHICH LEARN A GENDER, RACE, AGE... BIAS WE DON'T WANT OUR MODELS TO HAVE. EX. HIRING BASED ON GENDER, SENTENCING BASED ON RACE ETC.

## ADDRESSING BIAS



### 1. IDENTIFY BIAS DIRECTION

She → Cshe → find the average  
 Example: female → Efemale

### 2. NEUTRALIZE

FOR EVERY WORD THAT IS NOT DEFINITIONAL (GIRL, BOY, HE, SHE...) PROJECT TO GET RID OF BIAS

### 3. EQUALIZE PAIRS

THE ONLY DIFF BETWEEN EX GIRL/BOY SHOULD BE GENDER

Making sure that words like "girl" and "boy" are both exactly the same distance, or roughly same distance, from words that are gender neutral ("dishes", "baby stroller", ...)

HOW DO YOU KNOW WHICH WORDS TO NEUTRALIZE?

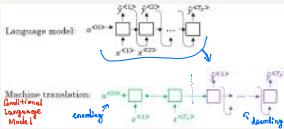
DOCTOR, BEARD, SEWING MACHINE?

A: BY TRAINING A CLASSIFIER TO FIND OUT IF A WORD IS DEFINITIONAL

URNS OUT THE # OF PAIRS IS FAIRLY SMALL SO YOU CAN EVEN HAND PICK THEM

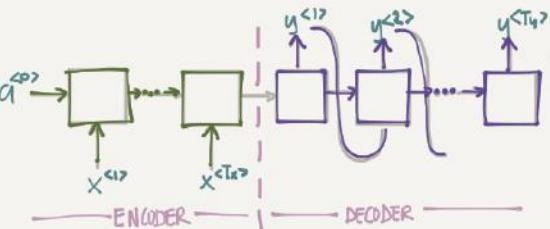
# SEQUENCE MODELS • COURSERA

# SEQUENCE TO SEQUENCE



## BASIC MODELS

JANE VISITE L'AFRIQUE  
EN SEPTEMBRE  $\rightarrow$  JANE IS VISITING AFRICA  
IN SEPTEMBER



$\rightarrow$  THIS IS A CAT  
ON A CHAIR

CNN  $\xrightarrow{\quad}$  RNN

HOW DO YOU PICK THE MOST LIKELY SENTENCE?

$$P(y^{<1>} | \dots | y^{<T_y>} | x^{<1>} | \dots | x^{<T_x>})$$

WE DON'T WANT A RANDOMLY GENERATED SENTENCE  
(WE WOULD SOMETIMES GET A GOOD, SOMETIMES BAD)  
INSTEAD WE WANT TO MAXIMIZE

$$\arg \max_{y^{<1>} \dots y^{<T_y>}} P(y^{<1>} \dots y^{<T_y>} | x)$$

## IDEA: USE GREEDY SEARCH

- PICK THE WORD WITH THE BEST PROBABILITY
  - REPEAT UNTIL CEOD
- WITH THIS WE COULD GET ✓

- JANE IS GOING TO BE VISITING AFRICA  
THIS SEPTEMBER

INSTEAD OF

- JANE IS VISITING AFRICA THIS SEPTEMBER

SOLUTION OPTIMIZE THE PROB OF THE WHOLE SENTENCE INSTEAD

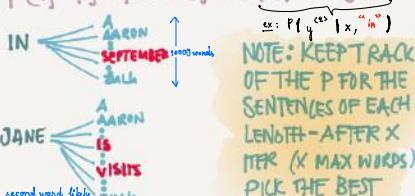
## BEAM SEARCH

- PICK THE FIRST WORD

PICK THE B (EX 3) BEST ALTERNATIVES  
(IN, JANE, SEPTEMBER)  $\rightarrow$  the 3 most likely first words in the translation

- FOR EACH B WORDS PICK THE NEXT WORD AND EVALUATE THE PAIRS TO END UP IN B PAIRS

$$P(y^{<1>} | \dots | y^{<2>} | x) = P(y^{<1>} | x) P(y^{<2>} | y^{<1>})$$



Because the second word belongs to the other two & the previous hypothesis is no longer the potential first word.

SEPTEMBER  $\xrightarrow{\text{only all this we sort down to the B again}}$  (IN SEPTEMBER, JANE IS, JANE VISITS)

- REPEAT TIL DONE

$$\arg \max_{y^{<1>} \dots y^{<T_y>}} \prod_{t=1}^{T_y} P(y^{<t>} | x, y^{<1>} \dots y^{<t-1>})$$

## OVERFLOWS

PROBLEM: MULTIPLYING PROBABILITIES  $(O(K^n))$   
RESULTS IN A VERY SMALL NUMBER

PROBLEM: IF WE OPTIMIZE FOR THE MULT WE WILL PREFER SHORT SENTENCES. SINCE EACH WORD WILL REDUCE PROB

INSTEAD WE CAN OPTIMIZE FOR THIS

$$\frac{1}{T_y} \sum_{t=1}^{T_y} \log(P(y^{<t>} | x, y^{<1>} \dots y^{<t-1>}))$$

$\downarrow$  a complete normalization by length  
 $\downarrow$  d = 0 not normalizing at all  $\Rightarrow d=97$  iteration  
HOW DO WE PICK B?

$\downarrow$  can get diminishing return

LARGE B: BETTER RESULT, SLOWER  
SMALL B: WORSE RESULT, FASTER

IN PRD YOU MIGHT SEE B=10.  
100 IS PROBABLY A BIT TOO HIGH -  
BUT ITS DOMAIN DEPENDENT

## ERROR ANALYSIS IN BEAMS.

HUMAN: JANE VISITS AFRICA IN SEPT...  $y^*$   
ALSO: JANE VISITED AFRICA LAST SEPTEMBER  $y^y$

HOW DO WE KNOW IF ITS OUR RNN OR OUR BEAM SEARCH WE SHOULD WORK ON?

LET THE RNN GIVE  $P(y^*) = P(y^*, x)$  &  $P(y^y) = P(y^y, x)$

IF  $P(y^*) > P(y^y)$ : beam search chose  $y^y$ . But  $y^*$  attains higher  $P(y^*, x)$   
BEAM PICKED THE WRONG ONE  
TRY A HIGHER B

ELSE:  $P(y^*) > P(y^y)$   $y^*$  is the better translation than  $y^y$ . But RNN predicted  $P(y^y, x) > P(y^*, x)$   
THE RNN PICKED THE WRONG PROBS. SO FOCUS ON THE RNN  
 $\Rightarrow$  RNN is fault

Human	Algorithm	$P(y^*, x)$ P(y^y, x)	Actual?
Jane visits Africa in Sept...	Beam search	$\frac{1}{10} \log(10) \approx 0.43$	✓

Fernandez

# SEQUENCE TO SEQUENCE

FRENCH: LE CHAT EST SUR LE TAPIS  
 HUMAN1: THE CAT IS ON THE MAT  
 HUMAN2: THERE IS A CAT ON THE MAT

HOW DO YOU EVALUATE THE MACHINE TRANSLATION WHEN MULTIPLE ARE RIGHT?

BLEU SCORE (\*optional)

IDEA: CHECK IF THE WORDS APPEAR IN THE REAL TRANSLATION

THE THE THE THE THE THE THE  
 SCORE: 7/7

IDEA: ONLY GIVE CREDIT FOR A WORD THE MAX # TIMES IT APPEARS IN A TARGET SENTENCE  
 SCORE: 2/7 COUNT CUP

THE CAT THE CAT ON THE MAT

COUNT	COUNT CLIP
2	1
1	0
1	1
1	1
6	4

BI-GRAM SCORE: 4/6  
 ↓ on the 2 references we have  
 BIGRAMS

## COMBINED BLEU SCORE

$$BP \cdot \exp\left(\frac{1}{4} \sum_{n=1}^4 P_n\right)$$

$P_1$  = SCORE SINGLE WORD  
 $P_2$  = SCORE BIGRAMS  
 ...

BP = BREVITY PENALTY

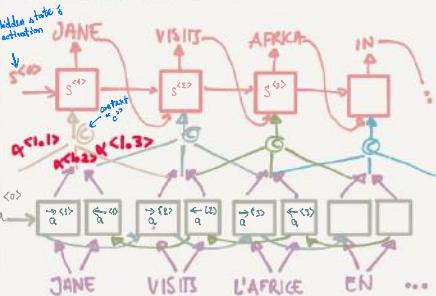
PENALIZES SENTENCES SHORTER THAN THE TARGET

A USEFUL SINGLE NUMBER EVAL METRIC

## ATTENTION MODEL



SOLUTION: TRANSLATE A LITTLE AT A TIME USING ONLY PARTS OF THE SENTENCE AS CONTEXT



$\alpha^{<t,t>}$  = HOW MUCH ATTENTION  $y^{<t>}$  SHOULD PAY TO  $x^{<t>}$

$$\alpha^{<2>} = \sum_{t'} \alpha^{<2,t>} \cdot a^{<t>} \quad | \quad \sum_t \alpha^{<1,t>} = 1$$

attention weight

$\alpha$  IS CALCULATED USING A SMALL NEURAL NETWORK



English word  $y^{<t>}$

French word  $x^{<t>}$

$$\alpha^{<2,t>} = \sum_{t'} \alpha^{<2,t>} \cdot a^{<t>} \quad | \quad \sum_t \alpha^{<1,t>} = 1$$

$\alpha$  IS CALCULATED USING A SMALL NEURAL NETWORK

$$\alpha^{<t,t>} = \exp(e^{<t,t>}) / \sum_{t'=1}^{10} \exp(e^{<t,t>})$$

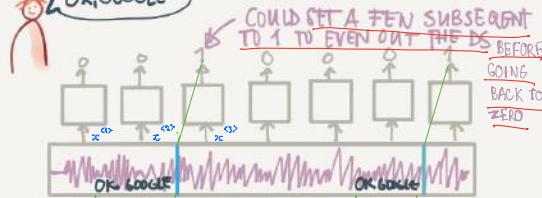
## SPEECH RECOGNITION



PROBLEM: 10s CLIP AT 100Hz = 1000 INPUTS BUT ONLY  $\approx 20$  OUTPUTS

SOLUTION: USE CTC COST (CONNECTION TEMPORAL CLASSIFICATION)  
 ttt-h-eee---u---999 →  
 COLLAPSE REPEATED CHARS NOT SEP BY BLANK

## TRIGGER WORD DETECTION



COULD SET A FEN SUBSEQUENT TO 1 TO EVEN OUT THE DS BEFORE GOING BACK TO ZERO