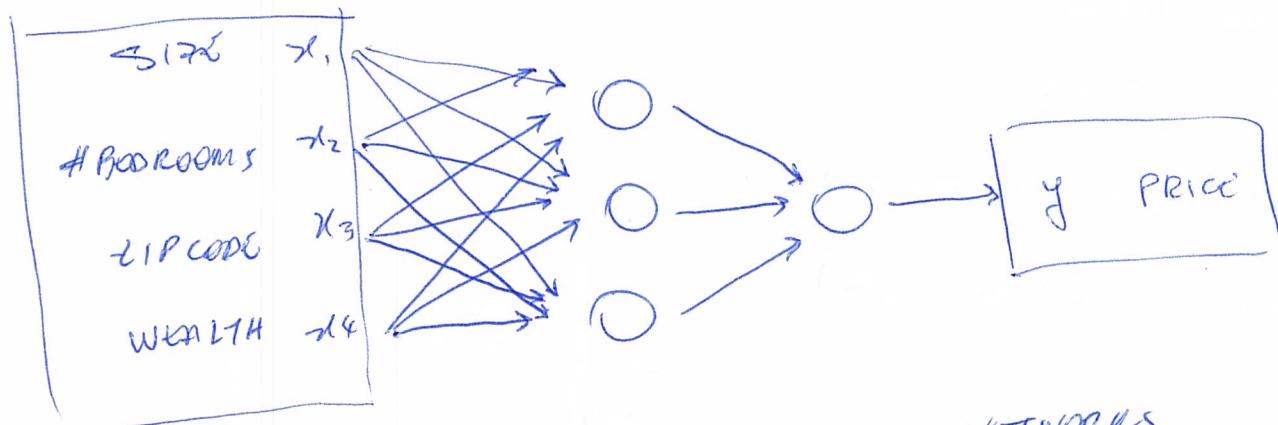
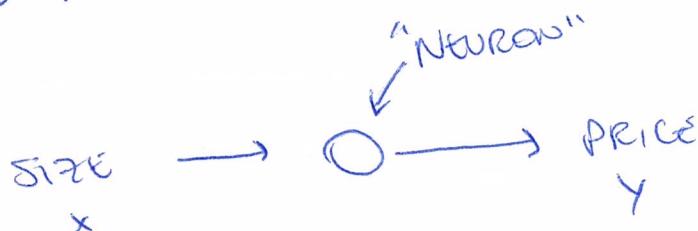
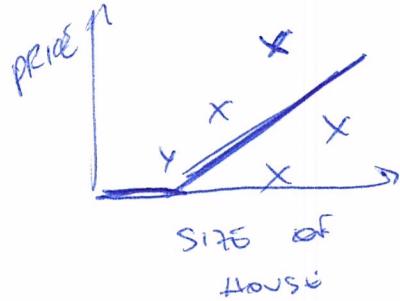


DEEP LEARNING • AI

COURSE 1 - NEURAL NETWORKS AND DEEP LEARNING

- WHAT IS A NEURAL NETWORK?



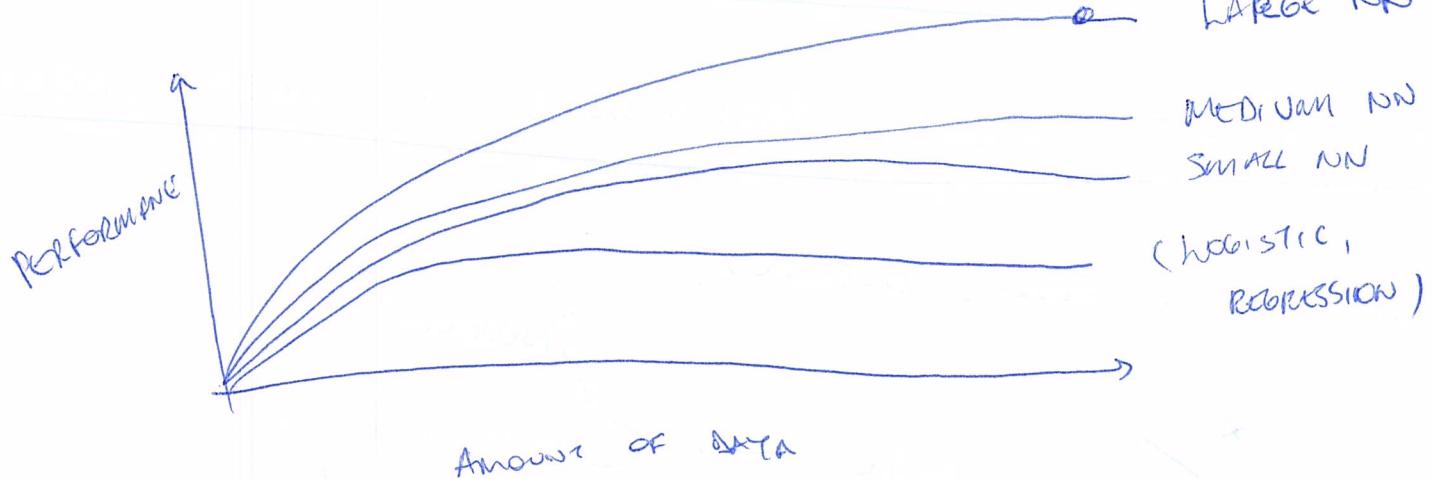
- SUPERVISED LEARNING WITH NEURAL NETWORKS

INPUT (x)	OUTPUT (y)	APPLICATION
Done FEATURES AD, USER INFO	PRICE	REAL ESTATE
IMAGE	CLICK ON AD	ONLINE ADVERTISING
AUDIO	OBJECT (1, ..., 100)	PHOTO TAGGING
ENGLISH	TEXT TRANSCRIPT	SPEECH RECOGNITION
IMAGE, RADAR INFO	CHINESE POSITION OF OTHER CARS	MACHINE TRANSLATION AUTONOMOUS DRIVING

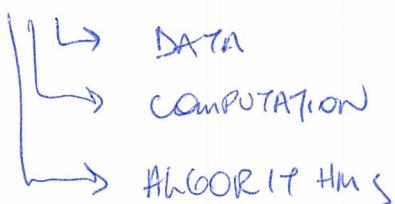
Legend:

- REAL ESTATE, ONLINE ADVERTISING, PHOTO TAGGING, SPEECH RECOGNITION, MACHINE TRANSLATION, AUTONOMOUS DRIVING: } CNN
- RNN: }
- Custom / Deep Hybrid NN: }

- WHY IS DEEP LEARNING TAKING OFF?

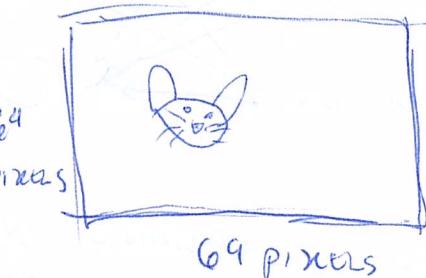
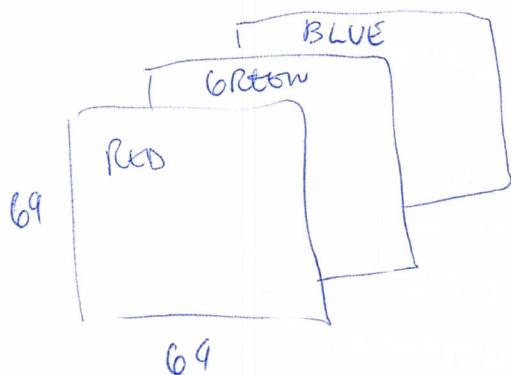


SCALE DRIVES DEEP LEARNING PROGRESS



- BINARY CLASSIFICATION

1 (car) vs. 0 (not car)



$$x = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 255 \\ 134 \\ \vdots \\ 255 \\ 134 \\ \vdots \end{bmatrix} \quad n = \text{length } x = 12288$$

$$64 \times 64 \times 3 = 12288$$

$$(x, y) \Rightarrow x \in \mathbb{R}^{n_x} \\ y \in \{0, 1\}$$

TRAINING EXAMPLES: $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

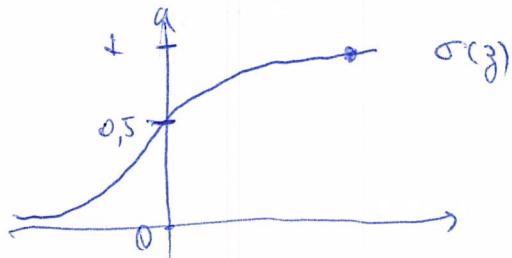
$$X = \begin{bmatrix} | & | & | & \dots & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} & | \\ | & | & \dots & | & \downarrow n_x \end{bmatrix}$$

- logistic Regression

$$x_i, \hat{y} = \text{IP}(y=1/x) , 0 \leq \hat{y} \leq 1 , x \in \mathbb{R}^{n_x}$$

Parameters: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

Output: $\hat{y} = \underbrace{\sigma(w^T x + b)}_{z}, z = w_1 x_1 + w_2 x_2 + b$ (for example)



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

(SIGMOID FUNCTION)

- logistic Regression COST FUNCTION

loss (error) function:

$$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2 \Rightarrow \text{U-shaped}$$

(GLOBAL MAXIMA)

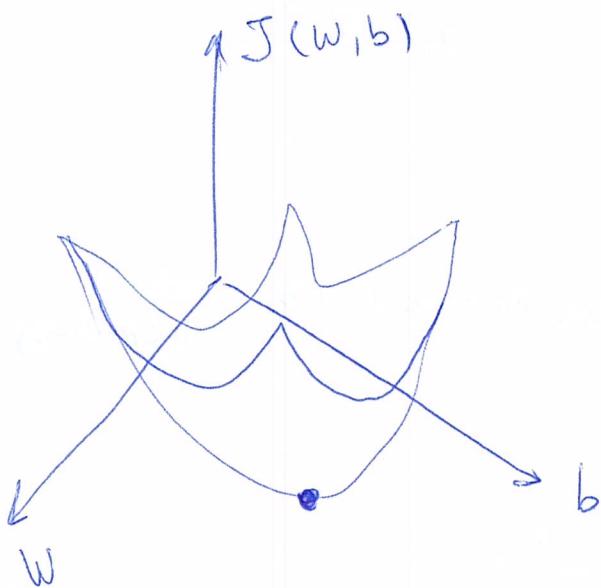
$$L(\hat{y}, y) = - [y \log \hat{y} + (1-y) \log(1-\hat{y})]$$

if $y=1$: $L(\hat{y}, 1) = -\log \hat{y}$ { want $\log \hat{y}$ large \Rightarrow \hat{y} large }

if $y=0$: $L(\hat{y}, 0) = \log(1-\hat{y})$ { want $\log(1-\hat{y})$ large \Rightarrow \hat{y} small }

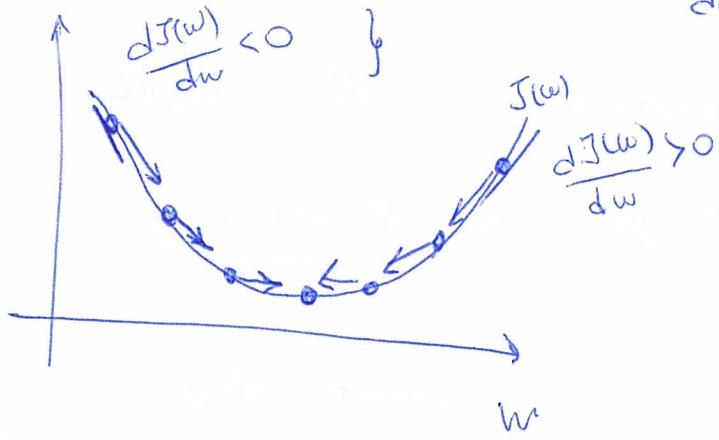
Cost Function: $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$

- GRADIENT DESCENT



REPEAT {

$$w_i := w_0 - \alpha \frac{dJ(w)}{dw}$$



- LOGISTIC REGRESSION ON m EXAMPLES

$$J=0 ; dw_1=0 ; dw_2=0 ; db=0$$

FOR i=1 TO m

$$z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J+ = -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 += x_1^{(i)} dz^{(i)} \quad \downarrow n=2$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

$$J+=m$$

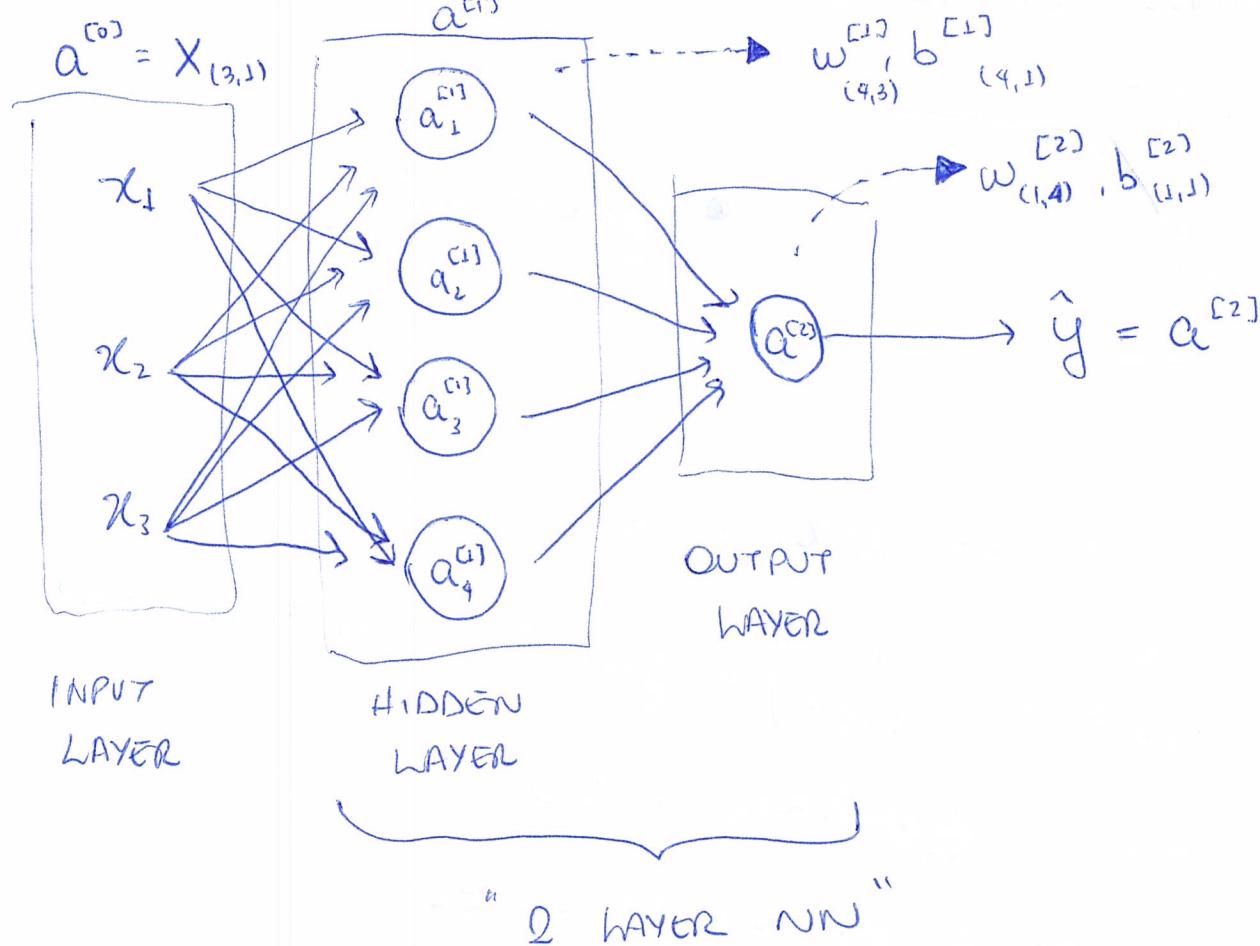
$$dw_1 /= m ; dw_2 /= m ; db /= m$$

-

NEURAL

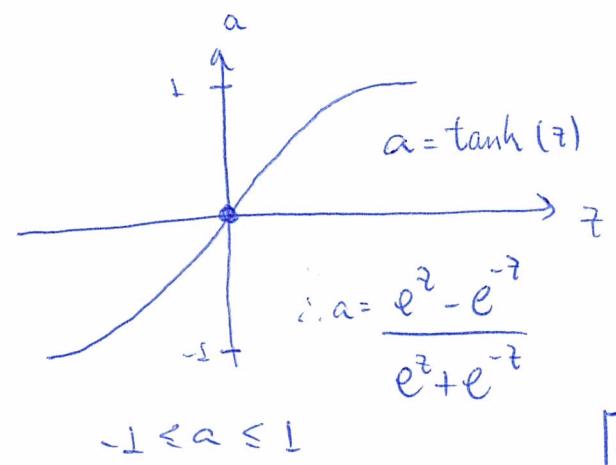
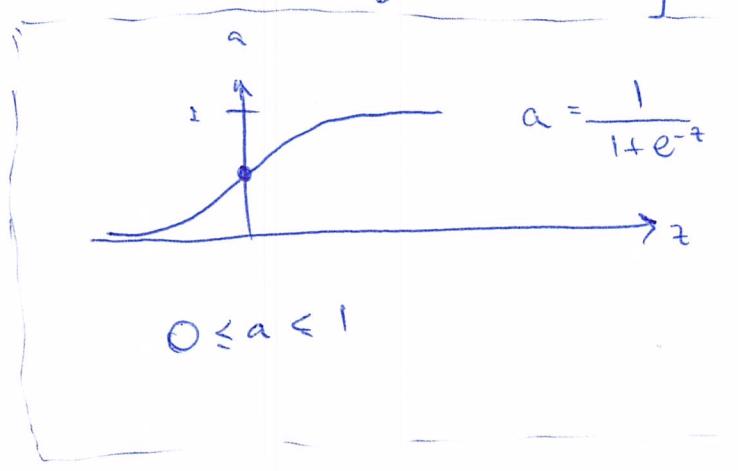
NET WORK

REPRESENTATION

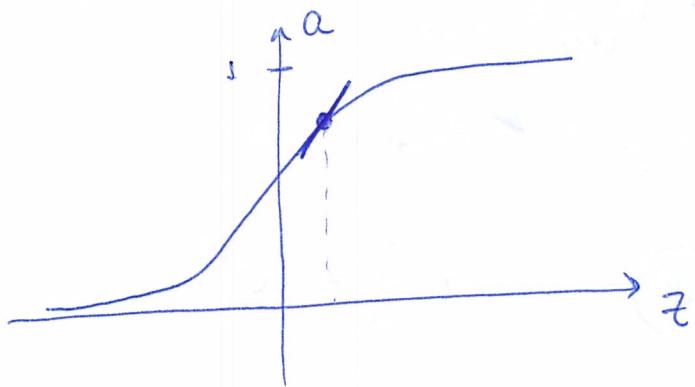


$$z^{[1]} = \begin{bmatrix} -w_1^{[1]T} \\ -w_2^{[1]T} \\ -w_3^{[1]T} \\ -w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} w_1^{[1]T}x + b_1^{[1]} \\ w_2^{[1]T}x + b_2^{[1]} \\ w_3^{[1]T}x + b_3^{[1]} \\ w_4^{[1]T}x + b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

$$\therefore a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix} = \begin{bmatrix} \sigma(z_1^{[1]}) \\ \sigma(z_2^{[1]}) \\ \sigma(z_3^{[1]}) \\ \sigma(z_4^{[1]}) \end{bmatrix} = \sigma(z^{[1]})$$



- DERIVATIVES OF ACTIVATION FUNCTIONS



$$\bullet \quad g(z) = \frac{1}{1+e^{-z}}$$

$\frac{d}{dz} g(z) = \text{Slope of } g(z) \text{ at } z$

$$= \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}} \right) = g(z) \cdot (1-g(z)) = a(1-a)$$

$$\underline{\underline{z=10}}: \quad g(z) \approx 1 \Rightarrow \frac{d}{dz}(g(z)) \approx 1 \cdot (1-1) \approx 0$$

$$\underline{\underline{z=-10}}: \quad g(z) \approx 0 \Rightarrow \frac{d}{dz}(g(z)) \approx 0 \cdot (1-0) \approx 0$$

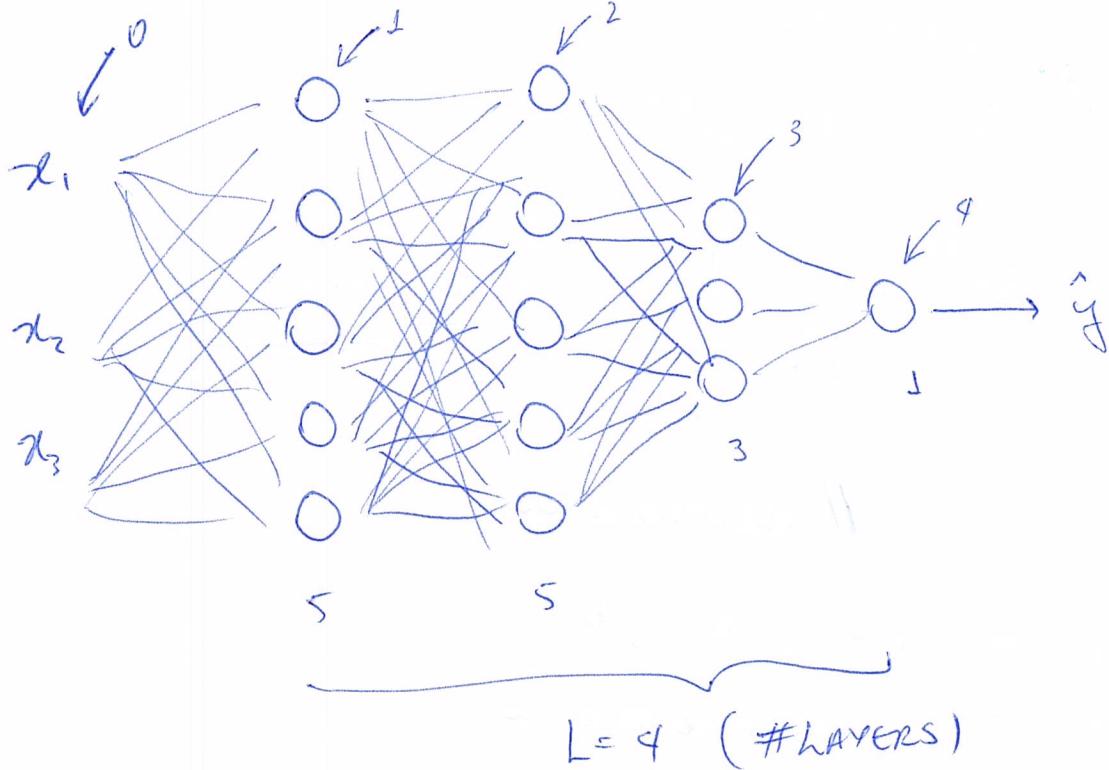
$$\underline{\underline{z=0}}: \quad g(z) = \frac{1}{2} \Rightarrow \frac{d}{dz}(g(z)) = \frac{1}{2} \cdot (1-\frac{1}{2}) = \frac{1}{4}$$

$$g'(z) = \frac{d}{dz} g(z) = a(1-a) \Rightarrow g'(z) = a(1-a)$$

$$\bullet \quad g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad g'(z) = \frac{d}{dz} g(z) = 1 - (\tanh(z))^2$$

$$a = g(z), \quad g'(z) = 1 - a^2$$

- DEEP L-LAYER NETWORK



$n^{[l]}$ = # UNITS IN LAYER l

$a^{[l]}$ = activations in LAYER l

$$n^{[1]} = 5, n^{[2]} = 5, n^{[3]} = 3, n^{[4]} = n^{[0]} = 1, n^{[0]} = 3$$

- FORWARD PROPAGATION IN A DEEP NETWORK

USING THE FAST EXAMPLE : $x = a^{[0]}$

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

$$a^{[1]} = g^{[1]}(z^{[1]})$$

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

$$a^{[2]} = g^{[2]}(z^{[2]})$$

\vdots

$$z^{[4]} = w^{[4]}.a^{[3]} + b^{[4]}$$

$$a^{[4]} = g^{[4]}(z^{[4]}) = \hat{y}$$

$$\left. \begin{array}{l} z^{[l]} = w^{[l]}.a^{[l-1]} + b^{[l]} \\ a^{[l]} = g^{[l]}(z^{[l]}) \end{array} \right\}$$

$$z^{(1)} = W^{(1)} \cdot x + b^{(1)}$$

$$(3,1) \leftarrow (3,2) \cdot (2,1) + (3,1)$$

$$(u^{(1)}, 1) \leftarrow (u^{(1)}, u^{(0)}) (u^{(0)}, 1) + (u^{(1)}, 1)$$

$$\begin{aligned} \therefore dW^{(L)} &= (u^{(L)}, u^{(L-1)}) \\ db^{(L)} &= (u^{(L)}, 1) \end{aligned} \quad \left. \right\} \text{DIMENSIONS}$$

- PARAMETERS VS. HYPERPARAMETERS

PARAMETERS: $W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}, \dots$

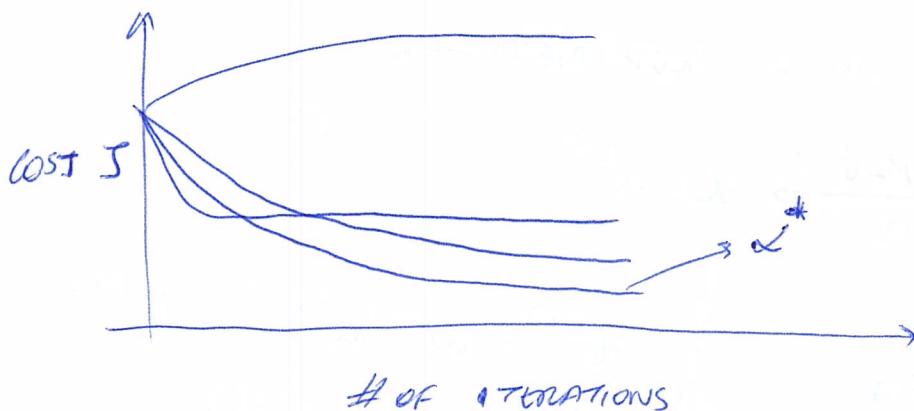
HYPERPARAMETERS: LEARNING RATE α

ITERATIONS

HIDDEN LAYERS L

HIDDEN UNITS $n^{(1)}, n^{(2)}, \dots$

CHOICE OF ACTIVATION FUNCTION



APPLIED DEEP LEARNING
IS A VERY EMPIRICAL
PROCESS

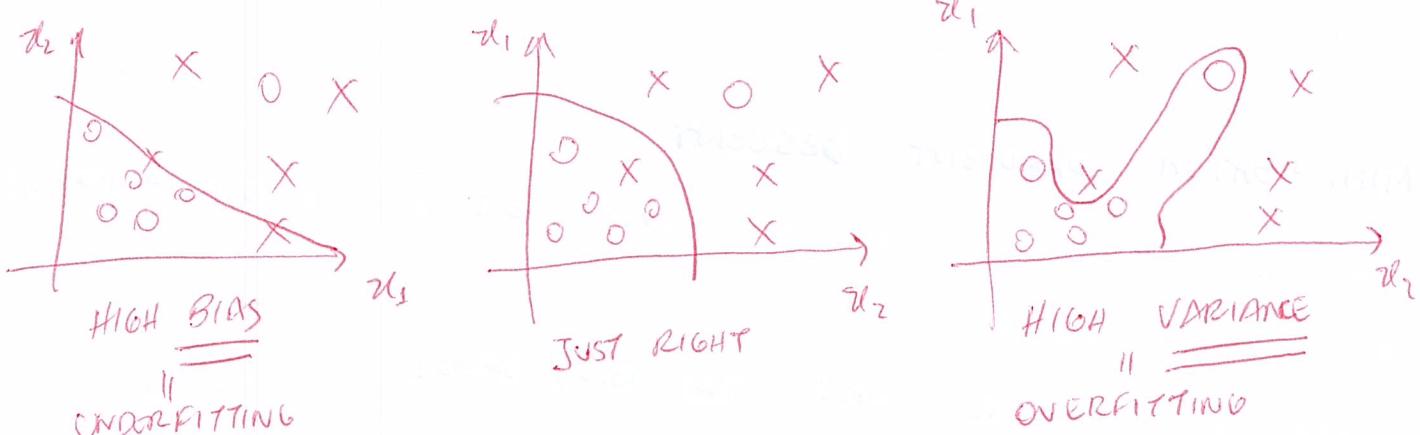
COURSE 2 - IMPROVING Deep Neural Networks

- TRAIN / VALID / TEST SETS SIZES

BUDGET : $\frac{\text{TRAIN}}{70\%} / \frac{\text{TEST}}{30\%}$ OR $\frac{\text{TRAIN}}{60\%} / \frac{\text{VALID}}{20\%} / \frac{\text{TEST}}{20\%}$

BIG DATA ERA : $1,000,000 + \Rightarrow 98\% / 1\% / 1\%$

- BIAS AND Variance



- SOLVE BIAS / VARIANCE PROBLEMS

HIGH BIAS? \rightarrow

- BIGGER NETWORK
- BIGGER TRAINING SET

HIGH VARIANCE? \rightarrow

- More DATA (more VARIABILITY)
- REGULARIZATION \downarrow

- REGULARIZATION (L2)

$$J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = \underbrace{\frac{1}{n} \sum_{i=1}^n L(\hat{y}^{(i)}, y^{(i)})}_{\text{COST FUNCTION}} + \frac{\lambda}{2n} \underbrace{\sum_{l=1}^L \|w^{[l]}\|_F^2}_{\text{REGULARIZATION FUNCTION}}$$

ELIMINATES SOME "EXTRA" NEURONS.

(λ := REGULARIZATION PARAMETER)

- NORMALIZING INPUTS

$x_1: 1 \dots 100$ $x_2: 0 \dots 1$ $\left\{ \begin{array}{l} w_1 \\ w_2 \end{array} \right\} \Rightarrow$ MORE TIME / ITERATIONS FOR GRADIENT DESCENT

- REGULARIZATION WITH DROPOUT

DURING TRAINING ELIMINATES (RANDOMLY) NOISES, USING A VARIABLE OF KEEP-PROB.

* IMPORTANT: ONLY IN TRAINING

- MINI-BATCH GRADIENT DESCENT

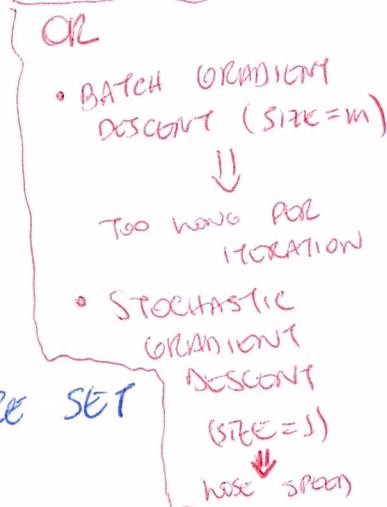
CREATES BATCHES FOR TRAINING SET (m OBSERVATIONS), AND DIVIDES SETS.

MINI-BATCH SIZE NOT TOO BIG / SMALL



FASTEST LEARNING

- VECTORIZATION
- MAKE PROGRESS WITHOUT PROCESSING ENTIRE SET



IF SMALL TRAINING SET: USE BATCH GRADIENT DESCENT
($m \leq 2000$)

TYPICAL SIZES OF MINI-BATCHES: $64, 128, 256, 512$
 $2^6, 2^7, 2^8, 2^9$

- ADAM OPTIMIZATION ALGORITHM

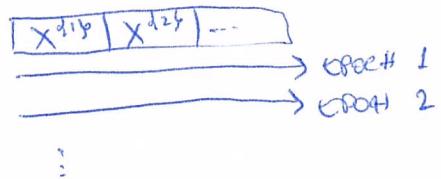


"GRADIENT DESCENT Momentum"

ADAM := ADAPTED MOMENT ESTIMATION

- LEARNING RATE DECAY

1 EPOCH = 1 PASS THROUGH THE DATA



$$\alpha = \frac{1}{1 + \text{Decay-Rate} * \text{Epoch-Num}} \alpha_0$$

$$\alpha_0 = 0.2$$

DECAY-RATE = 1

EPOCH	α
1	0.1
2	0.67
3	0.5
:	:

- HYPERPARAMETER TUNING

1^o α

2^o β

$\beta_1, \beta_2, \epsilon$
 $0.9 \quad 0.999 \quad 10^{-8}$

3^o # LAYERS

2^o # HIDDEN UNITS

3^o LEARNING RATE DECAY

2^o MINI-BATCH SIZE

BABYSITTING ONE MODEL

VS. TRAINING MANY MODELS IN PARALLEL

BATCH NORMALIZATION

FEATURE NORMALIZATION \Rightarrow SPEED UP LEARNING

USING BATCH NORMALIZATION:

- WE CAN USE HIGHER LEARNING RATES, CAUSE ACTIVATION VALUES ARE LOW
- IT REDUCES OVERFITTING BECAUSE IT HAS A SLIGHT REGULARIZATION EFFECT (ADD NOISE TO EACH HIDDEN LAYER)

SOFT MAX REGRESSION (MULTI-CLASS CLASSIFICATION)

1 - CAT

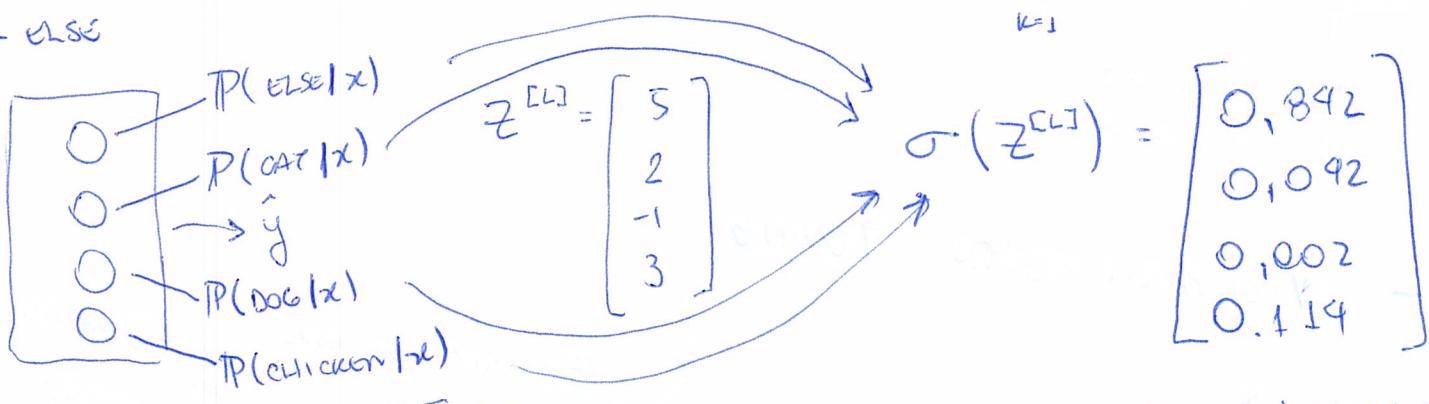
2 - DOG

3 - CHICKEN

0 - ELSE

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

$$\sigma(z)_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}, j=1, \dots, K$$



"HARD MAX" \Rightarrow $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ AND IF $C=2 \Rightarrow$ SOFT MAX = LOGISTIC REGRESSION

LOSS FUNCTION

$$L(\hat{y}, y) = - \sum_{j=1}^4 y_j \cdot \log \hat{y}_j$$

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$= - \log \hat{y}_2$$

$$\text{MIN } L(\hat{y}, y) = \text{MIN} - \log \hat{y}_2 \Rightarrow \text{MAX } \hat{y}_2$$

COURSE 3 - STRUCTURING ML PROJECTS

WAYS TO ANALYZE PROBLEMS OF ML IN A WAY TO FIND SOLUTIONS (GOOD ONES).

CHAIN OF ASSUMPTIONS IN ML:

- FIT TRAINING SET WELL
- " DEV " "
- " TEST " "
- PERFORMS WELL IN REAL WORLD

USING A SINGLE NUMBER EVALUATION METRIC:

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

$$F1 \text{ Score} = \frac{2}{\frac{1}{P} + \frac{1}{R}} \quad \text{"HARMONIC MEAN"}$$

↳ TO COMPARE PRECISION AND RECALL

IN ONLY ONE NUMBER

		PREDICTED	
		NEGATIVE	POSITIVE
ACTUAL	NEGATIVE	TN	FP
	POSITIVE	FN	TP

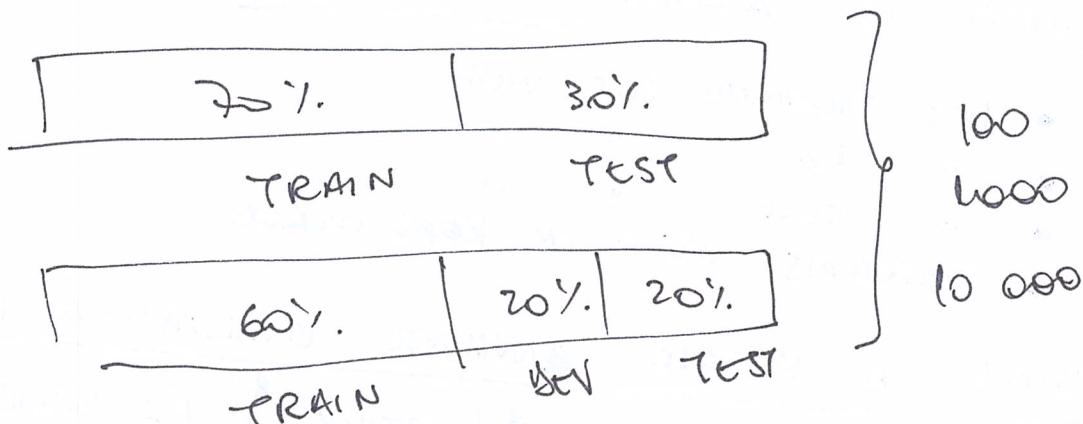
$$\text{PRECISION} = \frac{TP}{TP + FP}$$

$$\text{RECALL} = \frac{TP}{TP + FN} \quad (\text{SENSITIVITY})$$

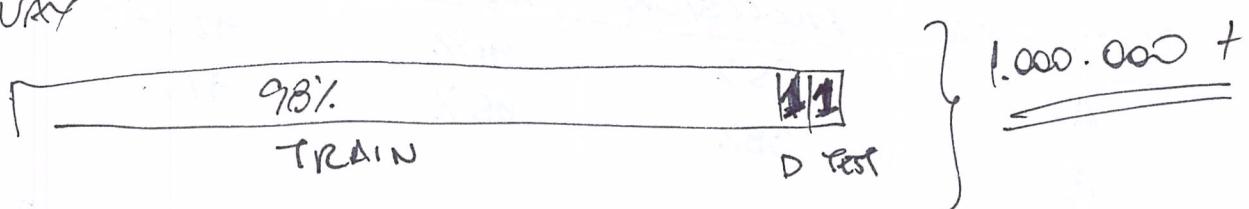
$$\text{SPECIFICITY} = \frac{TN}{TN + FP}$$

Size of dev/test set

- Old way of splitting



- New way



When to change dev/test sets and metrics

→ important to analyze error public over

TAX ERROR

$$\text{ERROR} = \frac{1}{\sum_{i=1}^{M_{\text{dev}}} w^{(i)}} \sum_{i=1}^{M_{\text{dev}}} w^{(i)} \cdot \begin{cases} 1 & \text{if } y_{\text{pred}}^{(i)} \neq y^{(i)} \\ 0 & \text{otherwise} \end{cases}$$

$$w^{(i)} = \begin{cases} 1, & \text{if } x^{(i)} \text{ is Non-PORN} \\ 10, & \text{if } x^{(i)} \text{ is PORN} \end{cases}$$

UNDERSTANDING

HUMAN - LEVEL

PERFORMANCE

- HUMAN - LEVEL ERROR
(PROXY for BAYES ERROR)

"AVOIDABLE ~~BEST~~ BIAS"

- TRAINING ERROR

"VARIANCE"

- Dev error

SUBPASSING HUMAN - LEVEL PERFORMANCE

- Team of humans

0.5%

→ PROBABLY IS
THE BAYES
ERROR

- ONE HUMAN

1%

0.5%
"AVOIDABLE BIAS"

- TRAINING ERROR

0.6%

0.2%
"VARIANCE"

- Dev error

0.8%

IMPROVING YOUR MODEL PERFORMANCE

HUMAN-LEVEL

AVOIDABLE BIAS

- TRAIN BIGGER MODEL
- TRAIN HONORABLE/BETTER OPTIMIZATION ALGORITHMS
- NN ARCHITECTURE/HYPERPARAMETERS SEARCH (RNN, CNN) *

TRAINING ERROR

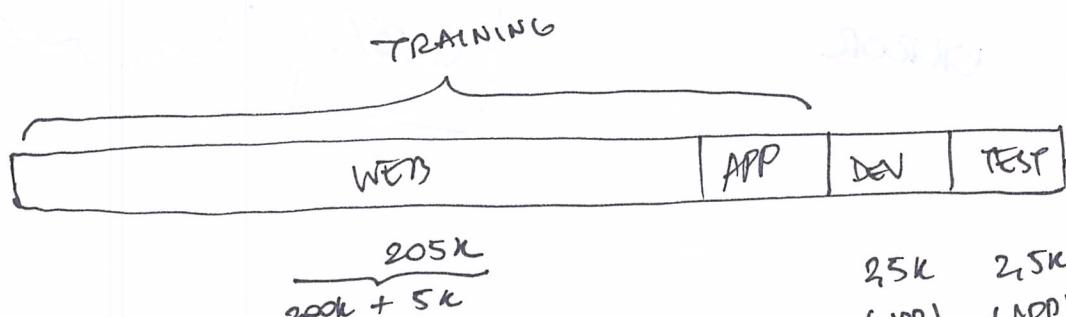
VARIANCE

- MORE DATA
- REGULARIZATION
L2, DROPOUT, DATA AUGMENTATION
- (*)

DEV ERROR

MISMATCHED TRAINING AND DEV/TEST SET

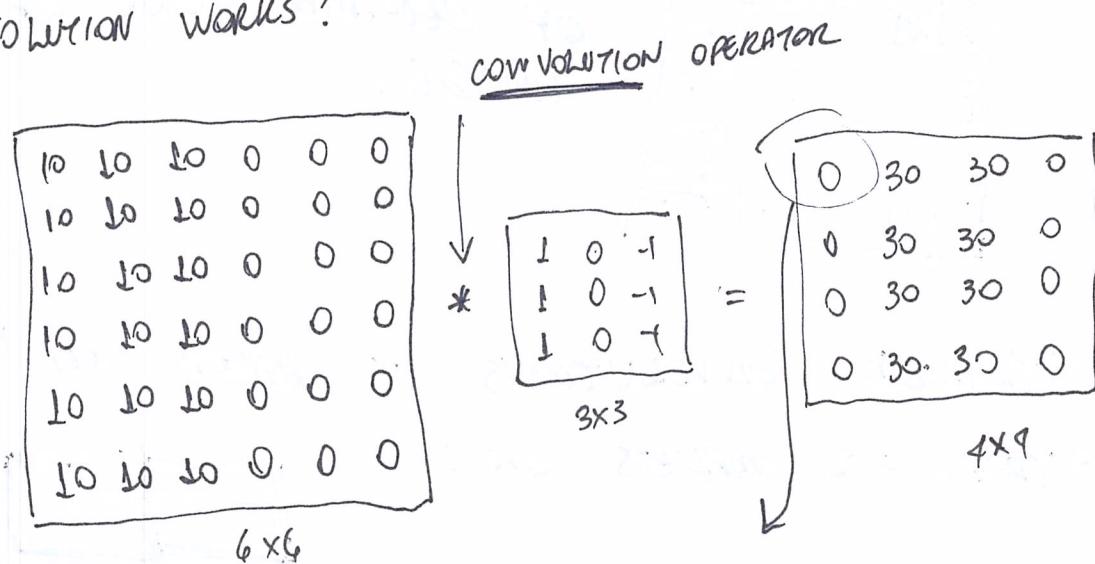
- 200K obs from WEB
- 10K obs from APP (THE ONES YOU WANT TO IDENTIFY)



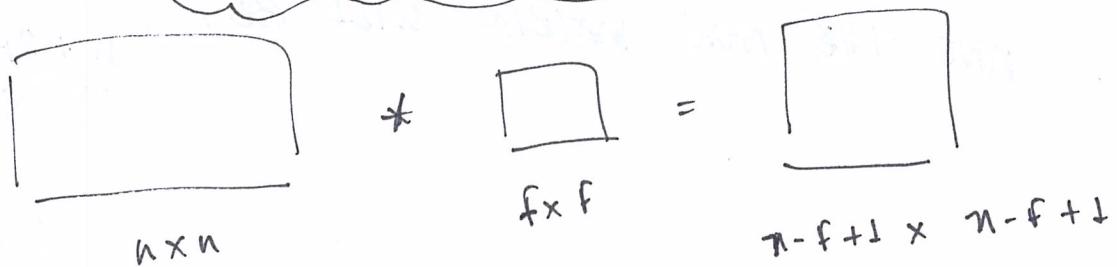
COURSE 4 - CONVOLUTIONAL NEURAL NETWORKS

Very used in computer vision

How convolution works?



$$\begin{aligned}
 & 10 \cdot 1 + 10 \cdot 1 + 10 \cdot 1 + \\
 & 10 \cdot 0 + 10 \cdot 0 + 10 \cdot 0 + \\
 & 10 \cdot (-1) + 10 \cdot (-1) + 10 \cdot (-1) = 0
 \end{aligned}$$



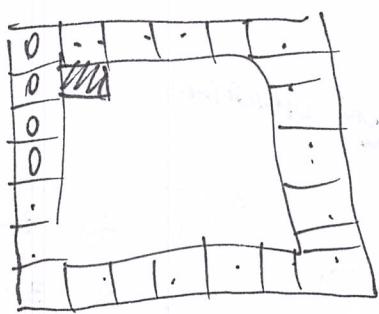
"VALID" AND "SAME" CONVOLUTIONS

- "VALID": $n \times n * f \times f \rightarrow n-f+1 \times n-f+1$
- "SAME": $\overset{\text{"PAD"} \text{ SO THE OUTPUT SIZE IS THE SAME AS}}{\text{THE INPUT SIZE.}}$

$$n + 2p - f + 1 = K \Rightarrow 2p = f - 1 \Rightarrow$$

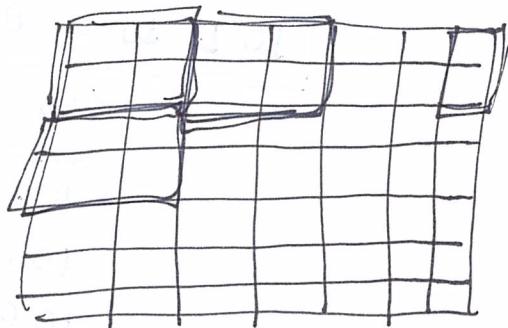
$$p = \frac{f-1}{2}$$

(1) PADDING IS TO ADD COLUMNS TO THE IMAGES



IN ORDER TO ENHANCE INFLUENCE OF CERTAIN VALUE OF THE IMAGE.

STRIDED CONVOLUTIONS IS BASED ON "JUMPS", SO STRIDE = 2 IMPLIES ON:

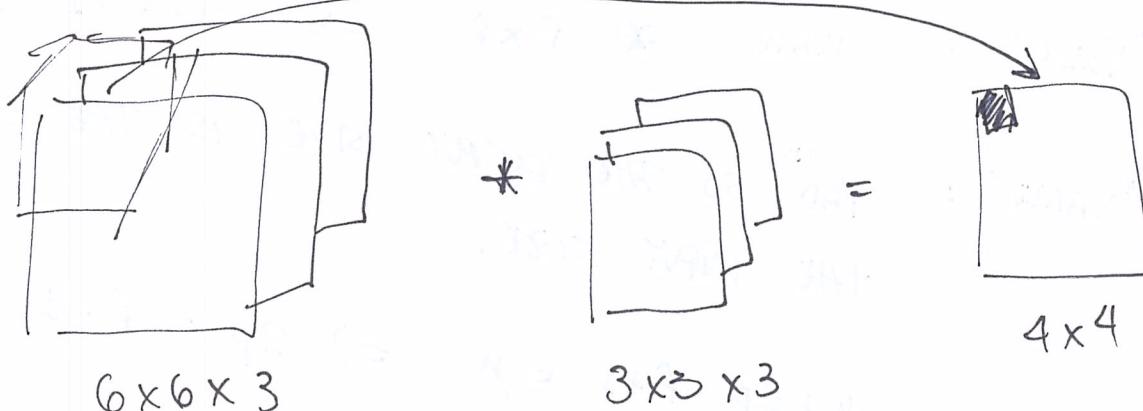


AND THE NEW MATRIX WILL BE:

$$\left(\frac{n+2p-f}{s} + 1 \right) \times$$

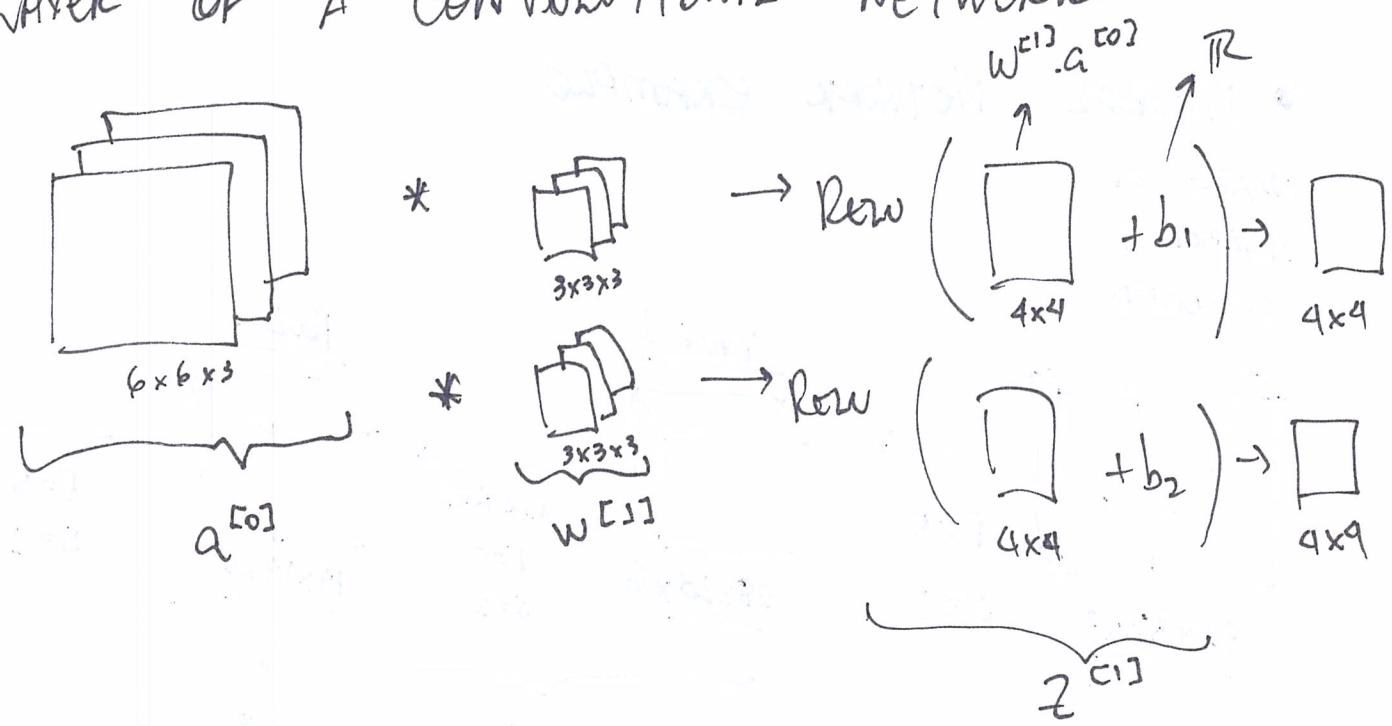
$$\left(\frac{n+2p-f}{s} + 1 \right)$$

CONVOLUTIONS OVER VOLUME:



"MULTIPLE FILTERS"

• LAYER OF A CONVOLUTIONAL NETWORK



$$\begin{cases} z^{c1} = W^{c1} a^{c0} + b^{c1} \\ a^{c1} = g(z^{c1}) \end{cases}$$

• POOLING LAYERS : MAX POOLING

1	3	2	1
2	9	1	1
1	3	2	3
5	6	1	2

4×4

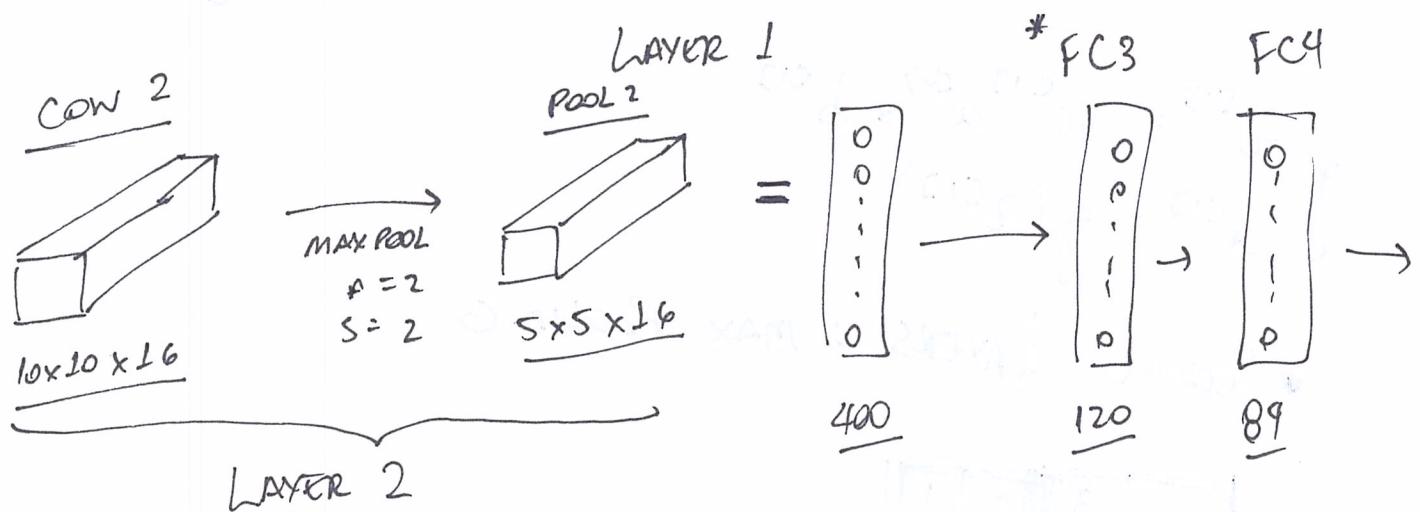
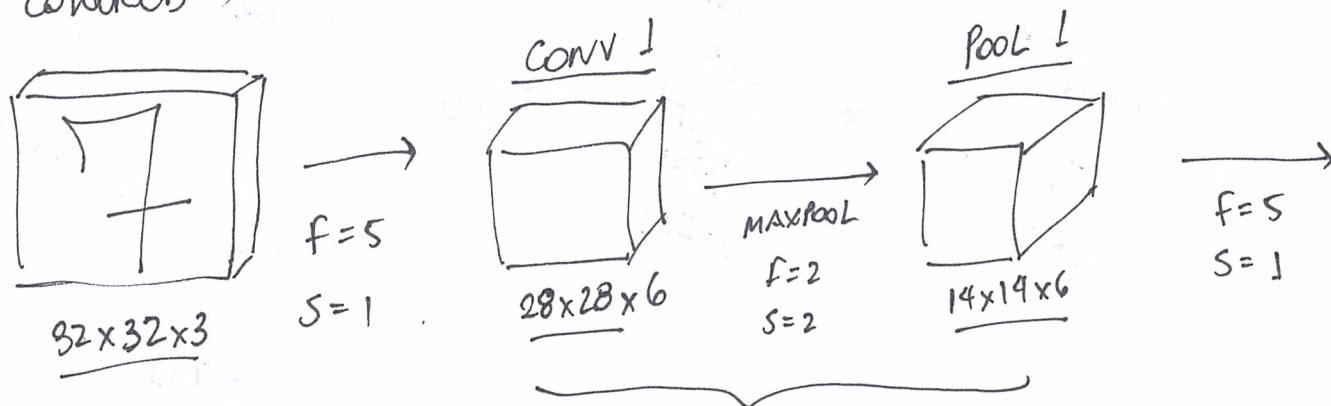
9	2
6	3

HYPER PARAMETERS:

$$\begin{aligned} f &= 2 && \text{(FILTER DIMENSION)} \\ s &= 2 && \text{(STRIDE)} \end{aligned}$$

• NEURAL NETWORK EXAMPLE

(IMAGE OF
NUMBER 7
CONVOLVED)



* FC = FULLY CONNECTED LAYER

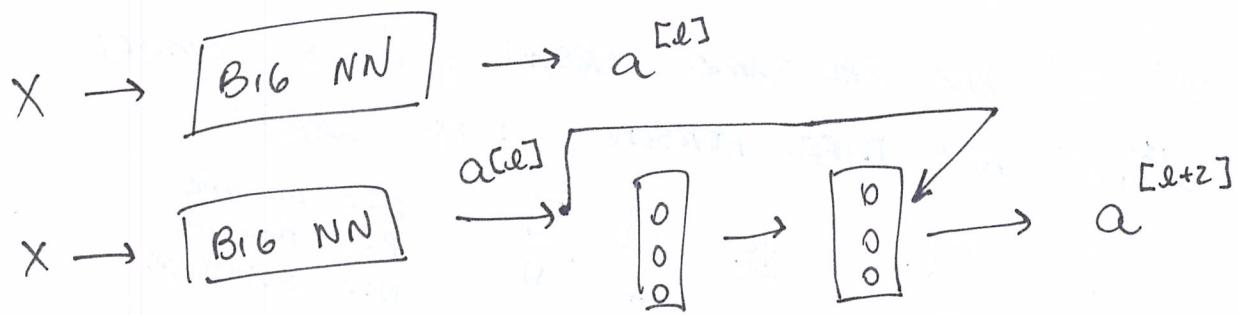
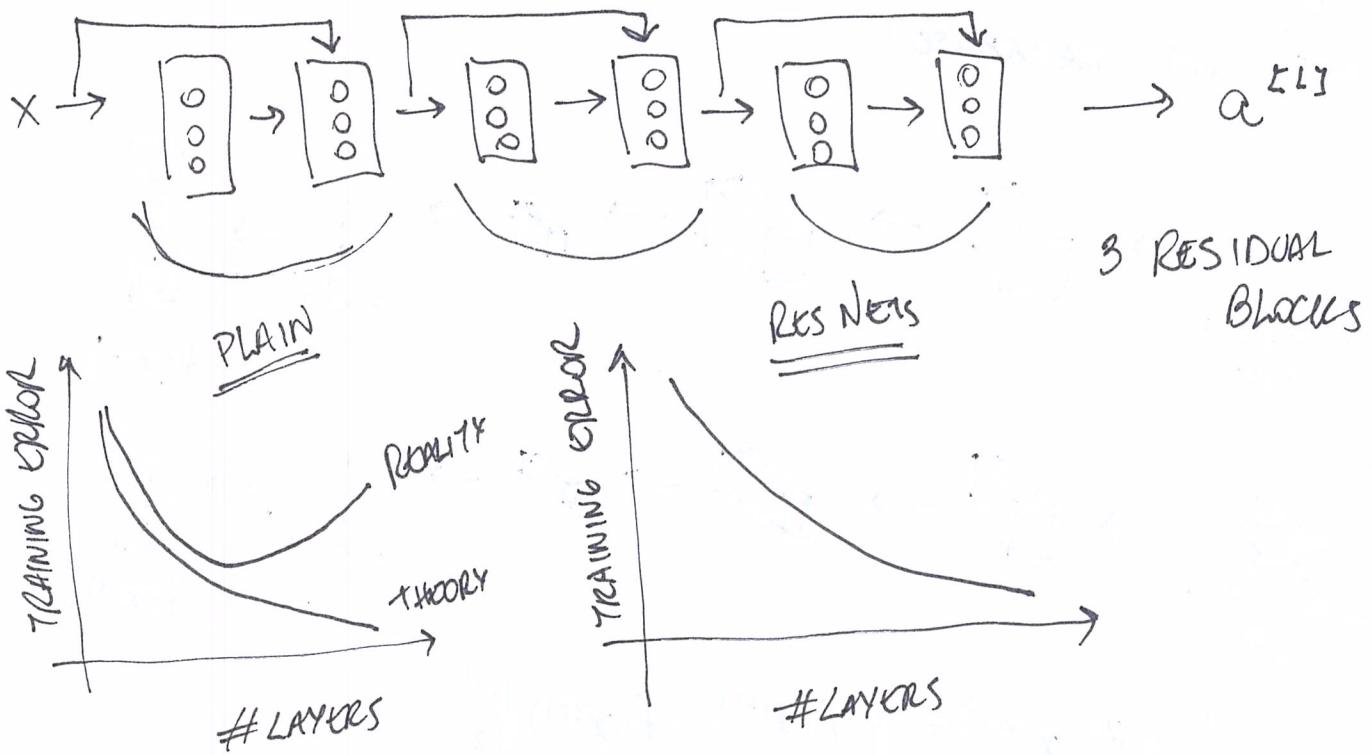
→ ○ SOFTMAX
(SO OUTPUTS = $0, 1, 2, \dots, 9$)

• WHY CONVOLUTIONS?

→ PARAMETER SHARING: A FEATURE DETECTOR THAT'S USEFUL IN ONE PART OF THE IMAGE MIGHT BE USEFUL IN ANOTHER PART OF THE IMAGE.

→ SPARSITY CONNECTIONS: IN EACH LAYER, EACH OUTPUT VALUE DEPENDS ONLY ON A SMALL NUMBER OF INPUTS.

• RESIDUAL NETWORK (ResNets)



$$\begin{aligned}
 a^{[e+2]} &= g(z^{[e+2]} + a^{[e]}) \\
 &= g(w^{[e+2]} a^{[e+1]} + b^{[e+2]} + a^{[e]}) \quad (= g(a^{[e]})) \\
 &\quad \underbrace{\qquad\qquad\qquad}_{=0} \\
 &= \cancel{a^{[e]}}
 \end{aligned}$$

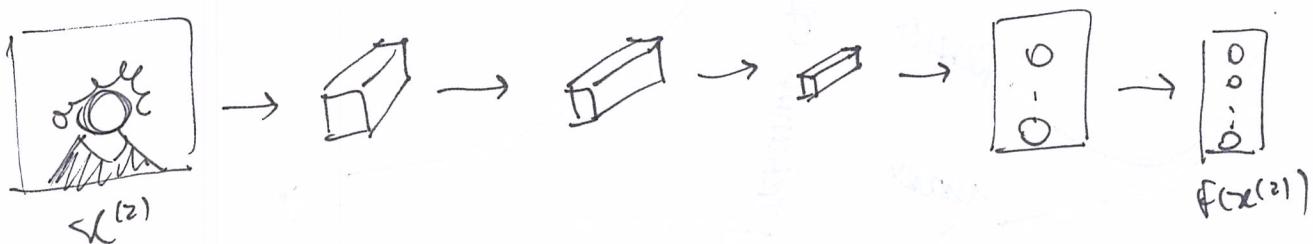
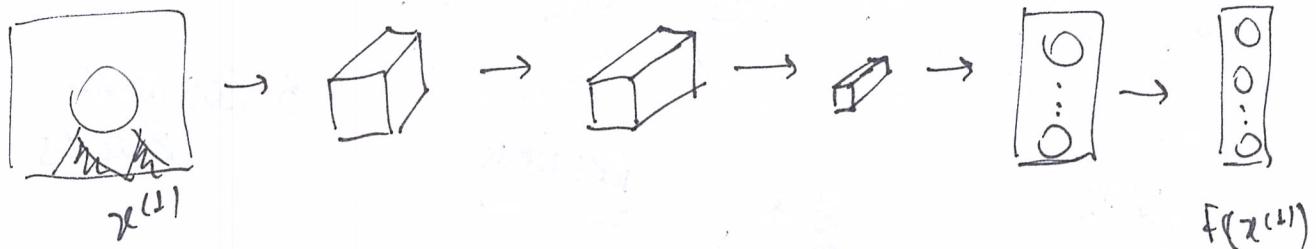
(*) RELU

RESIDUAL NETWORKS SKIP CONNECTIONS OR SHORT-CUTS TO JUMP OVER SOME LAYERS, THE MOTIVATION IS TO AVOID THE PROBLEM OF VANISHING GRADIENTS BY REUSING ACTIVATION FROM A PREVIOUS LAYER UNTIL THE LAYER NEXT TO THE CURRENT ONE HAVE LEARNED ITS WEIGHTS.

THE INTUITION ON WHY THIS WORKS IS THAT THE NEURAL NETWORK COLLAPSES INTO FEWER LAYERS INITIALLY, AND AFTER EXPANDS IN M, $\boxed{18}$

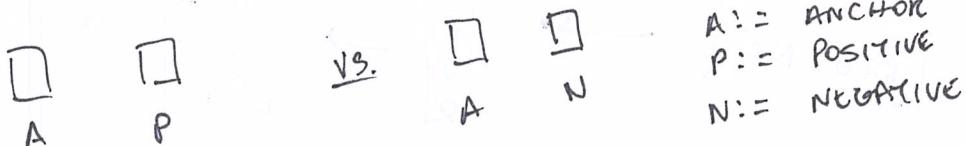
SIAMESE NET WORK

USED TO SOLVE ONE SHOT LEARNING (IDENTIFY IF PHOTO EXISTS ON DATABASE)



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

- IF $x^{(1)}, x^{(2)}$ ARE THE SAME PERSON, d IS SMALL
- IF $x^{(1)}, x^{(2)}$ ARE DIFF. PERSON, d IS LARGE



Loss function : $\mathcal{L}(A, P, N) = \max \left(\underbrace{\|f(A) - f(P)\|^2}_{d(A, P)} - \underbrace{\|f(A) - f(N)\|^2}_{d(A, N)} + \alpha, 0 \right)$

"TRIPLET LOSS"

$$\underbrace{\|f(A) - f(P)\|^2}_{d(A, P)} + \alpha \leq \underbrace{\|f(A) - f(N)\|^2}_{d(A, N)}$$

COURSE 5 - SEQUENCE MODELS

RECURRENT NEURAL NETWORKS

WHY RNNs?

SPEECH RECOGNITION

HUNNY HUNNY \Rightarrow "YOU ARE THE SUNSHINE OF MY LIV"

SENTIMENT CLASSIFICATIONS "HATED THIS MOVIE" \Rightarrow ★★☆☆☆

DNA SEQ. ANALYSIS

A G C C C F G T G A C C \Rightarrow PATTERN REC.

MACHINE TRANSLATION

"How ARE You?" \Rightarrow "COMO VOCÊ ESTA?"

NAME RECOGNITION

"HARRY POTTER INVENTED
A SPELL" \Rightarrow HARRY POTTER

E.G.: RECOGNIZING NAMES ON PHRASES

HARRY POTTER INVENTED A SPELL

NOTATION

$x: x^{<1>} \quad x^{<2>} \quad x^{<3>} \dots x^{<t>} \dots x^{<5>}$

$y: \downarrow \quad \downarrow \quad \quad \quad 0 \quad \dots \quad \dots \quad 0$
 $y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad \dots \quad \dots \quad y^{<5>}$

$x^{(i)<t>} :=$ WORD IN POSITION t OF SAMPLE i

$y^{(i)<t>} :=$ ANSWER IN POSITION t OF SAMPLE i

$T_x^{(i)} :=$ NUMBER OF WORDS OF SAMPLE i

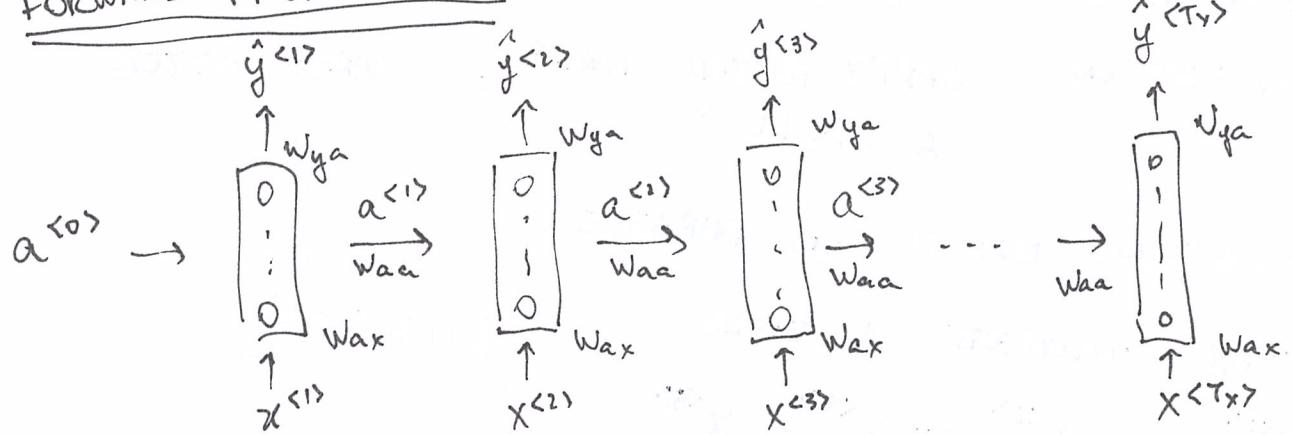
$T_y^{(i)} :=$ NUMBER OF ANSWERS OF SAMPLE i

Some cases? $T_x = T_y$ (like this example above). But for sentiment analysis, for example, $T_x = \#$ of words in TEXT, $T_y = \#$ answers (0 or 1) = good or bad.

EXPLANATORY FEATURES TO INPUT:

VOCABULARY =	A	$x^{<1>} =$ "HARRY"	$x^{<2>} =$ "POTTER"
	$\begin{bmatrix} A \\ \vdots \\ HARRY \\ \vdots \\ POTTER \\ \vdots \\ ZULU \end{bmatrix}$	$\begin{bmatrix} 1 \\ \vdots \\ 3675 \\ \vdots \\ 6800 \\ \vdots \\ 10000 \end{bmatrix}$	$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} \rightarrow 3675$
			$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} \rightarrow 6800$

FORWARD PROPAGATION



$$a^{<0>} = 0$$

$$a^{<1>} = g_1 (W_{aa} a^{<0>} + W_{ax} x^{<1>} + b_a)$$

$$\hat{y}^{<1>} = g_2 (W_{ya} a^{<1>} + b_y)$$

$g_1 \rightarrow \text{TANH/RELU}$

$g_2 \rightarrow \text{SIGMOID}$

$$a^{<t>} = g (W_{aa} a^{<t-1>} + W_{ax} x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g (W_{ya} a^{<t>} + b_y)$$

for SIMPLIFICATION:

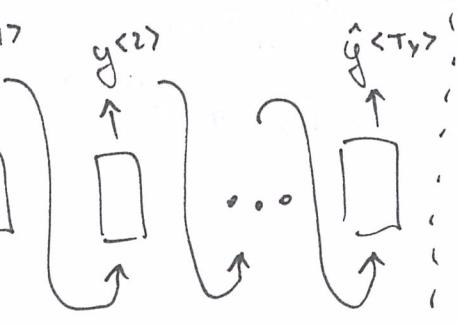
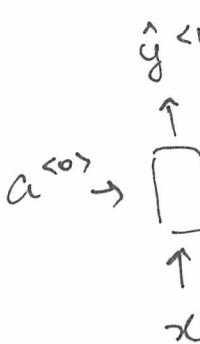
$$W_a = [W_{aa} : W_{ax}]$$

for BACK PROPAGATION:

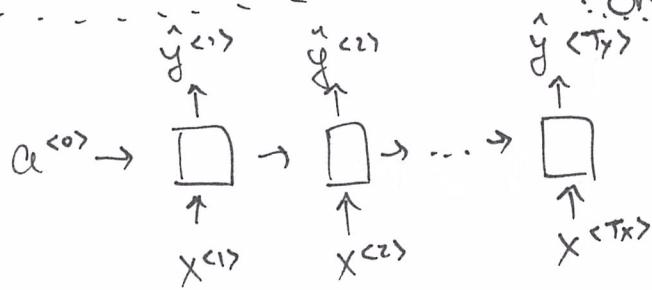
$$\delta^{<t>} (\hat{y}^{<t>}, y^{<t>}) = -y^{<t>} \cdot \log \hat{y}^{<t>} - (1 - \hat{y}^{<t>}) \log (1 - \hat{y}^{<t>})$$

$$\therefore \mathcal{L}(\hat{y}, y) = \sum_{t=1}^{T_y} \mathcal{L}^{<t>} (\hat{y}^{<t>}, y^{<t>})$$

RNN TYPES

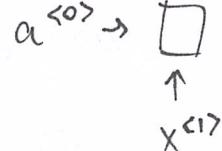


ONE TO ONE



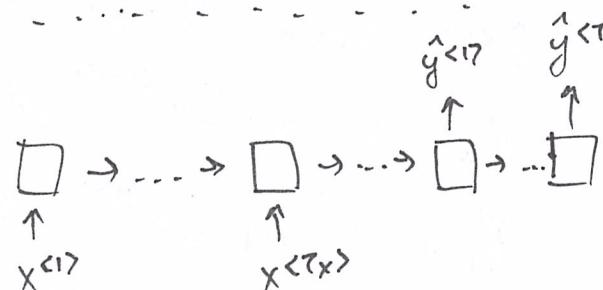
MANY TO MANY ($T_x = T_y$)

ONE TO MANY



MANY TO MANY

MANY TO ONE



GRU & LSTM

MOIVATION: "THE DOGS, WHO ATE A HOT YESTERDAY, WERE TIRED TODAY"

BUILD A MODEL THAT RELATED WORDS DISTANT FROM EACH OTHER.

GRU: GATED RECURRENT UNIT

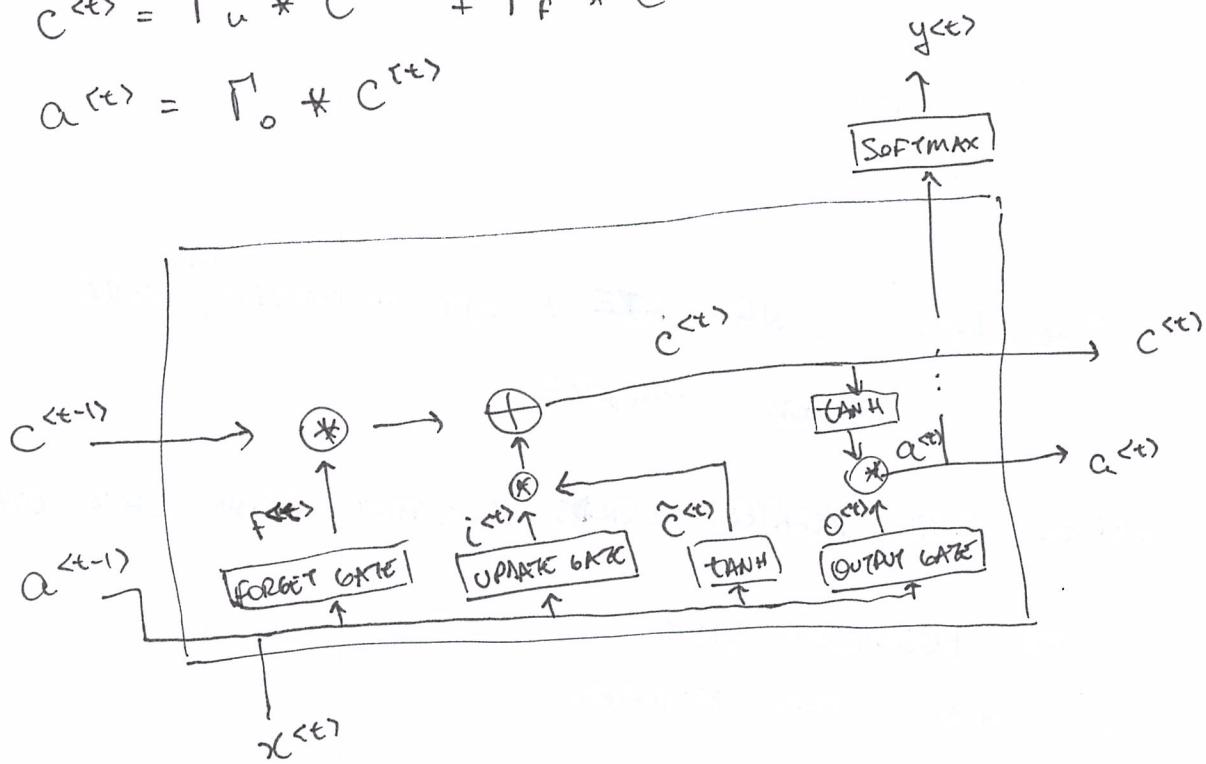
LSTM: LONG SHORT TERM MEMORY

GRU

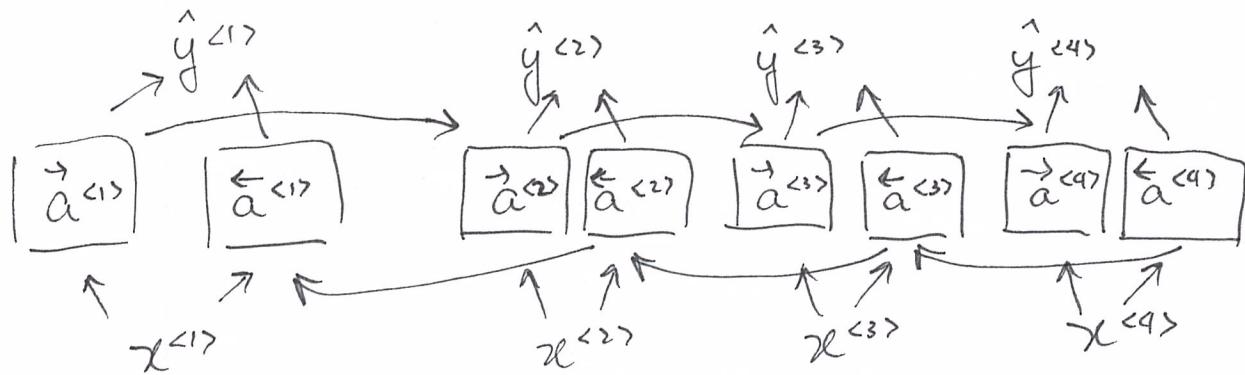
$$\begin{aligned}
 \tilde{c}^{(t)} &= \tanh(W_c[c^{(t-1)}, x^{(t)}] + b_c) \\
 \Gamma_u &= \sigma(W_u[c^{(t-1)}, x^{(t)}] + b_u) \\
 \Gamma_r &= \sigma(W_r[c^{(t-1)}, x^{(t)}] + b_r) \\
 c^{(t)} &= \Gamma_u * \tilde{c}^{(t)} + (1 - \Gamma_u) * c^{(t-1)} \\
 a^{(t)} &= c^{(t)}
 \end{aligned}
 \quad \left. \begin{array}{l} \\ \\ \\ \end{array} \right\} 2 \text{ GATES}$$

LSTM

$$\begin{aligned}
 \tilde{c}^{(t)} &= \tanh(W_c[a^{(t-1)}, x^{(t)}] + b_c) \\
 \Gamma_u &= \sigma(W_u[a^{(t-1)}, x^{(t)}] + b_u) \\
 \Gamma_f &= \sigma(W_f[a^{(t-1)}, x^{(t)}] + b_f) \\
 \Gamma_o &= \sigma(W_o[a^{(t-1)}, x^{(t)}] + b_o) \\
 c^{(t)} &= \Gamma_u * \tilde{c}^{(t)} + \Gamma_f * c^{(t-1)} \\
 a^{(t)} &= \Gamma_o * c^{(t)}
 \end{aligned}
 \quad \left. \begin{array}{l} \\ \\ \\ \end{array} \right\} 3 \text{ GATES}$$



Bi-DIRECTIONAL RNN



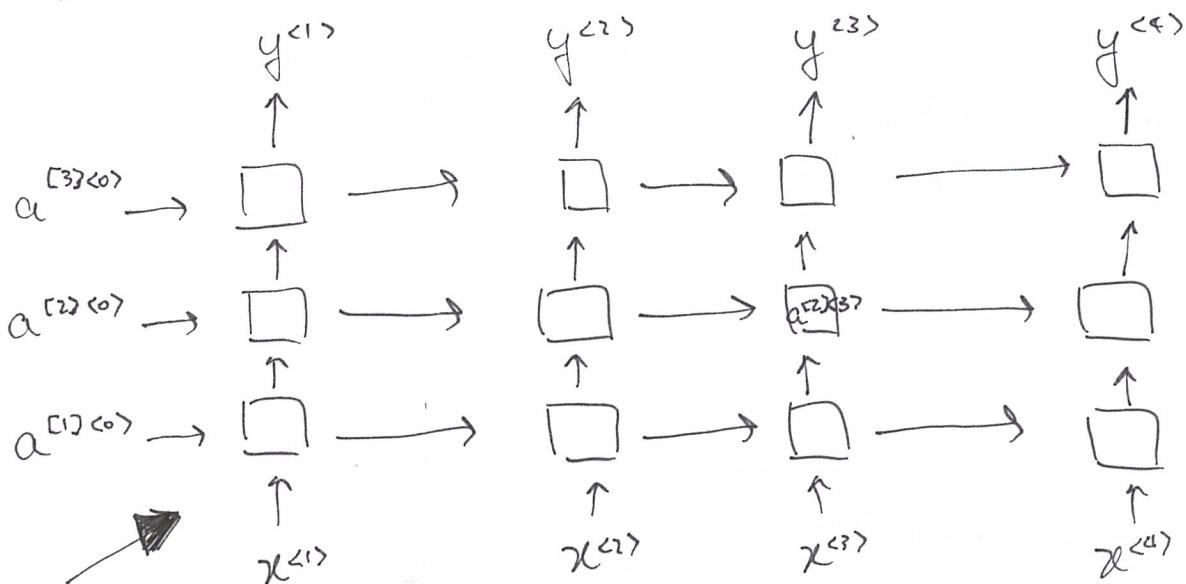
$$\therefore \hat{y}^{(t)} = g(W_y[\vec{a}^{(t)}, \underline{a}^{(t)}] + b_y)$$

THESE BLOCKS CAN BE GRU AND LSTM

ADVANTAGE: CAN PREDICT ANY TYPE OF TEXT (IN MIDDLE, BEGINNING OR END OF TEXT)

DISADVANTAGE: YOU NEED THE ENTIRE SEQUENCE

DEEP RNNs



$$a^{[2]<3>} = g(W_a^{[2]} [a^{[2]<2>}, a^{[1]<3>}] + b_a^{[2]})$$

RNN, GRU, LSTM, BRNN

WORD REPRESENTATION

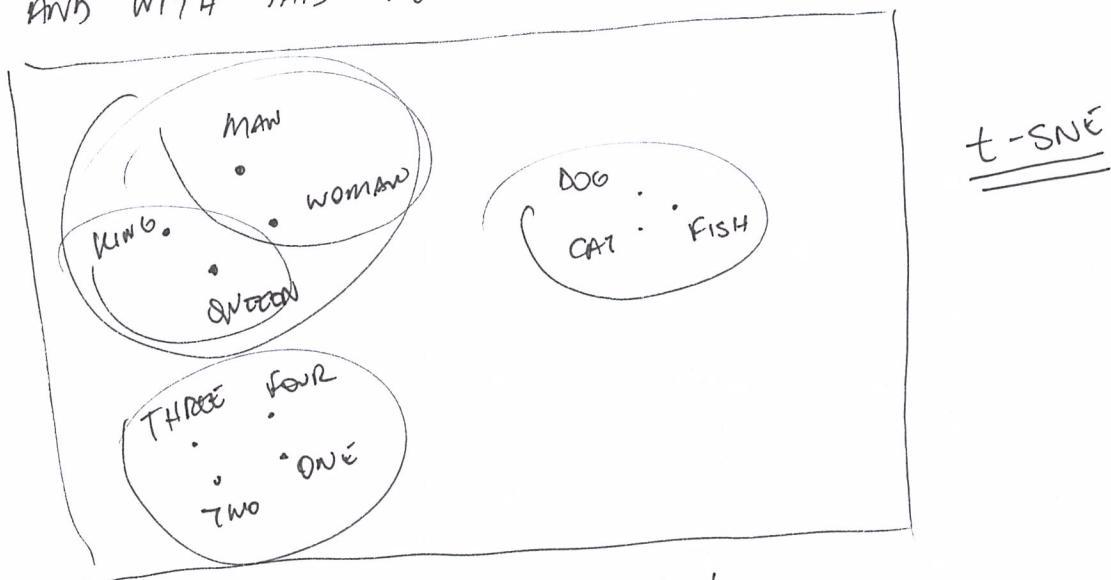
$V = [A, AARON, \dots, ZULU, \text{UNKNOWN}]$, $|V| = 10,000$

FEATURE REPRESENTATION

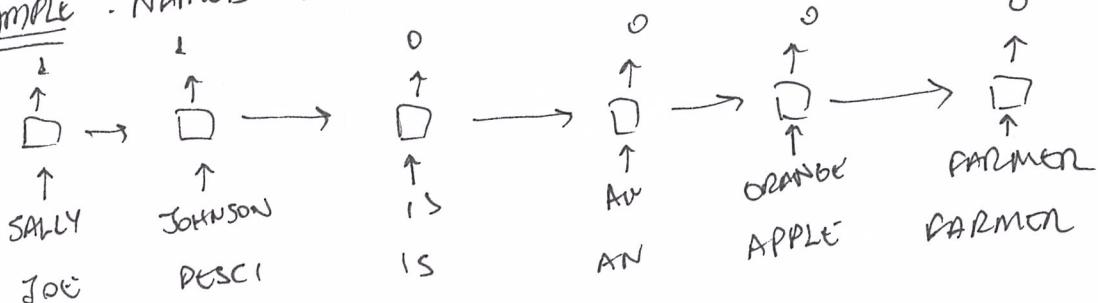
	MAN	WOMAN	KING	QUEEN	APPLE	ORANGE
GENDER	-1	1	-0.95	0.97	0.00	0.01
ROYAL	0.01	0.02	0.93	0.95	-0.01	0.00
AGE	0.03	0.02	0.7	0.69	0.03	-0.02
FOOD	0.04	0.01	0.02	0.01	0.95	0.97
:	:	:	:	:	:	:
SIZE	:	:				
COS1						

$\hookrightarrow e_{531}$ $\hookrightarrow e_{9853}$

AND WITH THIS REPRESENTATION CAN CREATE CLUSTERS:



EXAMPLE : NAMED ENTITY RECOGNITION



IT IS POSSIBLE FOR THE MODEL TO IDENTIFY CLOSURE FROM OTHER EXAMPLES (ANALOGIES). IN THAT MANNER:

$$\begin{aligned} \ell_{\text{MAN}} - \ell_{\text{WOMAN}} &\approx \left[\begin{array}{c} -2 \\ 0 \\ 0 \\ 0 \end{array} \right] \\ \ell_{\text{KING}} - \ell_{\text{QUEEN}} &\approx \left[\begin{array}{c} -2 \\ 0 \\ 0 \\ 0 \end{array} \right] \end{aligned}$$

$\ell_{\text{MAN}} - \ell_{\text{WOMAN}} \approx \ell_{\text{KING}} - \ell_{\text{QUEEN}}$

ℓ_{QUEEN}

$\vec{\theta}_1$ AND $\vec{\theta}_2$ ARE VERY SIMILAR

TO FIND WORD w : $\underset{w}{\text{ARG MAX}} \text{sim}(\ell_w, \ell_{\text{KING}} - \ell_{\text{MAN}} + \ell_{\text{WOMAN}})$

FOR FUNCTION $\text{sim}(\cdot)$: COSINE FUNCTION OF SIMILARITY

$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$

ENCODING MATRIX

	A	ARRON	...		ℓ_{6257}		(k) NUMBER OF FEATURES
	300				ORANGE ■■■ ■■■■ ■■■■■	... ZULU UNKNOWN	$E_{(300, 10000)}$

$$\begin{aligned} O_{6257} &= \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \rightarrow 6257 \quad \uparrow 30000 \quad ; E_{(300, 10000)} \cdot O_{6257, (10000, 1)} = \begin{bmatrix} ■■■ \\ ■■■■ \\ ■■■■■ \end{bmatrix} = \ell_{6257} \end{aligned}$$

THAT MEANS : $E \cdot O_f = e_f$:= EMBEDDING FOR WORD f .

WORD 2 VEC

"I WANT A GLASS OF ORANGE JUICE TO GO ALONG WITH MY CUPCAKE"

→ RANDOMLY PICK CONTEXT

TARGET

→ BY CHANCE CHOOSE $+N$ OR $-N$ WORD POSITION TO BE THE TARGET. ($N \in [1, 10]$)

ORANGE

JUICE

GLASS

ORANGE

MY

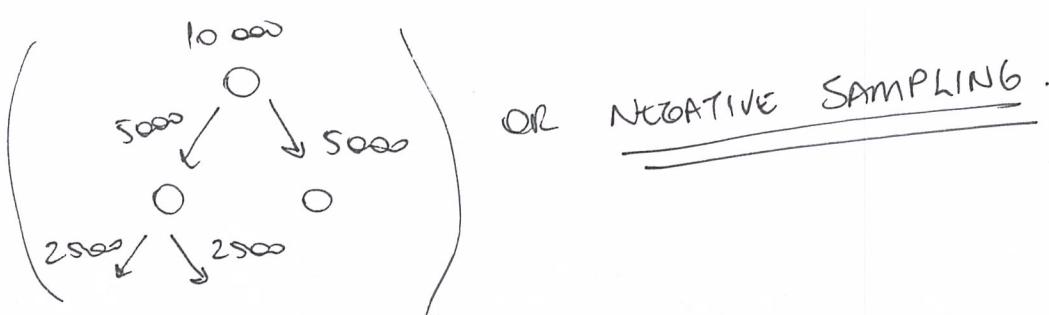
$x_c \rightarrow y$
CONTEXT (c) → TARGET (t)

$O_c \rightarrow E \rightarrow e_c \rightarrow \text{SOFTMAX} \rightarrow \hat{y}$
 $E \cdot O_c = e_c$

SOFTMAX : $P(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}}$, θ_t = PARAMETER ASSOCIATION WITH OUTPUT t

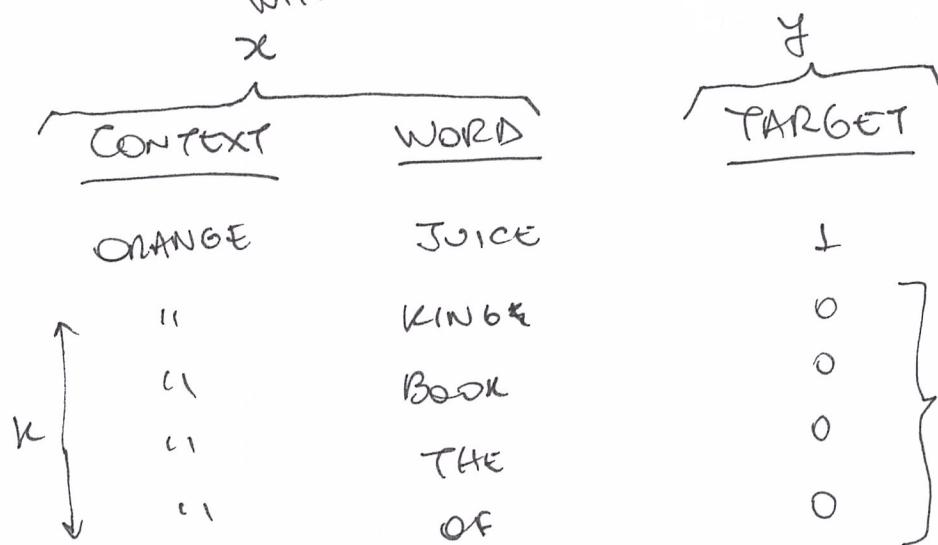
$$\mathcal{L}(\hat{y}, y) = - \sum_{i=1}^{10000} y_i \cdot \log \hat{y}_i, \quad y = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

TO TRAIN THE MODEL IT CAN BE USED HIERARCHICAL SOFTMAX

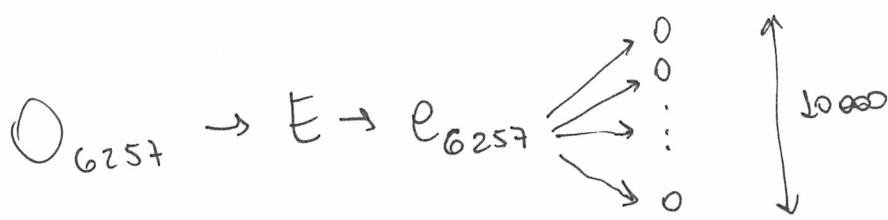


NEGATIVE SAMPLING

"I WANT A GLASS OF ORANGE JUICE TO GO ALONG WITH MY CEREAL"



THESE ARE ALL NEG. SAMPLES, THAT IS, NON-TRUE EXAMPLES



BUT ONLY CLASSIFY K+L AND THAT'S WHY IS FAST THAN HIERARCHICAL SOFTWARE

SELECTING NEG. EXAMPLES

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10000} f(w_j)^{3/4}}, \text{ where } f(w_i) \text{ is the frequency of word } i.$$

OBSERVATION

- THERE ARE TWO MAIN WORD 2 VEC MODELS :
- CONTINUOUS BAG OF WORDS (CBOW)
- SKIP-GRAM

IN THE CBOW MODEL, WE PREDICT A WORD GIVEN A CONTEXT (A CONTEXT CAN BE SOMETHING LIKE A SENTENCE).
IN THE SKIP-GRAM IS THE OPPOSITE: PREDICT THE CONTEXT GIVEN AN INPUT WORD.