

SUMMARY

E

T

P

- LEARNING = EXPERIENCE (DATA) + TASK + PERFORMANCE
- INTRODUCTION TO PROBABILITY (SKIPPED)
- UNSUPERVISED LEARNING:

GIVEN DATA x_1, \dots, x_n (NO LABELS y_i) \rightarrow OUTPUT: PATTERN BEHIND DATA
 \rightarrow ex. CLUSTERING

- CLUSTERING:

PROCESS OF GROUPING OBJECTS INTO CLASSES OF SIMILAR OBJECTS

- THE CONCEPT OF "SIMILARITY" WILL BE DEFINED BY DISTANCE

- 2 ALGORITHMS:

- K-MEANS ALGORITHM (\rightsquigarrow already covered) : HARD CLUSTERING METHOD
 \hookrightarrow 1 POINT \rightarrow 1 CLUSTER

- K-MEANS:

GIVEN A SET OF OBSERVATIONS $(x_1, \dots, x_n) / x_i \in \mathbb{R}^d$,
, PARTITION THE n OBSERVATIONS INTO K SETS $S = \{S_1, \dots, S_k\}$
, SUCH THAT THE SETS MINIMIZE THE WITHIN-CLUSTER SUM OF SQUARE:
$$\min_{S} \sum_{i=1}^k \sum_{x_j \in S_i} \|x_j - \mu_i\|^2 / \mu_i \cdot \text{mean of points in } S_i$$

• GAUSSIAN MIXTURE MODEL : <https://towardsdatascience.com/gaussian-mixture-models-explained-6986aaf5a95>

- NOT HARD CLUSTERING : GMM ASSIGN TO EACH POINT A PROBABILITY THAT THE POINT BELONGS TO A CLUSTER

$$\rightarrow \text{A POINT, } K \text{ CLUSTER} \rightarrow (\tilde{\pi}_1, \dots, \tilde{\pi}_K) / \sum_{i=1}^K \tilde{\pi}_i = 1$$

$\downarrow p \text{ OF BEARING TO CLUSTER } i$

- K "CLUSTERS" $\rightarrow K$ GAUSSIAN COMPONENTS $\sim N(\bar{\mu}_i, \Sigma_i)$

$$\rightarrow P(x) = \sum_{i=1}^K p(x | r=i) P(r=i)$$

$\underbrace{p(x | r=i)}_{\substack{\text{OBSERVED} \\ \text{DATA}}}$ $\underbrace{P(r=i)}_{\substack{\text{MIXTURE PROPORTION}}}$

$\hookrightarrow P \text{ THAT OUTPUT } \in \text{GAUSSIAN}$

- ASSUMING THAT $p(x | r=i) = N(\mu_i, \sigma^2 I)$:

$$\rightarrow \text{THE PARAMETERS TO FIND ARE } \underline{\mu_1, \dots, \mu_K}, \underline{\sigma^2}, \underline{\tilde{\pi}_1, \dots, \tilde{\pi}_K}$$

$\underbrace{\mu_1, \dots, \mu_K}_{\substack{\text{MEAN OF THE} \\ \text{K GAUSSIANS}}}, \underbrace{\sigma^2}_{\substack{\text{VAR.} \\ \text{OF} \\ \text{K GAUSS.}}}, \underbrace{\tilde{\pi}_1, \dots, \tilde{\pi}_K}_{P[r=i]}$

\rightarrow ALL THIS PARAMETERS CAN BE FOUND USING MLE:

$$\Theta = [\mu_1, \dots, \mu_K, \sigma^2, \tilde{\pi}_1, \dots, \tilde{\pi}_K]$$

$\hookrightarrow \text{IN GENERAL: } \sigma^2 \mapsto \Sigma_1, \dots, \Sigma_K$

\rightarrow EXPECTATION-MAXIMIZATION ALGORITHM:

- E-STEP: ASSIGN CLUSTERS TO EACH POINT
/ IN A SOFT WAY

- M-STEP: UPDATE CLUSTERS IN K-MEANS

E-step

Compute "expected" classes of all datapoints for each class

$$P(y_j = i | x_j, \theta^{t-1}) = \frac{\exp(-\frac{1}{2\sigma^2} \|x_j - \mu_i^{t-1}\|^2) \pi_i}{\sum_{i=1}^K \exp(-\frac{1}{2\sigma^2} \|x_j - \mu_i^{t-1}\|^2) \pi_i}$$

In K-means "E-step" we do hard assignment. EM does soft assignment

M-step

Compute Max. like μ given our data's class membership distributions (weights)

$$\mu_i^t = \sum_{j=1}^n w_j x_j \quad \text{where } w_j = \frac{P(y_j = i | x_j, \theta^{t-1})}{\sum_{l=1}^n P(y_l = i | x_l, \theta^{t-1})}$$

KNN, PERCEPTRONS:

KNN:

- A NEW TEST POINT GETS THE SAME LABEL AS THE K MOST CLOSE TRAINING POINTS, BY MAJORITY VOTING

S: SENSITIVE TO OUTLIERS

- KNN CAN BE USED FOR REGRESSION TOO:

GIVEN THE TEST POINT X:

- COMPUTE DISTANCE $D(x, x_i)$ TO EVERY TRAINING EXAMPLE, SELECT K CLOSEST AND THEIR y_i
- OUTPUT: MEAN OF $y_{i_1}, \dots, y_{i_K} \rightarrow \hat{y} = \frac{1}{K} \sum_{i=1}^K y_i$

- VALUE OF K: USE A VALIDATION SET AND VARY K
→ SELECT K THAT YIELDS, IN GENERAL, THE BEST PERFORMANCE

DISTANCE METRICS:

DEFINED FROM THE MINKOWSKI DISTANCE, $L_p(x, r) = \left(\sum_{i=1}^d |x_i - r_i|^p \right)^{1/p}$

$$\cdot p=1: L_1 \text{ norm (MANHATTAN DISTANCE)} = L_1(x, r) = \sum_{i=1}^d |x_i - r_i|$$

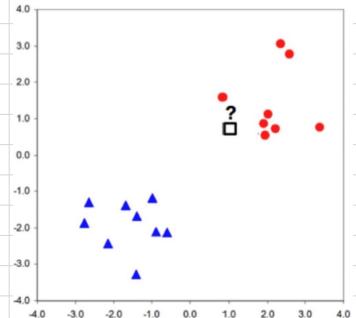
$$\cdot p=2: L_2 \text{ norm (EUCLIDEAN DISTANCE)} = L_2(x, r) = \sqrt{\sum_{i=1}^d (x_i - r_i)^2}$$

$$\cdot p=\infty: L_\infty \text{ norm (CHEBYSHEV DISTANCE)} = L_\infty(x, r) = \max_i |x_i - r_i|$$

$$\cdot p=0: L_H \text{ (HAMMING DISTANCE)} = L_H(x, r) = \sum_{i=1}^d 1_{x_i \neq r_i}$$

KNN TRICKS:

- PARTIAL DISTANCE CALCULATION: D COMPUTED ON SUBSET R COL OF THE FULL OF OBSERVATIONS
- PRE-STRUCTURING/SEARCHING FOR PROTOTYPES: EXPLORE SIMILARITIES IN SAMPLES → DATA AS SEARCH TREES
- EDITING: USELESS POINTS ARE PRUNED → COMPLEXITY REDUCTION
- FEATURE NORMALIZATION



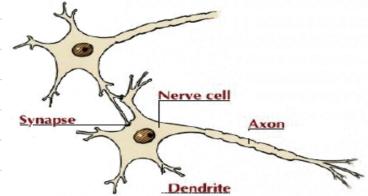
• PERCEPTRON :

• FROM BIOLOGY :

- GOOD BEHAVIOR \rightarrow REWARD, BAD BEHAVIOR \rightarrow PUNISHED (OR NO REWARD)

• NEURONS :

- SOMA (CPU) : CELL BODY, COMBINES SIGNALS
- DENDRITES (INPUT BUS) : COMBINES INPUTS FROM OTHER CELLS
- SYNAPSE (INTERFACE)
- AXON (CABLE) : TRANSPORTS ACTIVATION SIGNAL TO NEURONS AT DIFFERENT LOCATION



-> IMPLEMENTATION ON MACHINE -

• NEURON :

$$f(x) = \sum_i w_i \cdot x_i = \langle \bar{w}, \bar{x} \rangle$$

SCALAR PRODUCT

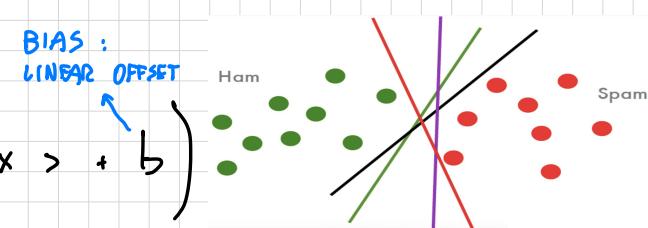
• PERCEPTRON :

NON LINEAR
DECISION FUNCTION
 \uparrow (e.g. SIGMOID, RELU)

BIAIS :
LINEAR OFFSET

$$f(x) = G \cdot \left(\langle \bar{w}, \bar{x} \rangle + b \right)$$

\uparrow IF BINARY OUTPUT
 \uparrow LINEAR SEPARATING HYPERPLAN



SEE ALSO
NOTION

$\rightarrow \bar{w}$ AND b ARE THE PARAMETERS TO LEARN

• LEARNING ALGORITHM : \rightsquigarrow FOR BINARY CLASSIFIER \rightarrow OUT = $\{-1, 1\}$

1. SET $\bar{w} = 0$, $b = 0$

TRAINING
DATASET

IF $\langle \bar{w}, \bar{x}_i \rangle + b \leq 0 \rightarrow x_i$ CLASSIFIED WRONG

2. REPEAT \forall POINT $x_i \in X$:

TRUE
LABEL

PREDICTED LABEL

\uparrow $\langle \bar{w}, \bar{x}_i \rangle + b \leq 0 \Leftrightarrow \begin{cases} y_i = -1, x_i = 1 \\ y_i = 1, x_i = -1 \end{cases}$ OR

\Leftrightarrow CLASSIFIED
WRONG \rightarrow NOT y_i , ONLY $i \in I$

IF $y_i [\langle \bar{w}, \bar{x}_i \rangle + b] \leq 0$:
 w, b UPDATE

STEP 0: $w = 0$
STEP 1: $w = 0 + y_1 x_1$
STEP 2: $w = 0 + y_1 x_1 + y_2 x_2$
 \vdots
STEP m : $w = 0 + \sum_{i=1}^m y_i x_i$

$\rightarrow w \leftarrow w + y_i x_i$ $\rightsquigarrow (ANALOG AS FOR W)$

UNTIL ALL POINTS ARE CLASSIFIED

- CONVERGENCE THEOREM :

IF \exists 2 VALUES (w^*, b^*) / $\forall_i [\langle x_i, w^* \rangle + b^*] \geq g, \forall_i$

\rightarrow PERCEPTRON CONVERGES TO A LINEAR SEPARATOR

AFTER A FINITE NUMBER OF STEPS BOUNDED BY :

$$((b^*)^2 + 1) (r^2 + 1) \tilde{s}^2, \|x_i\| \leq r$$

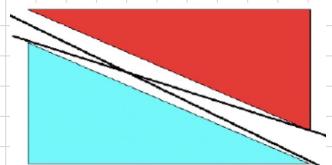
(PROOF ON NOTION, OR PAG. 11-12)

- IT IS NEEDED TO ONLY STORE ERRORS

- CONCEPTS AND VERSION SPACE :

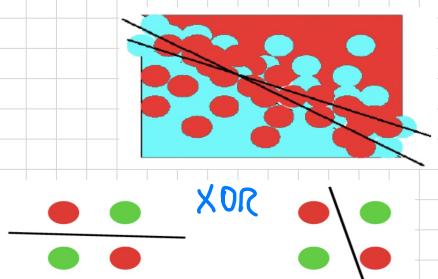
- REALIZABLE CONCEPT :

- \exists f THAT CAN SEPARATE DATA, f INCLUDED IN CONCEPT SPACE
- \rightarrow FOR PERCEPTRON, DATA ARE LINEARLY SEPARABLE



- UNREALIZABLE CONCEPT :

- DATA NOT SEPARABLE
- $\rightarrow \exists$ SUITABLE f : DATA NON DISTINGUISHABLE
- ex. XOR \rightarrow NOT LINEARLY SEPARABLE



NON-LINEARITY AND PRE-PROCESSING:

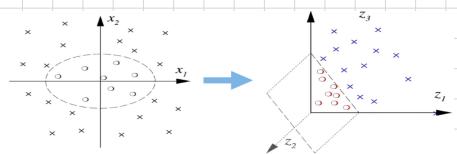
- DATA ARE MAPPED INTO FEATURE-SPACE: $x \mapsto \phi(x)$

\rightarrow SOLUTIONS BASED ON $\phi(x_i)$

ex. OR RECOGNITION

$\underline{\text{PIXELS} \mapsto \phi(\text{pixels})}$

\rightarrow RECOGNIZE CHAR BY FEATURES



	1	2	3	4	5	6	7	8	9	0
Loops	0	0	0	1	0	1	0	2	1	1
3 Joints	0	0	0	0	0	1	0	0	1	0
4 Joints	0	0	0	1	0	0	0	0	1	0
Angles	0	1	1	1	1	0	1	0	0	0
Ink	1	2	2	2	2	2	1	3	2	2

- PERCEPTRON ALGORITHM ON FEATURES:

$$1. \text{ SET } W = 0, b = 0$$

2. REPEAT \forall POINT $x_i \in X$:

$$\text{IF } V_i (W \cdot \phi(x_i) + b) \leq 0:$$

W, b UPDATE

$$\begin{aligned} W &\leftarrow W + V_i \phi(x_i) \\ b &\leftarrow b + V_i \end{aligned} \rightarrow W = \sum_{i \in I} V_i \phi(x_i)$$

UNTIL ALL POINTS ARE CLASSIFIED $\rightarrow V_i (W \cdot \phi(x_i) + b) > 0, \forall i$

$$\rightarrow \phi(x) = \sum_{i \in I} V_i \langle \phi(x_i), \phi(x) \rangle + b$$

- GOOD FEATURE ENGINEERING NEEDS DOMAIN EXPERTS, EXPENSIVE TO COMPUTE:

SOL:

COMPUTE MANY FEATURES AND HOPE THAT \exists GOOD ONES

, TRY TO SELECT FEATURES EFFICIENTLY: NO CORRELATION, NOT TOO MANY

- SOLVING XOR:

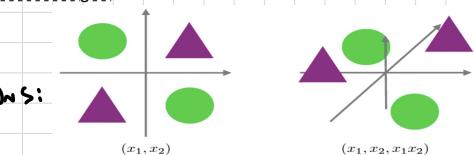
<https://medium.com/@lucaspereira0612/solving-xor-with-a-single-perceptron-34539f395182>

\rightarrow NOT LINEARLY SOLVABLE, BUT IN 3 DIMENSIONS:

$$(x_1, x_2) \mapsto (x_1, x_2, x_1 x_2)$$

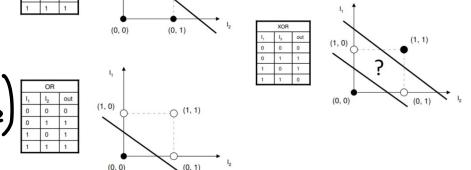
\rightarrow IN 3D IT IS LINEARLY SOLVABLE

- QUADRATIC FEATURE: $\phi(x) := (x_1^2, \sqrt{2} x_1 x_2, x_2^2)$



AND		
I_1	I_2	out
0	0	0
0	1	0
1	0	0
1	1	1

OR		
I_1	I_2	out
0	0	0
0	1	1
1	0	1
1	1	1



- KERNEL FUNCTION: <https://youtu.be/Q7vT0--5VII>

FEATURE EXTRACTING CAN BE VERY COSTLY; e.g. 2^{nd} ORDER FEATURES IN 1000 DIMENSIONS

SOL: DON'T COMPUTE $\langle \dots, \dots \rangle$, COMPUTE $K(x, x')$ KERNEL FUNCTION

DEF.: $K: X \times X \rightarrow \mathbb{R} / K(x, x') = \langle \phi(x), \phi(x') \rangle$

↳ SOMETIMES $K(x, x')$ IS MUCH CHEAPER THAN $\phi(x)$

- ALGORITHM: SCOR 66/86

$$f(x) = \sum_{i \in I} \gamma_i \langle \phi(x_i), \phi(x) \rangle + b = \sum_{i \in I} \gamma_i K(x_i, x) + b$$

POLYNOMIAL KERNELS

Ex.

Let's plug in some numbers to make this more intuitive: suppose $x = (1, 2, 3)$; $y = (4, 5, 6)$. Then:
 $f(x) = (1, 2, 3, 2, 4, 6, 3, 9)$
 $f(y) = (16, 20, 24, 20, 25, 30, 24, 36)$
 $\langle f(x), f(y) \rangle = 16 + 40 + 72 + 40 + 100 + 180 + 72 + 180 + 324 = 1024$

Simple Example: $x = (x_1, x_2, x_3); y = (y_1, y_2, y_3)$. Then for the function $f(x) = (x_1x_1, x_1x_2, x_1x_3, x_2x_1, x_2x_2, x_2x_3, x_3x_1, x_3x_2, x_3x_3)$, the kernel is $K(x, y) = \langle x, y \rangle^2$.

- MERCER'S THEOREM:

EXPLICIT CALCULUS
 USING $K(x, y)$

Now let us use the kernel instead:
 $K(x, y) = (4 + 10 + 18)^2 = 32^2 = 1024$
 Same result, but this calculation is so much easier.

\forall symmetric $K: X \times X \rightarrow \mathbb{R} / K$ is SQUARE INTEGRABLE IN $X \times X$

AND SAMPLERS:

$$\int_{X \times X} K(x, x') \cdot f(x) \cdot f(x') dx \times dx' \geq 0, \forall f \in L_2(X)$$

$$\Rightarrow \exists \phi_i: X \rightarrow \mathbb{R}, \exists \lambda_i \geq 0$$

$$K(x, x') = \sum_i \lambda_i \phi(x) \phi(x'), \forall x, x' \in X$$

→ IT GIVES INFO ON WHICH FUNCTIONS CAN BE KERNEL FUNCTIONS

• DISTANCE IN FEATURE SPACE :

$$d^2(x, x')^2 := \|\phi(x) - \phi(x')\|^2 = K(x, x) + K(x', x') - 2K(x, x')$$

• KERNEL MATRIX : $K_{ij} := \langle \phi(x_i), \phi(x_j) \rangle = K(x_i, x_j)$

$\hookrightarrow K_{ij}$ tells us the similarity of $\phi(x_i)$ and $\phi(x_j)$

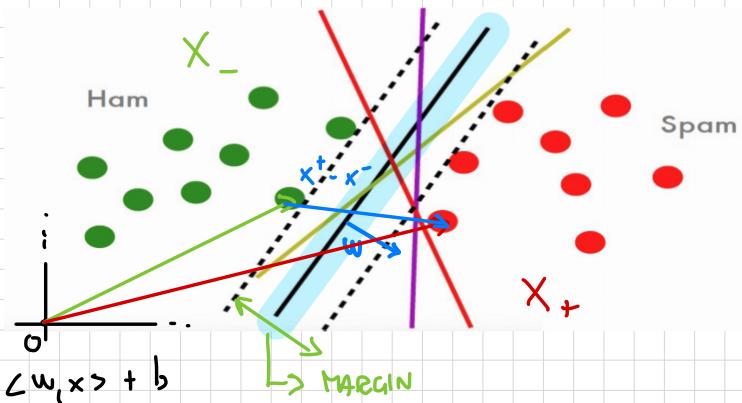
SVM:

LARGE MARGIN CLASSIFIER:

• $\langle w, x \rangle + b \leq -1$

• $\langle w, x \rangle + b \geq 1$

• LINEAR FUNCTION: $f(x) = \langle w, x \rangle + b$



MARGIN OPTIMIZATION:

• MARGIN: $\frac{\langle x_+ - x_-, w \rangle}{\|w\|}$ $\xrightarrow{\text{dot product is a measure of similarity}}$

OR $\frac{1}{\|w\|}$ half margin

$$= \frac{1}{\|w\|} \left[[\underbrace{\langle x_+, w \rangle + b}_{-1}] - [\underbrace{\langle x_-, w \rangle + b}_{-1}] \right] \approx \frac{2}{\|w\|}$$

OPTIMIZATION PROBLEM:

MARGIN "MOTH" MAXIMIZATION • MAXIMIZE $\frac{1}{\|w\|}$, CONSTRAINT: $\forall i [\langle x_i, w \rangle + b] \geq 1$

MINIMIZE L2 REGULARIZATION
↳ HOW WELL DATA ARE CLASSIFIED • MINIMIZE $\frac{1}{2} \|w\|^2$, CONSTRAINT: $\forall i [\langle x_i, w \rangle + b] \geq 1$

→ LAGRANGE METHOD:

• $L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_i \alpha_i [\forall_i [\langle x_i, w \rangle + b] - 1]$

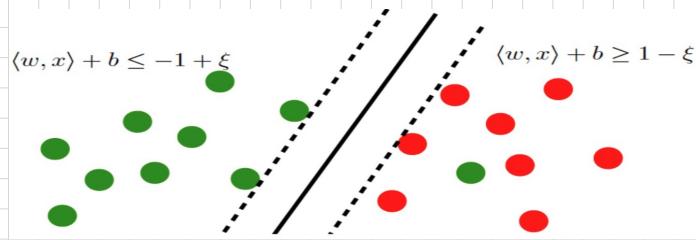
↳ ... $\rightarrow w = \sum_i \alpha_i x_i$, $\forall_i [\langle w, x_i \rangle + b] = 1$

COSTRAINT

• SOFT MARGIN CLASSIFIER :

- IN A DATASET LIKE IN FIGURE, HARD MARGINS ARE NOT POSSIBLE.

→ NEED TO "RELAX" MARGINS → ADDING A SLACK VARIABLE ξ .



$$\cdot \langle w, x \rangle + b \leq -1 \longrightarrow \langle w, x \rangle + b \leq -1 + \xi$$

$$\cdot \langle w, x \rangle + b \geq 1 \longrightarrow \langle w, x \rangle + b \geq 1 - \xi$$

- OPTIMIZATION PROBLEM:

$$\cdot \underset{w, b}{\text{MINIMIZE}} \quad \frac{1}{2} \|w\|^2 + C \sum_i \xi_i, \text{ constraint: } y_i [\langle w, x_i \rangle + b] \geq 1 - \xi_i, \xi_i \geq 0$$

\hookrightarrow ALWAYS FEASIBLE

→

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 + C \sum_i \xi_i - \sum_i \alpha_i [y_i (\langle x_i, w \rangle + b) + \xi_i] - \sum_i \eta_i \xi_i$$

\hookrightarrow

... → KARUSH-KHAN-TUCKER CONDITIONS:

$$\left\{ \begin{array}{l} \alpha_i = 0 \rightarrow y_i [\langle w, x_i \rangle + b] \geq 1 \\ 0 < \alpha_i < C \rightarrow y_i [\langle w, x_i \rangle + b] = 1 \\ \alpha_i = C \rightarrow y_i [\langle w, x_i \rangle + b] \leq 1 \end{array} \right.$$

- NON-LINEAR SEPARATION:

$$\rightarrow \text{KERNEL TRICK: } \begin{cases} x_i \mapsto \phi(x_i) \\ \langle x_i, x \rangle \mapsto \kappa(x_i, x), \langle x_i, x_j \rangle \mapsto \kappa(x_i, x_j) \end{cases}$$

→

- Linear soft margin problem

$$\underset{w, b}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$$

subject to $y_i [\langle w, x_i \rangle + b] \geq 1 - \xi_i$ and $\xi_i \geq 0$

- Dual problem

$$\underset{\alpha}{\text{maximize}} \quad -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_i \alpha_i$$

subject to $\sum_i \alpha_i y_i = 0$ and $\alpha_i \in [0, C]$

- Support vector expansion

$$f(x) = \sum_i \alpha_i y_i \langle x_i, x \rangle + b$$

- Linear soft margin problem

$$\underset{w, b}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$$

subject to $y_i [\langle w, \phi(x_i) \rangle + b] \geq 1 - \xi_i$ and $\xi_i \geq 0$

- Dual problem

$$\underset{\alpha}{\text{maximize}} \quad -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \kappa(x_i, x_j) + \sum_i \alpha_i$$

subject to $\sum_i \alpha_i y_i = 0$ and $\alpha_i \in [0, C]$

- Support vector expansion

$$f(x) = \sum_i \alpha_i y_i \kappa(x_i, x) + b$$

REGRESSION:

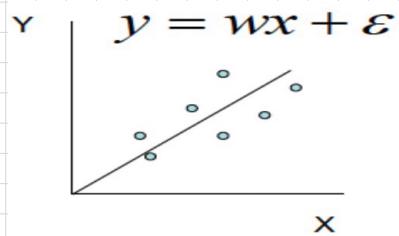
X : INPUT FEATURES

$$y = X \cdot w + \epsilon \quad \text{with } \epsilon \text{ noise / other parameters}$$

LINEAR REGRESSION:

GOAL: ESTIMATE w FROM $\langle x_i, y_i \rangle$ PAIRS

→ MINIMIZING LSE (LEAST SQUARES ERROR).



$$\underset{w}{\operatorname{arg\,min}} \sum_i (y_i - w x_i)^2$$

$$\rightarrow \frac{1}{\sqrt{m}} \sum_i (y_i - w x_i)^2 = 0 \rightarrow \dots \rightarrow w = \frac{\sum_i x_i y_i}{\sum_i x_i^2}$$

• BIAS TERM w_0 :

$$y = w_0 + w_1 x + \epsilon$$

$$\left\{ \begin{array}{l} w_0 = \frac{\sum_i y_i - w_1 x_i}{m} \\ w_1 = \frac{\sum_i x_i (y_i - w_0)}{\sum x_i^2} \end{array} \right.$$

MULTIVARIATE REGRESSION:

$$y = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_k x_k + \epsilon$$

BUT ALSO...

$$y = w_0 + w_1 x_1^2 + \dots + w_k x_k^2 + \epsilon$$

↳ FEATURES ARE POLYNOMIAL

BUT REGRESSION IS LINEAR

$$x \mapsto \phi(x) = x^k \quad (\text{in this case})$$

- Polynomial: $\phi_j(x) = x^j$ for $j=0 \dots n$

- Gaussian: $\phi_j(x) = \frac{(x - \mu_j)}{2\sigma_j^2}$

- Sigmoid: $\phi_j(x) = \frac{1}{1 + \exp(-s_j x)}$

GENERAL LINEAR REGRESSION PROBLEM:

$$\cdot \quad V = \sum_{j=0}^n w_j \phi_j(x)$$

→ MINIMIZE LSE :

$$\underset{w}{\text{argmin}} \quad J(w) = \sum_i (V_i - \sum_j w_j \phi_j(x^i))^2$$

GENERIC
POLYNOMIAL REGRESSION

similar to LINEAR REGRESSION

$$\rightarrow \frac{\partial}{\partial w} J(w) = 0 \rightarrow \dots \rightarrow \sum_i V^i \phi(x^i)^T = w^T \left[\sum_i \phi(x^i) \phi(x^i)^T \right]$$

$\sum_i \phi(x^i)^2$ IN MATRIX NOTATION

IF WE DEFINE $\phi = \begin{pmatrix} \phi_0(x^1) & \phi_1(x^1) & \dots & \phi_m(x^1) \\ \phi_0(x^2) & \phi_1(x^2) & \dots & \phi_m(x^2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x^n) & \phi_1(x^n) & \dots & \phi_m(x^n) \end{pmatrix}$

POWERS OF EACH FEATURE

DIFFERENT FUNCTION ON EACH FEATURE.

$$\rightarrow w = \begin{pmatrix} \phi^T & \phi \end{pmatrix}^{-1} \cdot \phi^T \cdot V$$

PSEUDO-INVERSE

ALSO MLE FOR w IS THE SOLUTION ABOVE

LOGISTIC REGRESSION:

<https://bit.ly/3UDyedq>

USEFUL WHEN
Y IS NOT CATEGORICAL

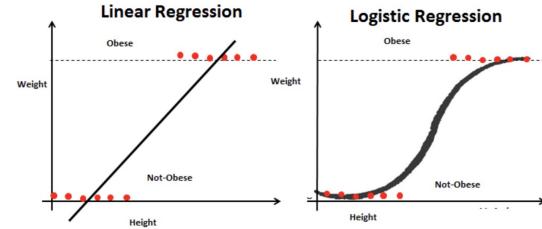
WHEN OUTPUT IS BINARY / DISCRETE $\rightarrow \exists$ MORE EFFICIENT MODELS THAN LINEAR REGRESSION
FOR CLASSIFICATION

• CLASSIFICATION : $p(Y|X; \theta) / \theta$: PARAMETERS \vec{w} OF REGRESSION

• WHEN Y IS BINARY \rightarrow LOGISTIC REGRESSION
IS BETTER THAN LINEAR REGRESSION

\rightarrow LINEAR FUNCTION \mapsto SIGMOID FUNCTION : $y(h) = \frac{1}{1 + e^{-h}}$

$$\text{L, } \begin{cases} p(Y=0|X; \theta) = g(w^T x) = \frac{1}{1 + e^{w^T x}} \\ p(Y=1|X; \theta) = 1 - g(w^T x) = \frac{e^{w^T x}}{1 + e^{w^T x}} \end{cases}$$



• HOW TO OBTAIN PARAMETERS $\theta = \vec{w}$? \rightarrow MLE

$$\rightarrow L(Y|X; w) = \prod_i (1 - y(x_i; w))^{Y_i} \cdot y(x_i; w)^{1 - Y_i}$$

• MLE : $\sum_{i=1}^N L(Y|X; w) = 0$ DO NOT LEAD TO A CLOSED-FORM SOLUTION

\hookrightarrow WE USE GRADIENT TO ADJUST w_i :

GRADIENT ASCENT STEP

$$w^{(t)} \leftarrow w^{(t-1)} + \varepsilon \sum_{i=1}^N X_i \{Y_i - (1 - y(x_i; w))\}$$

• ALGORITHM FOR LOGISTIC REGRESSION :

1. CHOOSE λ

2. START WITH A GUESS FOR \vec{w}

$$3. \forall t. \quad w^{(t)} \leftarrow w^{(t-1)} + \varepsilon \sum_{i=1}^N X_i \{Y_i - (1 - y(x_i; w))\}$$

4. IF NO IMPROVEMENT FOR $\log L(Y|X; \vec{w}) \rightarrow$ STOP

\hookrightarrow else : GO TO STEP 3.

PRIOR INFO

• IF NOT ENOUGH SAMPLE \rightarrow REGULARIZATION : $p(Y=1, \theta|X) \propto p(Y=1|X; \theta)p(\theta)$

MULTI LAYER PERCEPTRON :

- 1 PERCEPTRON \rightarrow LINEAR BOUNDARIES
- MULTI LAYER PERCEPTRON \rightarrow MORE COMPLEX BOUNDARIES : NON-LINEAR

• GIVEN A PERCEPTRON :

\Rightarrow WE NEED TO SOLVE THE PROBLEM :

$$\bar{w} = \underset{\bar{w}}{\operatorname{arg\,min}} \sum \frac{1}{2} (r_i - \hat{f}(x_i, \bar{w}))^2$$

PERCEPTRON LEARNING RULE :

- \forall TRAINING EXAMPLE $d \in D$:

1. COMPUTE GRADIENT $\nabla E_d [\bar{w}]$:

$$\bullet E_d [\bar{w}] = \text{TOT ERROR FUNCTION} = \frac{1}{2} \sum (t_d - o_d)^2$$

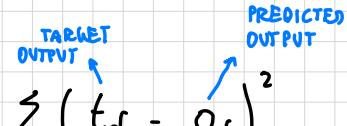
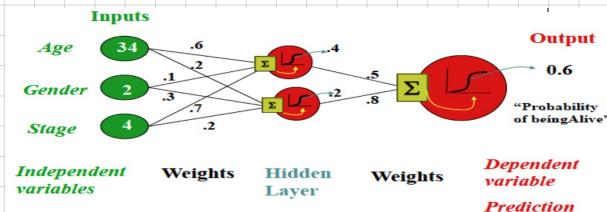
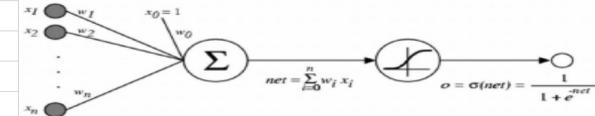
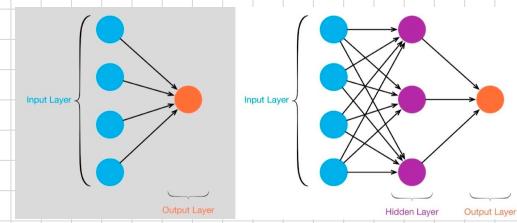
$$\hookrightarrow \nabla E_d [\bar{w}] = \left[\frac{\partial E_d}{\partial w_0}, \frac{\partial E_d}{\partial w_1}, \dots, \frac{\partial E_d}{\partial w_n} \right]$$

2. UPDATE \bar{w} USING THE TRAINING RULES :

PROOF:

$$\bullet \bar{w} \leftarrow \bar{w} - \eta \cdot \nabla E_d [\bar{w}] \quad / \eta: \text{LEARNING RATE}$$

$$\hookrightarrow \frac{\partial E_d}{\partial w_j} = \frac{\partial}{\partial w_j} \frac{1}{2} \sum_d (t_d - o_d)^2 = \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_j} (t_d - o_d) = \sum_d (t_d - o_d) \left(-\frac{\partial o_d}{\partial w_j} \right) = -\sum_d (t_d - o_d) \frac{\partial o_d}{\partial \text{net}_d} \frac{\partial \text{net}_d}{\partial w_j} = -\sum_d (t_d - o_d) o_d (1 - o_d) x_d^j$$



BACK PROPAGATION :

<https://towardsdatascience.com/how-does-back-propagation-work-in-neural-networks-with-worked-example-bc59dfb97f48>

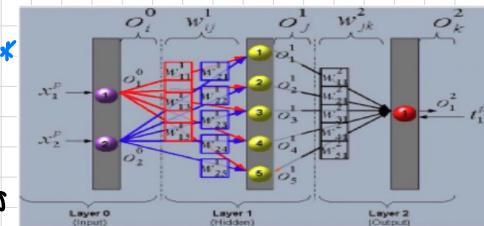
ONCE THE NETWORK OUTPUT IS COMPUTED, AN ERROR IS GENERATED ON NN, BETWEEN THE PREDICTED OUTPUT AND THE TRUE OUTPUT

→ THIS ERROR IS USED TO ADJUST THE \bar{w} OF THE NN.

WEIGHT \bar{w} UPDATE :

$$w_{i,j} \leftarrow w_{i,j} + \eta \cdot S_j \cdot x^j$$

UPDATE



* CONSIDERING THIS NN

WHERE:

~ UPDATE IS DIFFERENT BASED ON IF $w \in$ OUTPUT/MIDDLE LAYER

$$\cdot w_{i,s} \in \text{OUTPUT UNIT} : S_j = S_k \leftarrow O_n^2 (1 - O_n^2) (t - O_k^2)$$

$$\cdot w_{i,s} \in \text{MIDDLE UNIT} : S_j = S_h \leftarrow O_h^j (1 - O_h^j) \cdot \sum_{k \in \text{outputs}} w_{h,k} \cdot S_k$$

this BACK PROPAGATION IS DOING GRADIENT DESCENT OVER THE ENTIRE WEIGHTS OF THE NETWORK

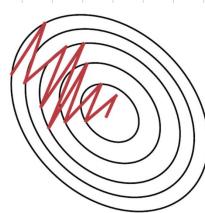
IT WILL FIND A LOCAL MINIMUM, NOT A GLOBAL ONE

A NEW PARAMETER γ called "WEIGHT MOMENTUM" CAN BE

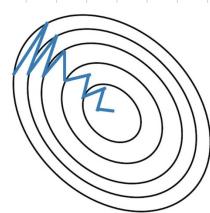
USED TO SPEED UP LEARNING AND NOT GET STUCK IN A LOCAL MINIMUM

→

$$\Delta w_{i,j} = \underbrace{\left(\eta \cdot \frac{S_E}{S_{w_{i,j}}} \right)}_{\text{BASIC UPDATE RULE}} + \underbrace{\left(\gamma \cdot \Delta w_{i,j}^{(t-1)} \right)}_{\text{WEIGHT UPDATE AT PREVIOUS ITERATION}}$$



Stochastic Gradient Descent without Momentum

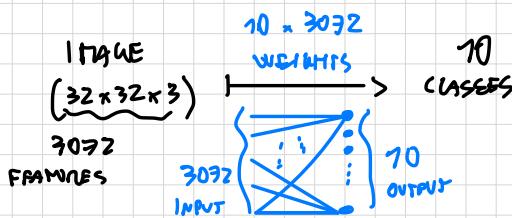


Stochastic Gradient Descent with Momentum

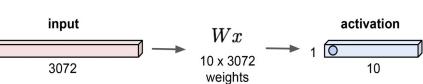
• CONVOLUTIONAL NEURAL NETWORKS (CNN) :

• RECAP :

FC \rightarrow FULL CONNECTED
UNFLATTEN :



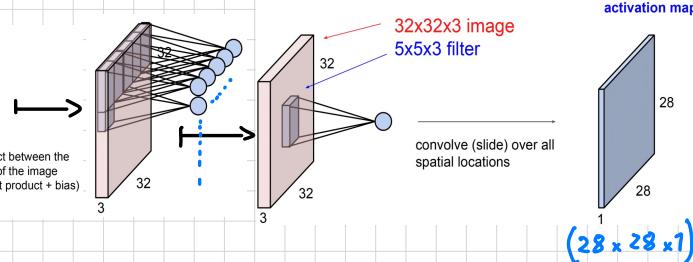
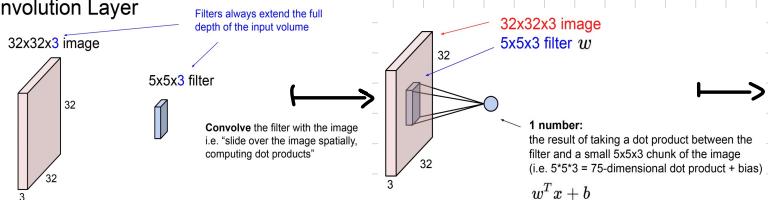
32x32x3 image \rightarrow stretch to 3072 x 1



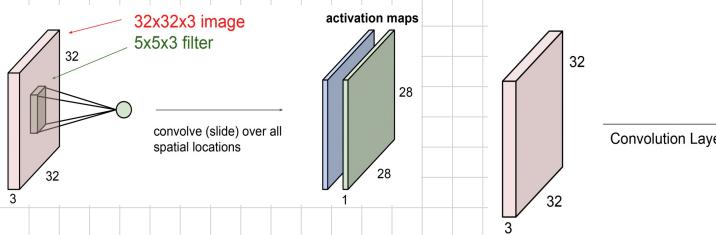
\rightarrow CONVOLUTIONAL LAYER :

A FILTER IS CONVOLVED OVER THE ORIGINAL IMAGE, AND THE FULL DEPTH OF THE INPUT VOLUME IS PRESERVED :

Convolution Layer

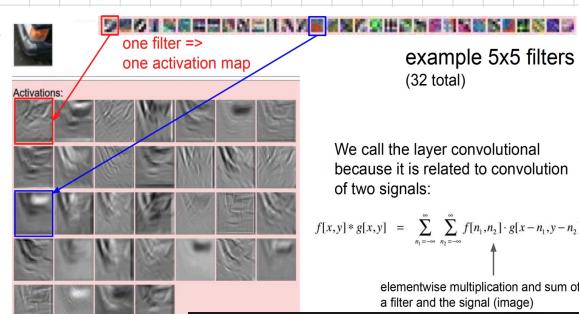
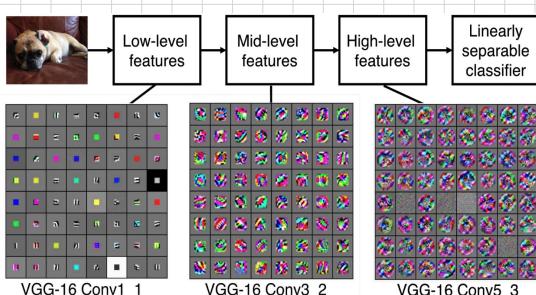


\rightarrow MANY FILTERS CAN BE APPLIED ON THE SAME IMAGE :



• 20. 6 FILTERS
OF $5 \times 5 \times 3$ APPLIED
ON SAME IMAGE $32 \times 32 \times 3$:
 \rightarrow NEW "IMAGE" $28 \times 28 \times 6$

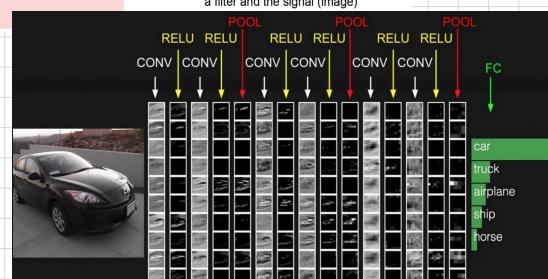
• AFTER A FILTER IS APPLIED AN ACTIVATION FUNCTION CAN BE USEFUL TO ACTIVATE ONLY MOST RELEVANT FRAMES:



We call the layer convolutional because it is related to convolution of two signals:

$$f[x, y] * g[x, y] = \sum_{n_1=0}^{m_1} \sum_{n_2=0}^{m_2} f[n_1, n_2] \cdot g[x - n_1, y - n_2]$$

elementwise multiplication and sum of a filter and the signal (image)



- CLOSER LOOK AT SPATIAL DIMENSION:

ex.

- INPUT: 7×7
- FILTER: 3×3

- STRIDE: OFFSET BETWEEN SUCCESSIVE POSITIONS OF A FILTER

• FILTER APPLIED WITH STRIDE = 2 ? → OK: 3×3 OUTPUT

• FILTER APPLIED WITH STRIDE = 3 ? → NOT POSSIBLE

- IN GENERAL:

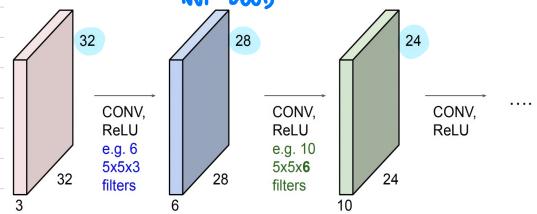
$$\text{OUTPUT SIZE} = \frac{(N - F)}{\text{STRIDE}} + 1 \rightarrow \text{IF OUTPUT SIZE } \notin \mathbb{Z} : \text{NOT POSSIBLE}$$

- IN PRACTICE IT IS COMMON TO ADD ZERO-PADDING TO THE BORDER

↳ USUALLY: STRIDE = 1, $P = \frac{1}{2}(F-1)$

$$\text{OUTPUT SIZE} = \frac{(N + 2P - F)}{\text{STRIDE}} + 1$$

... BUT SHRINKING TOO MUCH IS
NOT GOOD



ex.

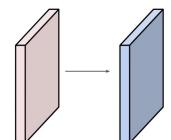
→ $N = 32, F = 5, \text{stride} = 1, P = 2$

- DEPTH = 10 ∵ 10 FILTERS APPLIED

• OUTPUT SIZE = $\frac{(32 + 2 \cdot 2 - 5)}{1} + 1 = 32$

→ OUT SHAPE = $(32 \times 32 \times 10)$

Examples time:



Input volume: $32 \times 32 \times 3$
10 ~~filters~~ filters with stride 1, pad 2

$5 \times 5 \times 3$

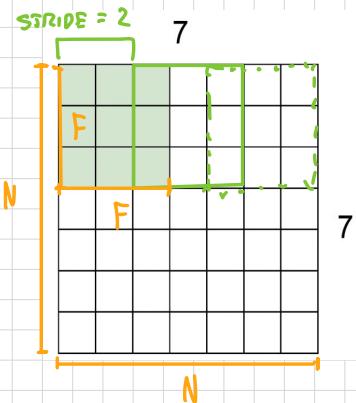
Output volume size: ?

Number of parameters in this layer?

- PARAMETERS: VALUES IN THE FILTER TO AGGREGATE

→ # FILTER: $5 \cdot 5 \cdot 3$ WEIGHTS + 1 BIAS = 76 PARAMETERS/FILTER

→ # FILTERS = 10 → 760 PARAMETERS IN TOTAL



• SUMMARY:

ASSUME INPUT = $W_1 \times H_1 \times C$

- CONV. LAYER NEEDS 4 HYPERPARAMETERS:

- K : n° OF FILTERS
- F : SIZE OF FILTER
- S : STRIDE
- P : ZERO-PADDING

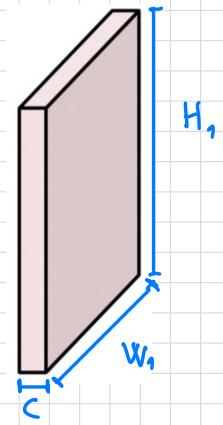
→ OUTPUT : $W_2 \times H_2 \times K$

→ n° OF PARAMETERS : $F^2 C \cdot K + K$
↳ BIAS

Common settings:

K = (powers of 2, e.g. 32, 64, 128, 512)

- $F = 3, S = 1, P = 1$
- $F = 5, S = 1, P = 2$
- $F = 5, S = 2, P = ?$ (whatever fits)
- $F = 1, S = 1, P = 0$



$$W_2 = \frac{(W_1 + 2P - F)}{S} + 1$$

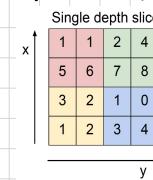
$$H_2 = \frac{(H_1 + 2P - F)}{S} + 1$$

INTRO TO DEEP LEARNING 1.

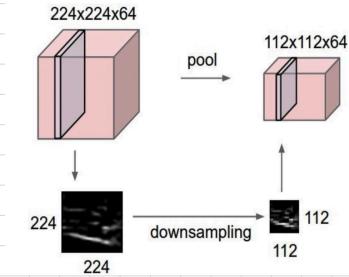
POOLING :

EXTRACTION OF ONLY MOST RELEVANT VALUES

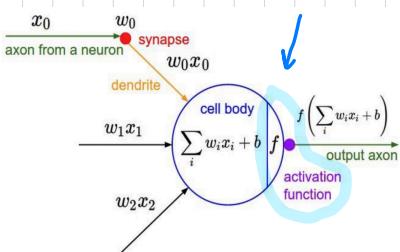
e.g. MAX-POOLING



max pool with 2x2 filters and stride 2

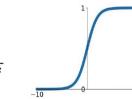


ACTIVATION FUNCTIONS :



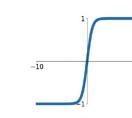
Sigmoid

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



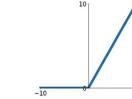
tanh

$$\tanh(x)$$



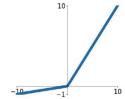
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

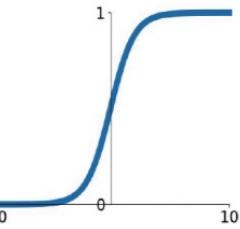


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



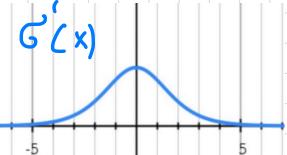
SIGMOID :

$$G(x) = \frac{1}{1+e^{-x}}, \quad G : \mathbb{R} \mapsto [0, 1]$$

PROBLEMS :

1. SATURATED NEURONS "KILL" THE GRADIENTS :

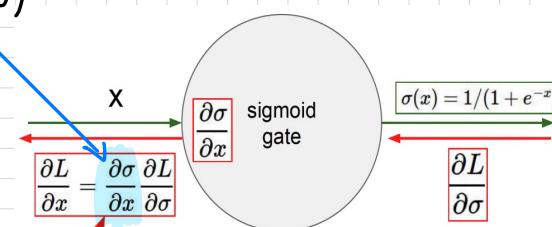
$$\cdot \frac{dG(x)}{dx} = G(x) \cdot (1 - G(x))$$



→ WEIGHT CAN NEVER CHANGE

IF ALL GRADIENTS

FLOWING BACK WILL BE 0



$$\text{INPUT } X = -10 : G(x) \approx 0, \frac{dG(x)}{dx} = 0 \cdot (1-0) = 0$$

$$\text{INPUT } X = 0 : G(x) \approx 1, \frac{dG}{dx} = 1 \cdot (1-1) = 0$$

$$\text{INPUT } X > 10 : G(x) = 0, \frac{dG}{dx} = 0 \cdot (1-0) = 0$$

→ $G(x)$ DOESN'T HELP IN LEARNING FAST :: OFTEN $\frac{dG}{dx}$ WILL BE 0

2. OUTPUT ARE NOT ZERO-CENTERED :

- CONSIDER WHAT HAPPENS WHEN ALL $X > 0$:

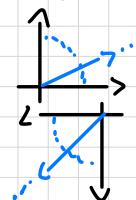
LOCAL GRADIENT OF $G(x)$ IS ALWAYS > 0

$$\frac{\partial L}{\partial w} = \sigma(\sum_i w_i x_i + b)(1 - \sigma(\sum_i w_i x_i + b))x \times \text{upstream_gradient}$$

$$\rightarrow \text{sign}\left(\frac{\partial L}{\partial w}\right) = \text{sign}\left(\text{upstream gradient}\right)$$

→ GRADIENT UPDATE WILL MOVE ONLY 2 DIRECTIONS:

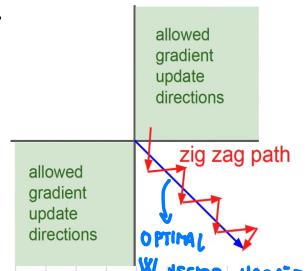
① ALL POSITIVE:



② ALL NEGATIVE:

(not shown)

→ THIS WILL LEAD TO A ZIG-ZAG PATH

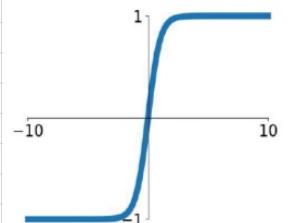


3. $\exp()$ IS A BIT COMPUTE EXPENSIVE

tanh():

V: ZERO-CENTERED

S: KILLS GRADIENT WHEN SATURATED



ReLU:

$$\text{ReLU}(x) = \max(0, x), \text{ ReLU}: \mathbb{R} \mapsto [0, +\infty]$$

V:

DO NOT SATURATE IN + REGION

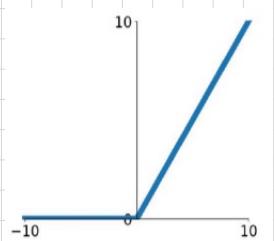
COMPUTATIONALLY EFFICIENT

CONVERGES MUCH FASTER THAN $G(x)$, $\tanh(x)$ ($6x$)

S:

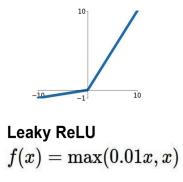
NOT ZERO-CENTERED

GRADIENT WHEN $X < 0$



- OTHERS:

LEAKY ReLU:



- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- will not "die".

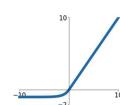
Parametric Rectifier (PReLU)
 $f(x) = \max(\alpha x, x)$

backprop into ' α ' (parameter)

EXPONENTIAL

ELU:

Exponential Linear Units (ELU)



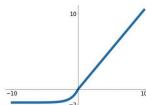
- All benefits of ReLU
- Closer to zero mean outputs
- Negative saturation regime compared with Leaky ReLU adds some robustness to noise

- Computation requires $\exp()$

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$

(Alpha default = 1)

Scaled Exponential Linear Units (SELU)



$f(x) = \begin{cases} \lambda x & \text{if } x > 0 \\ \lambda \alpha (\exp(x) - 1) & \text{otherwise} \end{cases}$

- Scaled version of ELU that works better for deep networks
- "Self-normalizing" property;
- Can train deep SELU networks without BatchNorm
- (will discuss more later)

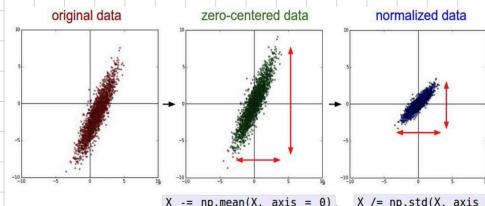
- IN PRACTICE:

- USE ReLU, IN CASE TRY LEAKY ReLU, SELU
- DO NOT USE Tanh(), G(x)

- DATA PREPROCESSING:

- NORMALIZATION : OR STANDARDIZATION

$$\frac{x - \mu}{\sigma}$$



WHITENING

- PCA OR WHITENING : $X / \text{VarCoV matrix} \rightarrow X / \text{VarCoV matrix} = I$

↳ Practice for images: SUBTRACT MEAN IMAGES, SUBTRACT PER CHANNEL MEAN AND \div STD

- WEIGHT INITIALIZATION:

- \bar{W} = "small random numbers". OK FOR SMALL NETWORK, PROBLEMS WITH DEEPER NETWORKS

→ ALL ACTIVATIONS TEND TO ZERO FOR DEEPER NETWORK LAYERS

→ GRADIENT $\frac{\partial L}{\partial w} \rightarrow 0 \rightarrow$ NO LEARNING

- \bar{W} = "XAVIER INITIALIZATION": $W = np.random.randn(Din, Dout) / np.sqrt(Din)$

→ STD = $\sqrt{D_{in}} \rightsquigarrow$ SIZE OF \bar{X} / FOR CONV. LAYERS: $D_{in} = F^2 \cdot \text{INPUT_CHANNELS}$

↳ ACTIVATIONS ARE NICELY SCALED FOR ALL LAYERS

- IT DOES NOT WORK WELL FOR ReLU: XAVIER ASSUMES ZERO-CENTERED ACTIVATION

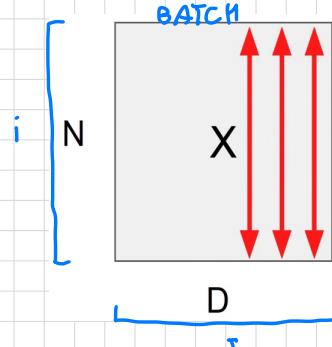
- \bar{W} = "KAIMING / MSRA INITIALIZATION": $W = np.random.randn(Din, Dout) * np.sqrt(2/Din)$

→ STD = $\sqrt{\frac{2}{D_{in}}}$

→ IT SOLVES XAVIER PROBLEMS WITH NON-ZERO CENTERED FUNCTIONS

• BATCH NORMALIZATION :

INPUT $X : N \times D$



$$\text{PER CHANNEL MEAN : } M_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

shape $(D,)$

$$\text{PER CHANNEL VAR : } G_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - M_j)^2$$

shape $(D,)$

$$\rightarrow \text{NORMALIZED } X : \hat{x}_{i,j} = \frac{x_{i,j} - M_j}{\sqrt{G_j^2 + \epsilon}} \quad \text{shape } (N \times D)$$

- OUTPUT : $\hat{y}_{i,j} = \gamma_j \cdot \hat{x}_{i,j} + \beta_j \quad \text{shape } (N \times D)$

\hookrightarrow LEARNING $\gamma = \sigma$, $\beta = \mu$ WILL RECOVER THE IDENTITY FUNCTION

- AT TEST TIME M_j AND G_j^2 COMES FROM VALUES SEEN DURING TRAINING

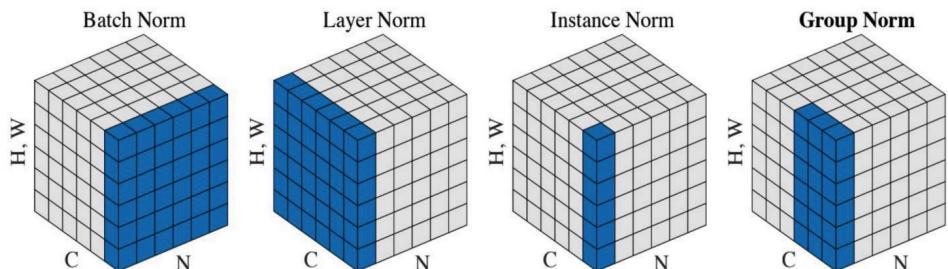
- BN IS USUALLY INSERTED AFTER FC OR CONVOLUTIONAL LAYERS
AND BEFORE NON-LINEARITY (ACTIVATION FUNCTIONS)

V:

- IT MAKES DEEP NETWORKS MUCH EASIER TO TRAIN
- IMPROVEMENTS ON GRADIENT FLOW
- ALLOWS HIGHER LEARNING RATES \rightarrow FASTER CONVERGENCE
- NETWORKS MORE ROBUST TO INITIALIZATION

S:

- BEHAVES DIFFERENTLY DURING TRAINING AND TESTING (COMMON SOURCE OF BUGS)



• TRANSFER LEARNING :

USE PARAMETERS LEARNED FROM A NN USED IN A TASK, INTO A NEW SIMILAR NN, FOR SIMILAR PURPOSES

- so. TRAINING ON IMAGES:

1. TRAIN NN ON IMAGENET

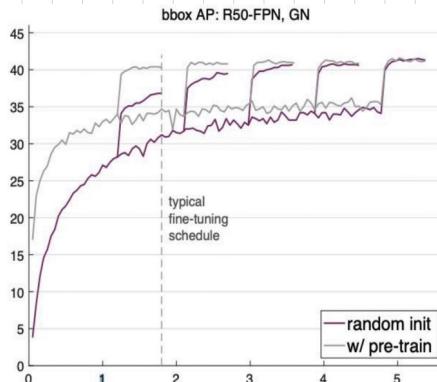
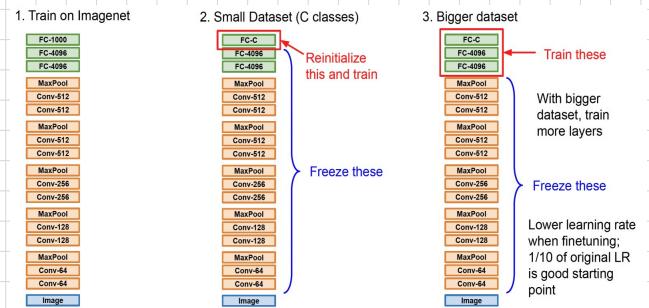
2.

IF SMALL DATASET:

-> FREEZE ALL LAYERS AND
REINITIALIZE LAST ONE

IF BIG DATASET:

-> USE LOWER LEARNING RATE ($\approx \frac{1}{10}$ η) AND REINITIALIZE AND TRAIN MORE LAYERS



Training from scratch can work just as well as training from a pretrained ImageNet model for object detection

But it takes 2-3x as long to train.

They also find that collecting more data is better than finetuning on a related task

INTRO TO DEEP LEARNING 2 :

LEARNING THEORY:

(slides 1-32, covered section in MML)

OPTIMIZATION:

HINGE LOSS:

IN SVMs, IT PENALIZES DATA POINTS WHICH ARE INSIDE THE MARGIN. SVMs MINIMIZE HINGE LOSS

$$\cdot L = \max(0, 1 - y \cdot f(x))$$

• IN ORDER TO REACH A MINIMUM IN LOSS WE NEED TO TAKE SMALL STEPS IN DIRECTION OF THE NEGATIVE GRADIENT: $-\nabla_w L(w)$

$$\rightarrow \text{WEIGHTS UPDATE: } w^{(t+1)} = w^t - \alpha \cdot \nabla_w L(w)$$

↳ 6-D PROBLEM: SADDLE POINTS → HERE $\nabla(\cdot)$ VANISHES

• GD IS OFTEN COMPUTED ON MINI-BATCHES, NOT FULL DATASET

↳ THIS APPROACH IS CALLED SGD (STOCHASTIC GRADIENT DESCENT)

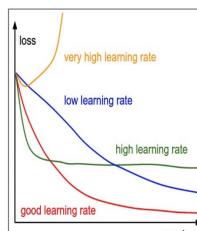
↳ IT IS "STOCHASTIC" BECAUSE IT USES THE WEIGHTS OF A MINI-BATCH INSTEAD OF THE ONE OF THE FULL DATASET

• MOMENTUM: TAKES INTO ACCOUNT PREVIOUS STEPS OF THE WEIGHTS (ϵ_{-i})

- RMS PROP

- ADA GRAD

- ADAM



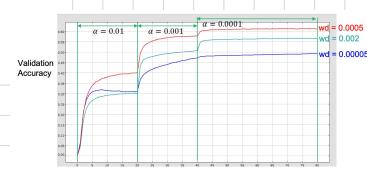
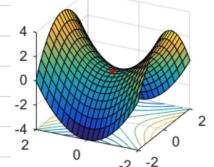
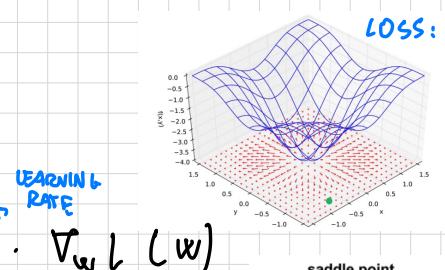
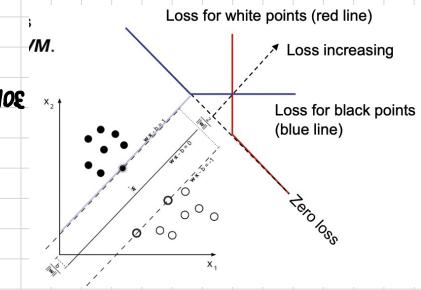
=> Learning rate decay over time!

step decay: e.g. decay learning rate by half every few epochs.

exponential decay: $\alpha = \alpha_0 e^{-kt}$

$1/t$ decay: $\alpha = \alpha_0 / (1 + kt)$

$1/\sqrt{t}$ decay (ADAGRAD): $\alpha = \alpha_0 / \sqrt{1 + kt}$



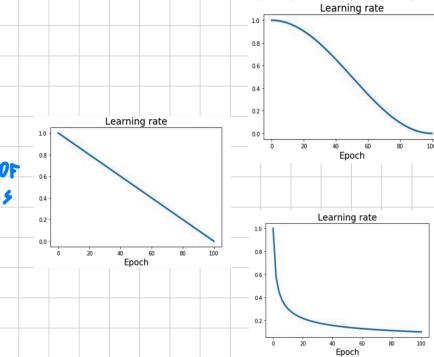
• LEARNING RATE SCHEDULES:

IT IS USEFUL TO CHANGE LR OVER EPOCHS:

- COSINE: $\alpha_t = \frac{1}{2} \alpha_0 \left(1 + \cos \left(\pi \frac{t}{T} \right) \right)$
- LINEAR: $\alpha_t = \alpha_0 \left(1 - \frac{t}{T} \right)$
- INVERSE SQRT: $\alpha_t = \alpha_0 / \sqrt{t}$

INITIAL LR

TOT N OF EPOCHS



- PRACTICAL TIP: BATCH SIZE += N \rightarrow LR * = N
- SECOND ORDER OPTIMIZATION: STEP TO MINIMUM OF PARABOLA

• IMPROVE TEST ERRORS:

- STOP TRAINING MODE: WHEN ACC ON VALIDATION SET STARTS DECREASING
- MODEL ENSEMBLES: TRAIN MULTIPLE INDEPENDENT MODELS AND AVERAGE THEIR RESULTS AT TEST TIME
- IMPROVE SINGLE MODEL PERFORMANCE:
 - REGULARIZATION: L2, L1, ELASTIC NET ($L2 + L1$)
 - DROPOUT: 1 FORWARD PASS \rightarrow RANDOMLY SET SOME NEURONS TO 0
 - \rightarrow AT TRAINING: ADD SOME RANDOMNESS $\rightarrow V = f_W(x, z)$
 - \rightarrow AT TESTING: AVERAGE OUT RANDOMNESS $\rightarrow V = f(x) = \mathbb{E}_z [f(x, z)]$
 - DATA AUGMENTATION: CREATE NEW VERSION OF SAME DATA
 - \hookrightarrow es. IMAGES: HORIZONTAL FLIP, CROPS, CHANGING BRIGHTNESS, etc.
 - DROPCONNECT: RANDOMLY SET SOME $W=0$ IN NN AT TRAINING TIME
 - etc. (see 117 - 121)

DROPOUT MASK

• CHOOSING HYPERPARAMETERS:

• STEPS:

1. CHECK INITIAL LOSS L

2. OVERFIT A SMALL SAMPLE:

TRY TO TRAIN 100% ACC ON A

SMALL SAMPLE OF DATA

LOSS NOT GOING DOWN

$\uparrow L$: LR too low

$\downarrow L \rightarrow +\infty_{N \rightarrow N}$: LR too high

3. FIND LR / $\downarrow L$:

USE ARCHITECTURE FROM PREVIOUS STEP, USE 100% TRAINING DATA,

TURN ON SMALL WEIGHT DECAY (L_2, L_1)

\rightarrow FIND LR / $\downarrow L$ WITHIN ~ 100 ITERATIONS

\hookrightarrow ex. TRN LR = $10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}$

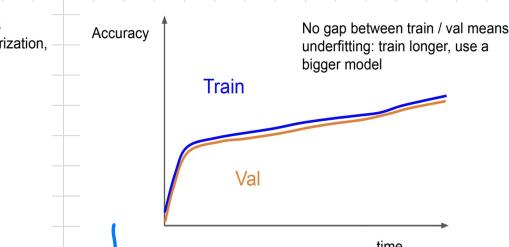
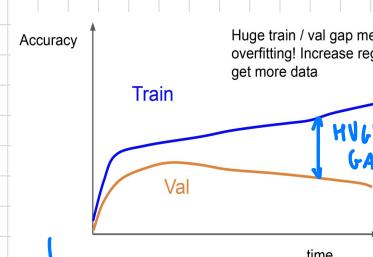
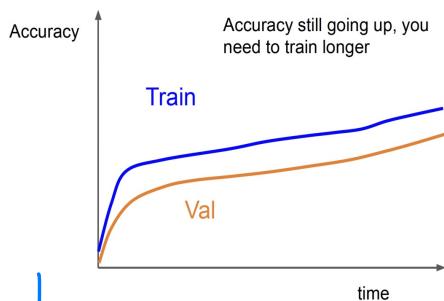
4. CHOOSE SOME VALUES OF LR AND WEIGHT DECAY THAT WORKED WELL

REPEAT AND CREATE A GRID, TRAIN FOR 1-5 EPOCHS

\hookrightarrow ex. WEIGHT DECAY: $10^{-4}, 10^{-5}, 0$

5. REDEFINING LR AND TRAIN FOR LONGER $\sim 10-20$ EPOCHS, NO WEIGHT DECAY

6. LOOK AT LOSS AND ACC CURVES



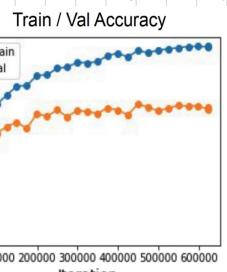
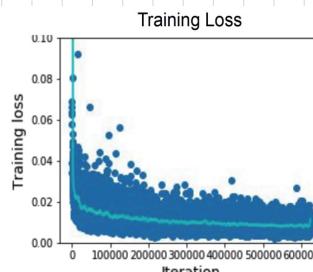
\hookrightarrow TODO: TRAIN LONGER

\hookrightarrow TODO: INCREASE REGULARIZATION GET MORE DATA

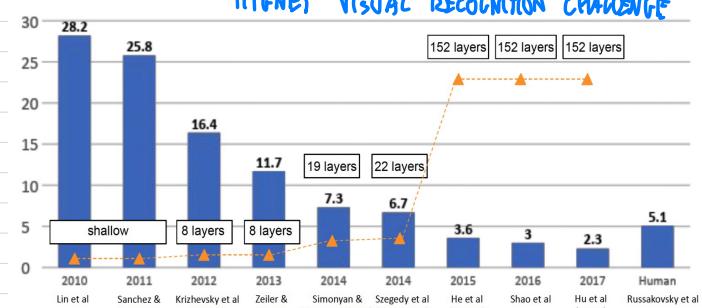
\hookrightarrow TODO: TRAIN LONGER USE BIGGER MODEL

7. GOTO STEP 5

GOOD CURVES \sim



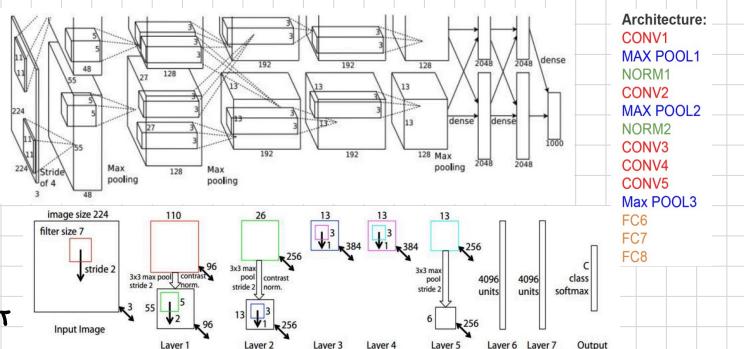
DEEP LEARNING ARCHITECTURES:



ALEXNET (2012):

1st CNN for images

with big improvements



ZF NET (2013):

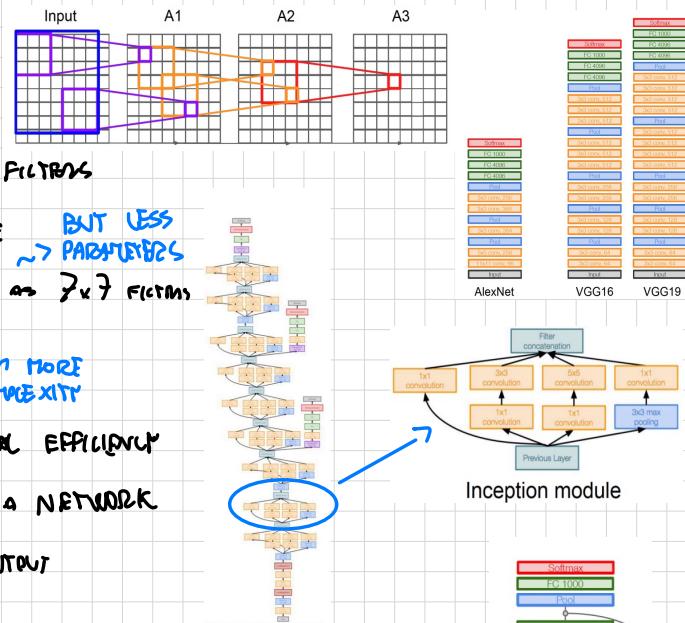
parameters improved over AlexNet

VGG-NET (2014):

smaller filters, deeper network

↳ 7×7 filters $\mapsto 3 \times 3$ filters

$\rightarrow 3 \times 3$ filters has same BUT LESS parameters
effective receptive field as 7×7 filters



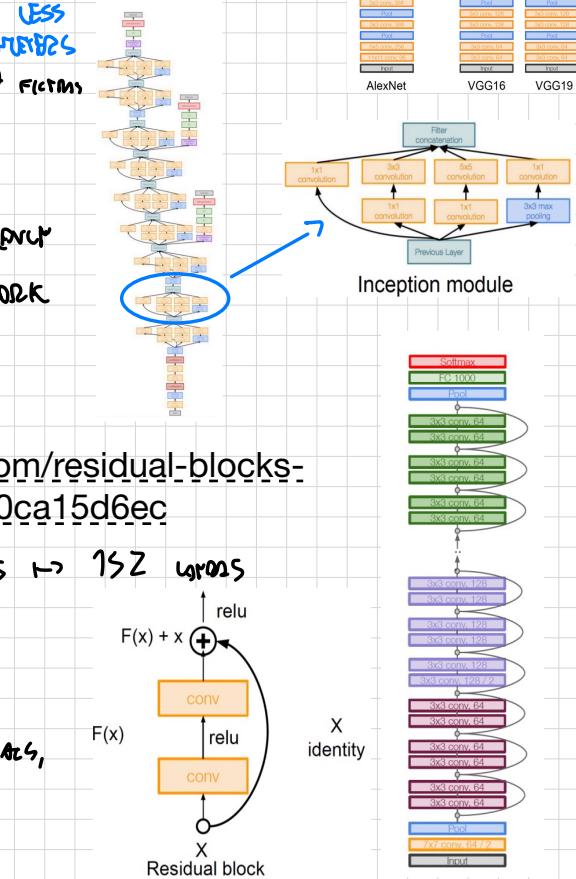
GOOGLE NET:

BUT MUCH MORE COMPLEXITY?
?

DEEPER NETWORK, MORE COMPUTATIONAL EFFICIENCY

INCEPTION MODULE: NETWORK IN A NETWORK

↳ FILTERS APPLIED IN PARALLEL, OUTPUT CONCATENATED CHANNEL-WISE



RESNET:

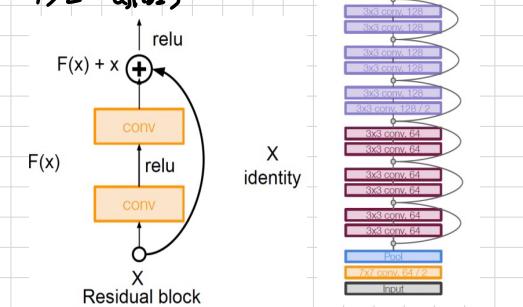
<https://towardsdatascience.com/residual-blocks-building-blocks-of-resnet-fd90ca15d6ec>

VG2N BIG JUMP IN DEPTH: 22 LAYERS \mapsto 152 layers

USE OF RESIDUAL CONNECTIONS;

Network tries to fit residual

MAPPING: LAYERS TRY TO FIT RESIDUALS, NOT THE OUTOTS



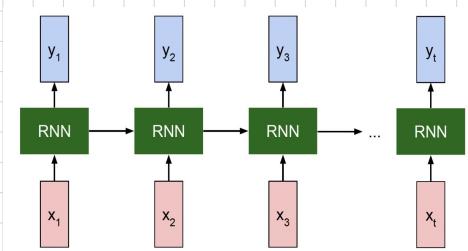
SENET:

SENSESE AN EXCITATION NETWORK;

ADDS A "FEATURE RECALIBRATION" MODULE: ADAPTIVELY RE-WEIGHTS MAPS

• RECURRENT NN (RNN):

RNNs have an internal state that is updated as a sequence is processed



$$h_t = f_W(h_{t-1}, x_t)$$

new state
 old state
 some function with parameters W
 input vector at some time step

$$y_t = f_{W_{hy}}(h_t)$$

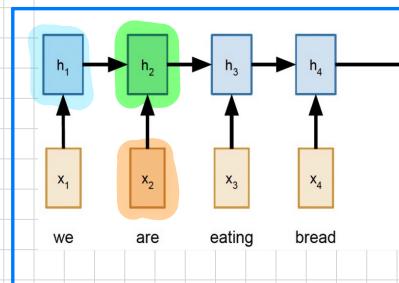
output
 new state
 another function with parameters W_o

TRANSFORMERS :

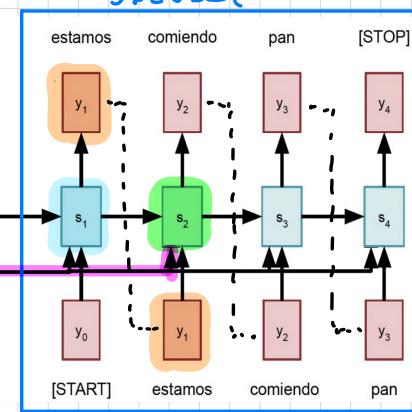
- RNN (RECALL) :

EN. RNN FOR ENG/ESP TRANSLATION

ENCODER



INITIAL DECODER STATE s_0
↑
CONTEXT VECTOR



- INPUT SEQUENCE: x_1, \dots, x_T / x_6 : GENERIC INPUT

- ENCODER:

$$\forall \text{ INPUT } x_t \rightarrow \text{STATE } h_t \text{ COMPUTED} : h_t = f_w(x_t, h_{t-1})$$

\rightarrow WHEN FINAL h_T COMPUTED \rightarrow OUTPUTS = $\begin{cases} s_0: \text{INITIAL DECODER STATE} \\ c: \text{CONTEXT VECTOR} \end{cases}$

/ OFTEN $c = h_T$, c HAS FIXED LENGTH \rightarrow BOTTLENECK

- DECODER:

OUTPUT SEQUENCE: y_1, \dots, y_T

$$\rightarrow y_t \text{ GENERATED FROM } s_t = y_v(s_{t-1}, c)$$

\hookrightarrow RNN PROBLEM: c HAS FIXED LENGTH

\hookrightarrow SOL: USE NEW CONTEXT VECTOR \forall STEP OF DECODE

\rightarrow A SINGLE c VECTOR REQUIRES TO SUMMARIZE WHICH WORDS ARE MORE IMPORTANT THAN OTHERS IN A (PIXEL) LENGTH REPRESENTATION

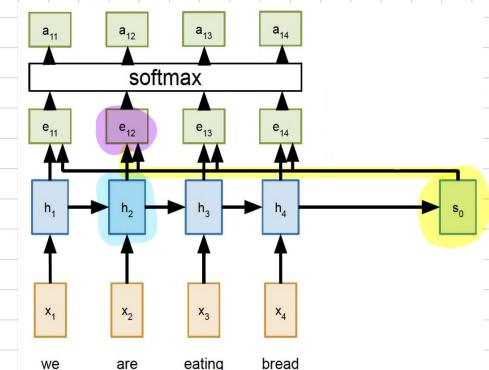
\rightarrow IT IS BETTER TO HAVE $\forall y_t \rightarrow c$; DIFFERENT SUMMARY OF INPUT SENTENCE $\forall y_t$

ATTENTION CONCEPT :

- ALIGNMENT SCORES: $e_{t,i} = f_{ATT}(s_{t-1}, h_i)$
 \uparrow scalar
 \hookrightarrow changes based on state too

ATTENTION WEIGHTS :

$$a_{t,i} = \text{softmax}(e_{t,i}) = \frac{\exp(e_{t,i})}{\sum_i \exp(e_{t,i})}$$



CONTEXT VECTOR:

$$c_t = \sum_i a_{t,i} \cdot h_i$$

- IT ATTENDS THE RELEVANT PART OF THE INPUT SEQUENCE

so -

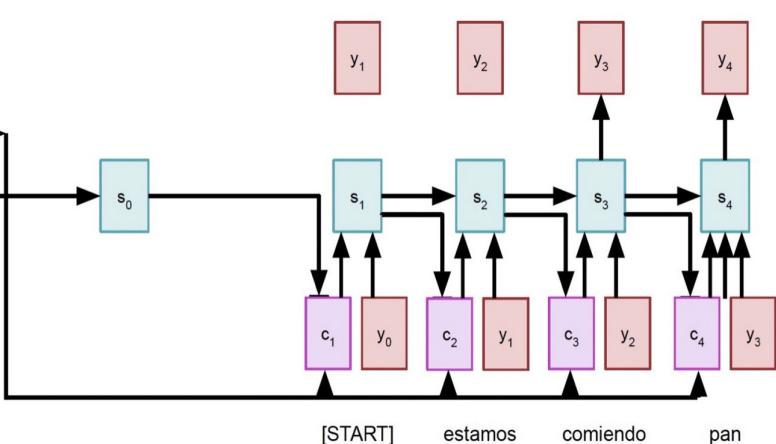
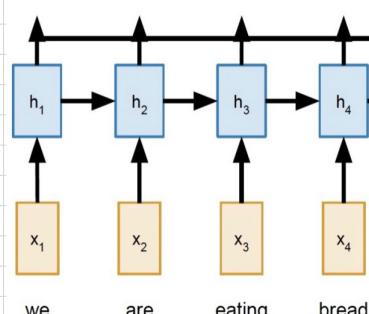
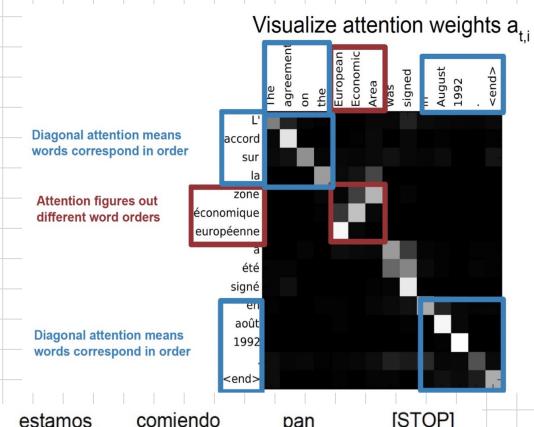
"ESTAMOS" = "WE ARE" \rightarrow IT IS POSSIBLE THAT :

$$\left\{ \begin{array}{l} \text{WE} \quad \text{ARE} \\ Q_{11} = Q_{12} = 0.45 \\ \text{EATING} \quad \text{BREAD} \\ Q_{13} = Q_{14} = 0.05 \end{array} \right.$$

Different from static RNN

DECODER :

$$s_t = y_u (N_{t-1}, s_{t-1}, c_t)$$



- IMAGE CAPTIONING WITH RNN AND ATTENTION : (SKIPPED, SEE SLIDES 26 - 37)

• TRANSFORMERS ARCHITECTURE:

IDEA: GET RID OF RNN, JUST USE ATTENTION CONCEPT \rightsquigarrow

"ATTENTION IS ALL YOU NEED", 2017

• GENERAL ATTENTION LAYER:

INPUT SEQUENCE = $\{x_0, x_1, x_2\}$

1. h : QUERY VECTOR, LAST DECODER STATE s_t IN TRANSLATION EXAMPLE

• $h.\text{shape} = D$ # same
• INPUT VECTORS $\bar{x}.\text{shape} = N \times D$ # features
 \downarrow $V_{\text{Word}} \rightarrow \text{Embedding}$

$$\rightarrow \text{ATTENTION: } e_i = \sum_{\text{ATT}} (h, x_i)$$

2. ATTENTION: $\bar{Q} = \text{softmax}(\bar{e}) \rightsquigarrow \text{NORMALIZATION}$

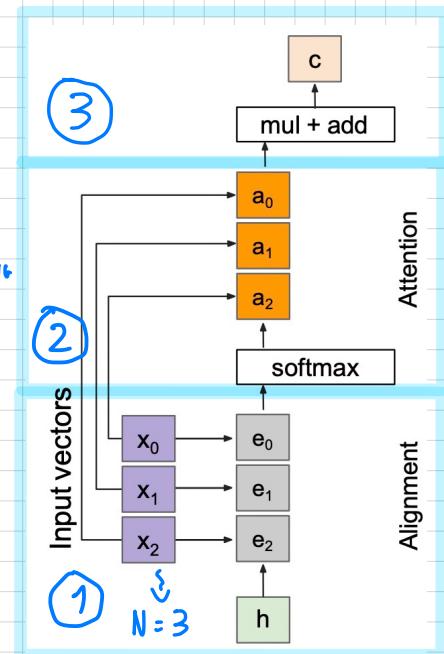
3. OUTPUT CONTEXT VECTOR: $\bar{C} = \sum_i Q_i \cdot x_i$

• $\bar{C}.\text{shape} = D$

• N.B.: BY NOW, ATTENTION MECHANISM: $h \mapsto C$ / C is almost a linear function of inputs

\hookrightarrow "ATTENTION IS ALL YOU NEED" PAPER PROPOSED SOME

CHANGES IN ORDER TO MAKE A FLEXIBLE NON-LINEAR MAPPING



"ATTENTION IS ALL YOU NEED" COMPUTATIONS:

ALIGNMENT:

$$e_i = f_{\text{ATT}}(h, x_i) \mapsto e_i = \frac{h \cdot x_i}{\sqrt{D}}$$

C. SHAPE,
FRAMES

- $f_{\text{ATT}}: \text{MLP} \mapsto \langle \dots, \dots \rangle$ DOT PRODUCT
 ↳ ONLY WORKS FOR WITH KEY & VALUE TRANSFORMATION TRICK

DIVIDING BY D REDUCES EFFECT OF LARGUE FRAMEWORK VECTORS

MULTIPLE QUERY VECTORS:

$$\bar{h} \text{ (shape } 0) \mapsto \bar{q} \text{ (shape } M \times 0)$$

$$\text{outputs: } \bar{c} \mapsto \bar{r} = \{r_0, r_1, r_2\}$$

N.B.:

IN THIS WAY: FROM $\bar{x} \mapsto c_i, q_i$

→ IT IS POSSIBLE TO ADD MORE EXPRESSIVITY TO THE INPUT
 ADDING FC LAYERS BEFORE ALIGNMENT AND ATTENTION:

$$\begin{aligned} \bar{x} &\xrightarrow{\text{FC}} \text{KEY VECTOR: } \bar{K} = \bar{x} \cdot W_K, \quad \bar{x} \xrightarrow{\text{FC}} \bar{K} \mapsto \text{ALIGNMENT } e_i; \\ &\xrightarrow{\text{FC}} \text{VALUE VECTOR: } \bar{V} = \bar{x} \cdot W_V, \quad \bar{x} \xrightarrow{\text{FC}} \bar{V} \mapsto \bar{V}(\bar{e}) \mapsto \text{ATTENTION } q_i \end{aligned}$$

→ USING KEY AND VALUES WE GET:

QUESTIONS

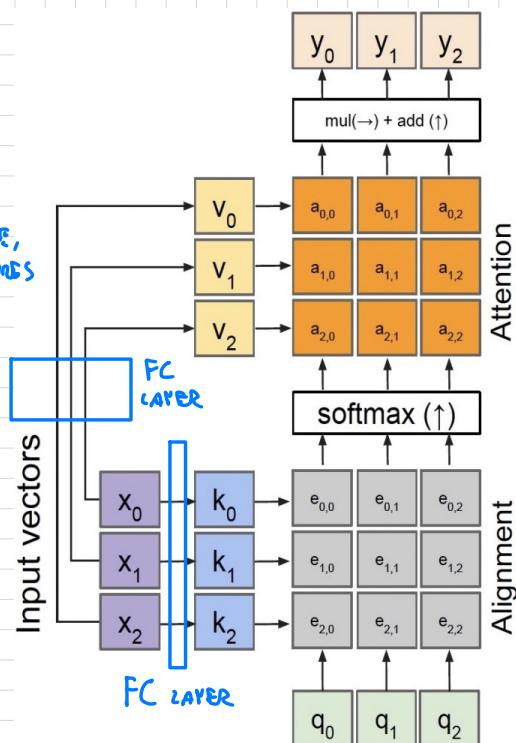
$$\cdot q. \text{shape} = M \times D \mapsto q. \text{shape} = M \times D_K$$

NOW INPUT AND OUTPUT DO NOT NEED TO HAVE SAME DIMENSIONS
 → THEY CHANGE BASED ON \bar{K} AND \bar{V}

$$\cdot \text{outputs: } V. \text{shape} = D \mapsto V. \text{shape} = M \times D_V$$

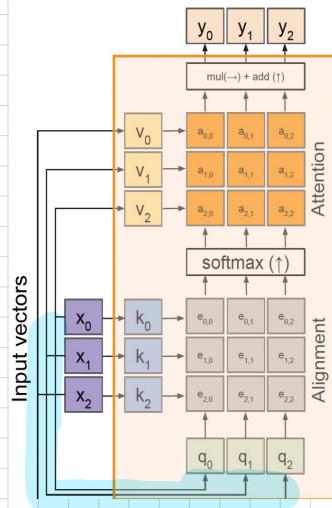
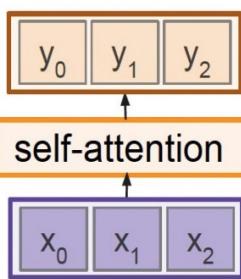
$$\hookrightarrow V_j = \sum_i q_{i,j} \cdot v_i$$

$$\cdot \text{ALIGNMENT: } e_i = \frac{h \cdot x_i}{\sqrt{D}} \mapsto e_i = \frac{q_i \cdot k_i}{\sqrt{D}}$$



• SELF ATTENTION LAYER:

\bar{q} is calculated directly from \bar{x}



• ATTENTION IN FORMULAS:

$$\text{ATTENTION : } A = \text{SOFTMAX} \left(\frac{(Q \cdot K^T)}{\sqrt{d_k}} \right)$$

→ HERE $d_k \approx D_k$

queries • Let's consider $D = 3 \rightarrow 3$ when \bar{q} , 3 key vector \bar{k} , 3 value vector \bar{v}

$$Q = \begin{pmatrix} q_{11} & q_{12} & q_{13} \\ q_{21} & q_{22} & q_{23} \\ q_{31} & q_{32} & q_{33} \end{pmatrix} = \begin{pmatrix} \bar{q}_1 \\ \bar{q}_2 \\ \bar{q}_3 \end{pmatrix} \quad K = \begin{pmatrix} k_{11} & k_{12} & k_{13} \\ k_{21} & k_{22} & k_{23} \\ k_{31} & k_{32} & k_{33} \end{pmatrix} = \begin{pmatrix} \bar{k}_1 \\ \bar{k}_2 \\ \bar{k}_3 \end{pmatrix} \quad V = \begin{pmatrix} v_{11} & v_{12} & v_{13} \\ v_{21} & v_{22} & v_{23} \\ v_{31} & v_{32} & v_{33} \end{pmatrix} = \begin{pmatrix} \bar{v}_1 \\ \bar{v}_2 \\ \bar{v}_3 \end{pmatrix}$$

$M \times D_K$ $M \times D_K$ $M \times D_V$

↑ # queries ↑ # queries ↑ # queries

• ALIGNMENT :

$$e_i = \bar{q}_i \cdot \bar{k}_i \quad / \text{in matrix: } E = Q \cdot K^T$$

rows of Q and K

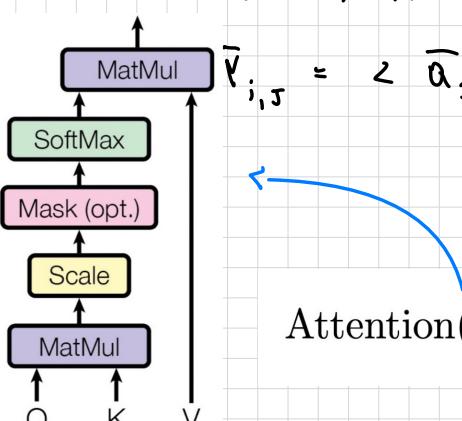
$$\hookrightarrow e_{i,j} = \langle \bar{q}_{i1}, \bar{k}_{j2} \rangle$$

$$QK^T = \begin{pmatrix} \bar{q}_1 \cdot \bar{k}_1 & \bar{q}_1 \cdot \bar{k}_2 & \bar{q}_1 \cdot \bar{k}_3 \\ \bar{q}_2 \cdot \bar{k}_1 & \bar{q}_2 \cdot \bar{k}_2 & \bar{q}_2 \cdot \bar{k}_3 \\ \bar{q}_3 \cdot \bar{k}_1 & \bar{q}_3 \cdot \bar{k}_2 & \bar{q}_3 \cdot \bar{k}_3 \end{pmatrix} = \begin{pmatrix} e_{11} & e_{12} & e_{13} \\ e_{21} & e_{22} & e_{23} \\ e_{31} & e_{32} & e_{33} \end{pmatrix} = \begin{pmatrix} \bar{e}_1 \\ \bar{e}_2 \\ \bar{e}_3 \end{pmatrix}$$

• ATTENTION :

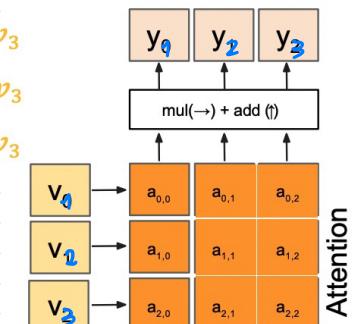
$$A : E = QK^T \longrightarrow A = \text{SOFTMAX} \left(\frac{E}{\sqrt{d_k}} \right) = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} = \begin{pmatrix} \bar{a}_1 \\ \bar{a}_2 \\ \bar{a}_3 \end{pmatrix}$$

• OUTPUTS VALUES:



$$\begin{aligned} \bar{y}_1 &= a_{11} \bar{v}_1 + a_{12} \bar{v}_2 + a_{13} \bar{v}_3 \\ \bar{y}_2 &= a_{21} \bar{v}_1 + a_{22} \bar{v}_2 + a_{23} \bar{v}_3 \\ \bar{y}_3 &= a_{31} \bar{v}_1 + a_{32} \bar{v}_2 + a_{33} \bar{v}_3 \end{aligned}$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



POSITIONAL ENCODING :

USED TO ENCODE ORDERED SEQUENCES, SINCE TRANSFORMERS ARE PERMUTATION EQUIVARIANT

- $\forall x_j \rightarrow \exists \text{ POSITION } p_j = \text{pos}(j)$

/ $\text{pos}: N \xrightarrow{\text{SCALAR}} \mathbb{R}^d \rightarrow \text{OUTPUT IS A } d\text{-DIMENSIONAL VECTOR}$

- A LOOK UP TABLE NEEDS TO BE LEARNT / LUT, SHAPE = $T \times D$

L → GIVEN AN INPUT TIME t , IT OUTPUTS THIS d -dimensional POSITION

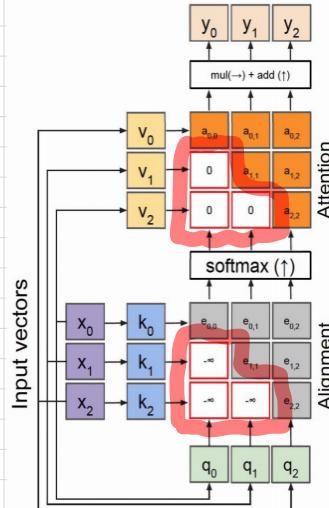
$$\rightarrow \text{pos}(t) = \begin{cases} \left[\begin{array}{c} \text{sim}(w_1, t), \text{sim}(w_2, t), \dots, \text{sim}(w_{d/2}, t) \\ \cos(w_1, t), \cos(w_2, t), \dots, \cos(w_{d/2}, t) \end{array} \right] & \text{if } d \text{ is even} \\ \left[\begin{array}{c} \text{sim}(w_1, t), \text{sim}(w_2, t), \dots, \text{sim}(w_{d/2}, t) \\ \cos(w_1, t), \cos(w_2, t), \dots, \cos(w_{d/2}, t) \\ w_k = \frac{1}{10000^{\frac{2k}{d}}} \end{array} \right] & \text{if } d \text{ is odd} \end{cases}$$

MASKED SELF-ATTENTION LAYER:

IT PREVENTS VECTORS TO LOOK AT FUTURE VECTORS

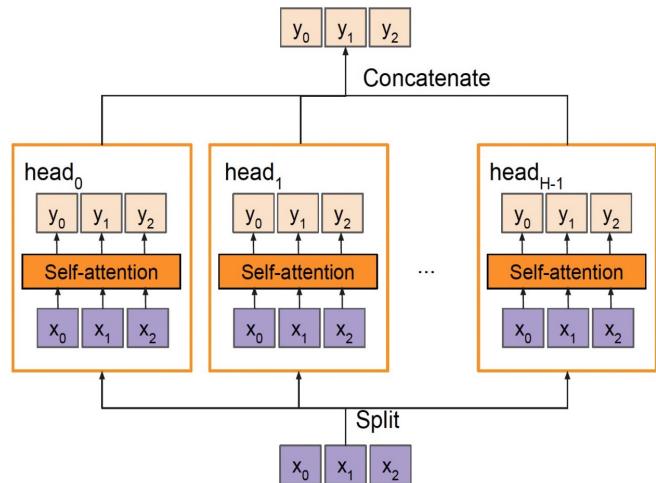
- SOME $e_{i,j} = -\infty$

$$\rightarrow a_{i,j} = 0$$

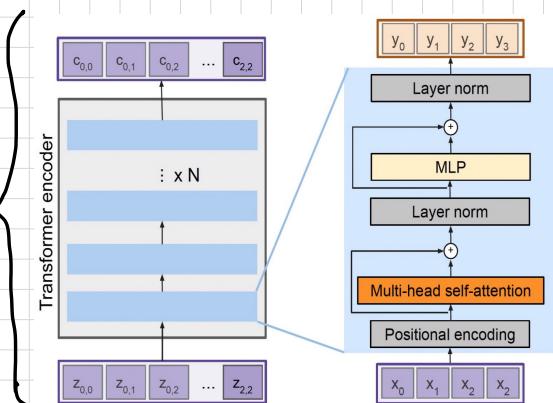


MULTI-HEAD SELF-ATTENTION LAYER:

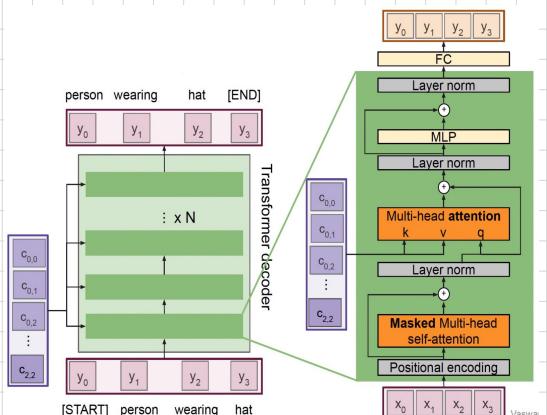
MULTIPLE SELF-ATTENTION HEADS IN PARALLEL



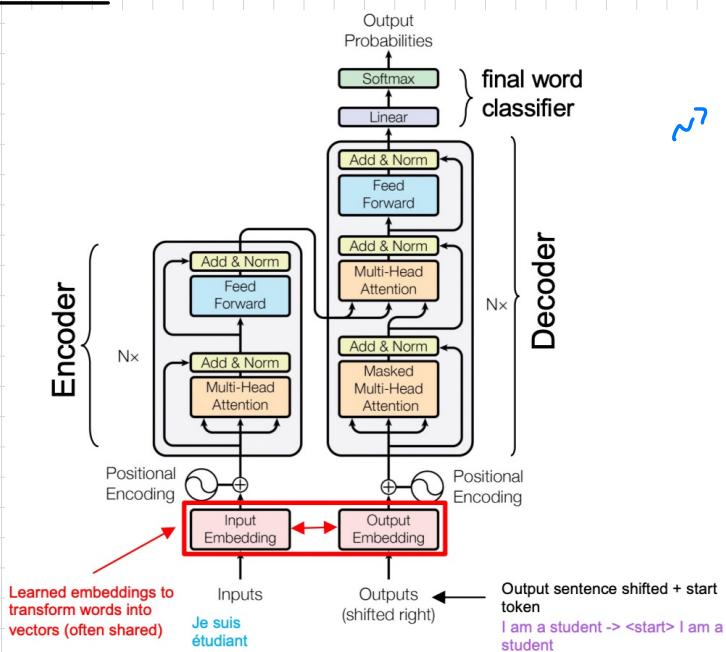
• TRANSFORMER ENCODER BLOCK K :



• TRANSPORTER DECODER BLOCK :



• FULL ARCHITECTURE :



• RNNs VS TRANSFORMERS :

RNNs

- (+) LSTMs work reasonably well for long sequences.
- (-) Expects an ordered sequences of inputs
- (-) Sequential computation: subsequent hidden states can only be computed after the previous ones are done.

Transformer:

- (+) Good at long sequences. Each attention calculation looks at all inputs.
- (+) Can operate over unordered sets or ordered sequences with positional encodings.
- (+) Parallel computation: All alignment and attention scores for all inputs can be done in parallel.
- (-) Requires a lot of memory: $N \times M$ alignment and attention scalers need to be calculated and stored for a single self-attention head. (but GPUs are getting bigger and better)

HELPFUL
GUIDE

<https://jalammar.github.io/illustrated-transformer/>