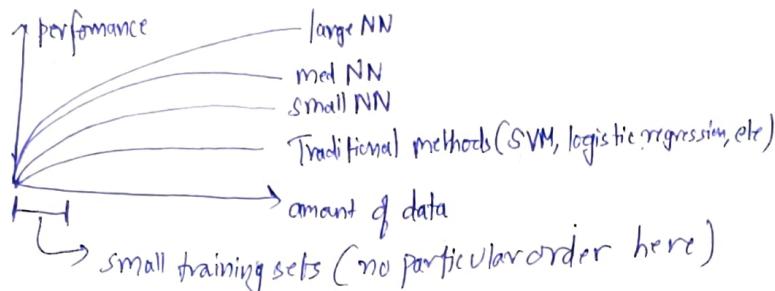


* ReLU \rightarrow Rectified Linear Unit



(15)

↳ Why deep learning taking off \rightarrow Scale drives deep learning progress



\rightarrow Logistic Regression as a Neural Network (week-2 / DHS-1)

↳ Binary classification (Notations)

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

↳ m training Examples $\Rightarrow \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

↳ $X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}^T \quad m_x \quad Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$

$$X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(m)} \end{bmatrix} \quad m_x \quad Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}] \quad Y \in \mathbb{R}^{1 \times m}$$

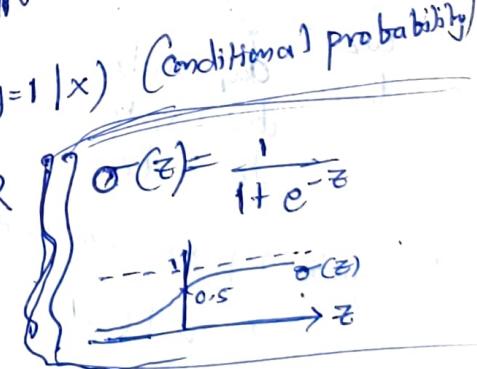
$m_x \Rightarrow \# \text{ of features}$
(e.g. for an image $64 \times 64 \times 3$)
 $\therefore m_x = 64 \times 64 \times 3$

↳ Logistic Regression

↳ given x (feature vector) want $\hat{y} = P(y=1|x)$ (Conditional probability)

↳ Parameters to be used: $w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$

↳ output: $\hat{y} = \sigma(w^T x + b)$



↳ Logistic regression cost function

* here we don't take squared error as loss function because it turns out to give a non convex optimization problem (multiple local minimas)

↳ loss function

$$\mathcal{L}(\hat{y}, y) = -[y \log(\hat{y}) + (1-y) \log(1-\hat{y})] \quad \xrightarrow{\text{Single training example}}$$

↳ Cost function

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

\rightarrow Cost of the parameters
↳ need to minimize it to get best parameters.

→ Gradient descent

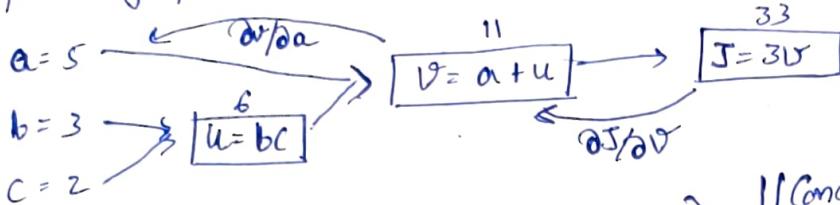
↳ want to find w, b that minimizes $J(w, b)$

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

(16)
Learning rate
represented as "dw"

↳ Computing Derivatives Using Computation graph

a computation graph is like this: (to compute $J = 3(a + bc)$)



→ forward propagation
← backward propagation
(to calculate derivative)

↳ Concept-1 $\Rightarrow \frac{\partial J}{\partial v} \approx \frac{\Delta J}{\Delta v} \approx \left(\begin{array}{l} \text{how much increase in } J \\ \text{(due to small increase in } v) \end{array} \right)$ || Concept-2: chain rule of derivation

↳ Notation: repeatedly need to calculate $d(J)$ wrt all the parameters $d(\text{var})$ of J is the final output || denote: $\frac{dJ}{d\text{var}} = d\text{var}$

× Backward Propagation for derivatives: \Rightarrow dvar of the rightmost variables are found first which then helps to find the dvar of left variable

→ Logistic Regression derivatives

$$\begin{aligned} & z = w_1x_1 + w_2x_2 + b \\ & \frac{dz}{dw_1} = x_1 \\ & \frac{dz}{dw_2} = x_2 \\ & \frac{db}{dz} = 1 \end{aligned}$$

$$a = \sigma(z) \rightarrow \mathcal{L}(a, y)$$

$$\frac{da}{dz} = \left[\frac{-y}{a} + \frac{1-y}{1-a} \right]$$

$$\begin{aligned} & w_1 := w_1 - \alpha dw_1 \\ & w_2 := w_2 - \alpha dw_2 \\ & b := b - \alpha db \end{aligned}$$

↳ Logistic regression on m examples: $\Rightarrow J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(a^{(i)}, y^{(i)})$

$$\frac{\partial J(w, b)}{\partial w_1} = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_1} \mathcal{L}(a^{(i)}, y^{(i)}) = \frac{1}{m} \sum_{i=1}^m dw_1^{(i)}$$

→ Logistic Regression Implementation

↳ a single step of gradient descent:

init: $\Rightarrow J=0; dw_1=0; dw_2=0; db=0$

for $i = 1$ to m

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

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

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

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

$$\begin{aligned} dw_1 &+= x_1^{(i)} dz^{(i)} \\ dw_2 &+= x_2^{(i)} dz^{(i)} \\ db &+= dz^{(i)} \end{aligned}$$

[End of for loop]

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

$$\begin{aligned} w_1 &:= w_1 - \alpha dw_1; w_2 := w_2 - \alpha dw_2 \\ b &:= b - \alpha db \end{aligned}$$

(17)

↳ Vectorization:> Calculations using matrix algebra can be 100x faster in computation than explicit for loops.

↳ Vectorizing Logistic Regression

$$\text{we want to compute: } z^{(1)} = w^T x^{(1)} + b \dots z^{(i)} = w^T x^{(i)} + b \dots z^{(m)} = w^T x^{(m)} + b$$

$$a^{(1)} = \sigma(z^{(1)}) \dots a^{(i)} = \sigma(z^{(i)}) \dots a^{(m)} = \sigma(z^{(m)})$$

$$\left\{ \begin{array}{l} X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ 1 & 1 & \dots & 1 \end{bmatrix}; \quad Z = \begin{bmatrix} z^{(1)}, z^{(2)}, \dots, z^{(m)} \end{bmatrix} \end{array} \right.$$

$$Z = w^T X + [b, b, \dots, b] \Rightarrow Z = \begin{bmatrix} w^T x^{(1)} + b, \dots, w^T x^{(m)} + b \end{bmatrix}$$

$$A = [a^{(1)} \ a^{(2)}, \dots, a^{(m)}] = \sigma(Z)$$

↳ Implementing Logistic Regression by vectorization:

$$Z = w^T X + b = np.dot(w^T, X) + b$$

$$A = \sigma(Z)$$

$$dZ = A \cdot Y$$

$$dw = \frac{1}{m} X dZ^T$$

$$db = \frac{1}{m} np.sum(dZ)$$

→ Python Broadcasting vectors

$$\left\{ \begin{array}{l} w := w - \alpha dw \\ b := b - \alpha db \end{array} \right.$$

} single iteration of gradient descent
(it can't be)
vectorized

↳ while loop
Condition based

$$\left[\begin{array}{c|c} \begin{bmatrix} (m,1) \\ \vdots \\ (1,m) \end{bmatrix} & \begin{bmatrix} (1,1) \\ \vdots \\ (1,n) \end{bmatrix} \\ \hline \begin{bmatrix} (m,n) \\ \vdots \\ (1,n) \end{bmatrix} & \begin{bmatrix} (\#_1, \#_2, \dots, \#_n) \\ \vdots \\ (\#_1, \#_2, \dots, \#_n) \end{bmatrix} \end{array} \right] + \begin{bmatrix} (\#_1, \#_2, \dots, \#_n) \\ \vdots \\ (\#_1, \#_2, \dots, \#_n) \end{bmatrix} = \begin{bmatrix} (\#_1, \#_2, \dots, \#_n) \\ \vdots \\ (\#_1, \#_2, \dots, \#_n) \end{bmatrix} + \begin{bmatrix} (\#_1, \#_2, \dots, \#_n) \\ \vdots \\ (\#_1, \#_2, \dots, \#_n) \end{bmatrix}$$

↳ Why such a cost function of logistic regression (Maximum likelihood)

$$\left\{ \begin{array}{l} \hat{y} = \sigma(w^T x + b) = P(y=1|x) \end{array} \right.$$

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

$$1-\hat{y} = P(y=0|x)$$

→ to maximize the probability $P(y|x)$ maximized $\log(P(y|x))$ monotonically increasing

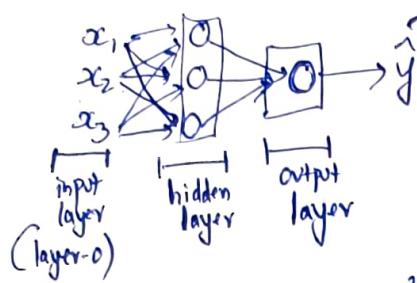
$$\text{so, } P(y|x) = (\hat{y})^y (1-\hat{y})^{1-y}$$

$$\text{for } m\text{-samples: } \text{IID} \Rightarrow P(y|x) = \prod_{i=1}^m P(y^{(i)}|x^{(i)})$$

$$\rightarrow \text{maximize } \log(P(y|x)) = \sum_{i=1}^m -\mathcal{L}(\hat{y}^{(i)}, y^{(i)}) \rightarrow \text{minimize loss function}$$

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

→ Neural Network Representation :



two layer NN (1 hidden layer)

output from layer 0 $\Rightarrow a^{[0]} = x$
 output from layer 1 $\Rightarrow a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \end{bmatrix}$
 output layer $\Rightarrow a^{[2]} = \hat{y}$

Parameter in each layers :-

layer-1 $\Rightarrow W^{[1]}, b^{[1]}$ || layer-2 $\Rightarrow W^{[2]}, b^{[2]}$

$$\begin{aligned} z_1^{[1]} &= w_1^{[1]T} x + b_1 \\ a_1^{[1]} &= \sigma(z_1^{[1]}) \\ z_2^{[1]} &= w_2^{[1]T} x + b_2 \\ a_2^{[1]} &= \sigma(z_2^{[1]}) \end{aligned}$$

and so on ..

and,

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

Given input x (1-hidden layer)

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

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

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

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

$$z^{[1]} = W^{[1]} X + b^{[1]}$$

$$A^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(z^{[2]})$$

for m training examples

$\text{Nt} \Rightarrow \text{super}[i] \rightarrow i^{\text{th}} \text{ layer}; \text{super}(i) \rightarrow i^{\text{th}} \text{ example}$

$$X^{(i)} \rightarrow a^{[2](i)} = \hat{y}^{(i)}; X^{(2)} \rightarrow a^{2} = \hat{y}^{(2)}; \dots; X^{(m)} \rightarrow a^{[2](m)} = \hat{y}^{(m)}$$

$$\text{take} \Rightarrow X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}; Z^{[1]} = \begin{bmatrix} z^{1} & z^{[1](2)} & \dots & z^{[1](m)} \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} a^{1} & a^{[1](2)} & \dots & a^{[1](m)} \end{bmatrix} \text{ and so on - .}$$

Activation functions :-

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

almost always $\tanh(z)$ better than $\sigma(z)$

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

as mean of $\tanh(z)$ centered around zero which helps for predictions in the next layers.

$\sigma(z)$ only used at output layer when $\hat{y} = \{0, 1\}$

Demerit of $\tanh(z)$

≈ 0 gradient for larger $|z|$ values which may slow down the gradient descent

Use ReLU (Rectified Linear Unit)

$$a = \max(0, z)$$



also sometimes used : Leaky ReLU

$$a = \max(0.01z, z)$$

0.01 can be made an hyperparameter

→ Why we need Activation function

If get rid of activation function: $a^{[1]} = z^{[1]} = W^{[0]}X + b^{[1]}$; $a^{[2]} = z^{[2]} = W^{[1]}a^{[1]} + b^{[2]}$

So, even if we would have a large number of hidden layers all we are doing here is just calculating a linear function at the end

only place to use only linear functions: output layer of a regression problem so that the output $\hat{y} \in (-\infty, \infty)$

→ Derivatives of Activation functions :-

$$\begin{aligned} & \text{Sigmoid } g(z) = \frac{1}{1+e^{-z}} \rightarrow \frac{d}{dz}g(z) = g(z)[1-g(z)] = g'(z) \\ & \text{tanh } g(z) = \tanh(z) \rightarrow g'(z) = 1 - (g(z))^2 \end{aligned}$$

$$\text{ReLU: } g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 0 & z < 0 \\ 1 & z \geq 0 \\ \text{n.d.} & z=0 \end{cases}$$

$$\text{Leaky ReLU: } g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01 & z < 0 \\ 1 & z \geq 0 \end{cases}$$

→ Gradient Descent for Neural Networks (One-hidden-layer)

Parameters: $W^{[0]}, b^{[0]}, W^{[1]}, b^{[1]}$ || $n_x = n^{[0]}, n^{[1]}, n^{[2]} = 1$
 $(n^{[0]}, n^{[0]})$ $(n^{[0]}, 1)$ $(n^{[0]}, n^{[1]})$ $(n^{[1]}, 1)$

Cost function: $J(W^{[0]}, b^{[0]}, W^{[1]}, b^{[1]}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}_i, y_i)$ || $\hat{y} = a^{[2]}$
 ↪ for binary classification loss function is same as logistic one

Gradient descent: Repeat { Compute predicts (\hat{y}_i , $i=1, \dots, m$) }

$$\frac{\partial J}{\partial W^{[1]}} = \frac{\partial J}{\partial a^{[2]}} \cdot \frac{\partial a^{[2]}}{\partial W^{[1]}} \quad \frac{\partial J}{\partial b^{[1]}} = \frac{\partial J}{\partial a^{[2]}} \cdot \frac{\partial a^{[2]}}{\partial b^{[1]}}$$

$$W^{[1]} := W^{[1]} - \alpha \frac{\partial J}{\partial W^{[1]}}, \quad b^{[1]} := b^{[1]} - \alpha \frac{\partial J}{\partial b^{[1]}}$$

formulas for computing derivatives

↓ Backpropagation: $d_z^{[2]} = A^{[2]} - Y$ || $Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$

forward

$$z^{[1]} = W^{[0]}X + b^{[1]}$$

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

$$z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$

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

$$= \sigma(z^{[2]})$$

(as, binary classification)

$$y \in \{0, 1\}$$

$$dW^{[2]} = \frac{1}{m} d_z^{[2]} (A^{[1]})^T$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(d_z^{[2]}), \quad \text{axis}=1, \quad \text{keepdims=True}$$

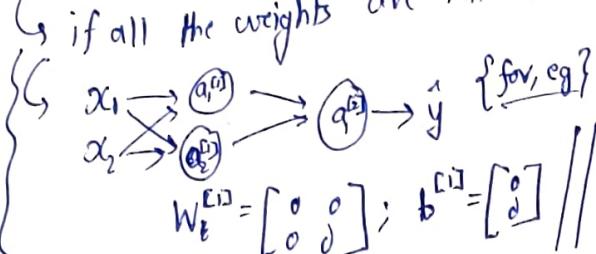
$$dz^{[1]} = W^{[2]T} d_z^{[2]} * g^{[2]'}(z^{[2]})$$

↑ element wise product

$$dW^{[1]} = \frac{1}{m} d_z^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(d_z^{[1]}, \text{axis}=1, \text{keepdims=True})$$

Initialization of weights

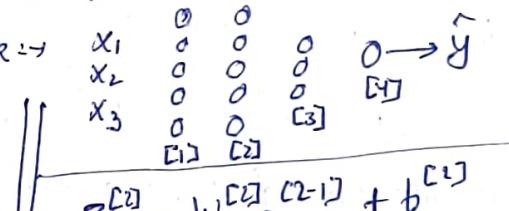
- ↳ if all the weights are initialized to zero \Rightarrow all the units of the hidden layer become identical due to symmetrical situation \Rightarrow they compute the exact same function
- ↳  {for eg?}
- $$W_1^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}; \quad b_1^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad || \quad W_2^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \Rightarrow a_1^{[1]} = a_2^{[1]} \Rightarrow dz_1^{[1]} = dz_2^{[1]} \Rightarrow dw = \begin{bmatrix} u & v \\ u & v \end{bmatrix}; \quad W^{[1]} = W^{[1]} - \alpha dw$$

↳ Hence initialize the parameters Randomly {ok to initialize b to zero}

- ↳ $W^{[1]} = np.random.randn(2, 2) * 0.01$ take a small multiplier because don't want to end-up near the edges of a $\sigma(z)$ or $tanh(z)$ function for faster convergence.
- $b^{[1]} = np.zeros(2)$
- and so on...
- $W^{[2]} = np.random.randn(1, 2) * 0.01 \quad || \quad b^{[2]} = 0$

↳ Deep Neural Networks \sim having more than one hidden layer.

$$\eta^{[2]} = \# \text{ of units in layer 2} \quad || \quad 2 = \# \text{ of layers in neural net}$$

* Forward propagation in a deep network \Rightarrow 

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

$a^{[1]} = g^{[1]}(z^{[1]}) \quad || \quad a^{[2]} = g^{[2]}(z^{[2]})$

and so on... $z^{[4]} = W^{[4]}a^{[3]} + b^{[4]}; \quad a^{[4]} = \hat{y} = g^{[4]}(z^{[4]})$

needs to be implemented using for loop only.

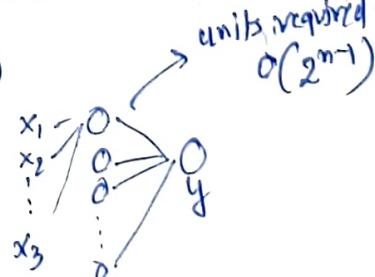
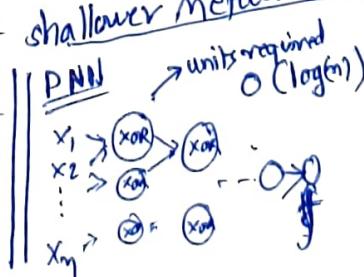
for m-training examples \Rightarrow $Z^{[l]} \in \left\{ \frac{1}{m} \sum_{i=1}^m z^{[l]}_i \right\} \quad Z^{[l]} = W^{[l]}A^{[l-1]} + b^{[l]} \quad \left\{ \begin{array}{l} A^{[0]} = \left[a^{[0]} \right] \\ A^{[1]} = \left[a^{[1]} \right] \dots \left[a^{[m]} \right] \\ A^{[2]} = X = \left[x^{[1]} \right] \dots \left[x^{[m]} \right] \end{array} \right.$

Intuition about Deep Neural Networks

- ↳ why we need the neural networks to be deep?
- ↳ because, the initial layers try to learn the simpler features in the dataset and then the deeper layers tries to compose these simpler features to learn complex features in the dataset.

* Circuit theory and deep learning \Rightarrow There are functions that you can compute with a "small" L-layered deep NN that shallower networks require exponentially more hidden units to compute.

$f = x_1(x_{OR})x_2(x_{OR})x_3 \dots (x_{OR})x_m \rightarrow$



Building blocks of DNN

Layer 2:

forward: Input $a^{[2-1]}$, output $a^{[2]}$

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

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

Cache $z^{[2]}$

backward: Input $da^{[2]}$ ($\text{Cache } z^{[2]}$) → output $da^{[2-1]}$ ($dW^{[2]}, db^{[2]}$)

Layer 2

$$a^{[2-1]} \rightarrow [W^{[2]}, b^{[2]}] \rightarrow a^{[2]}$$

Cache $\downarrow z^{[2]}$ (copy $W^{[2]}, b^{[2]}$)

$$da^{[2-1]} \leftarrow [W^{[2]}, b^{[2]}] \text{ Compute: } dZ^{[2]}$$

$$dW^{[2]}, db^{[2]}$$

$$dZ^{[2]} = da^{[2]} * g^{[2]}'(z^{[2]})$$

$$dW^{[2]} = dZ^{[2]} \cdot a^{[2-1]T}$$

$$db^{[2]} = dZ^{[2]}$$

$$da^{[2-1]} = W^{[2]T} dZ^{[2]}$$

$$dZ^{[2]} = dA^{[2]} * g^{[2]}'(z^{[2]})$$

$$dW^{[2]} = (\frac{1}{m}) dZ^{[2]} A^{[2-1]T}$$

$$db^{[2]} = (\frac{1}{m}) \text{np.sum}(dZ^{[2]}, \text{axis}=1, \text{keepdims=True})$$

$$dA^{[2-1]} = W^{[2]T} dZ^{[2]}$$

initialize backward prop with $da^{[2]} = \frac{d\mathcal{L}(g, y)}{dy} [g \neq y, g = a]$

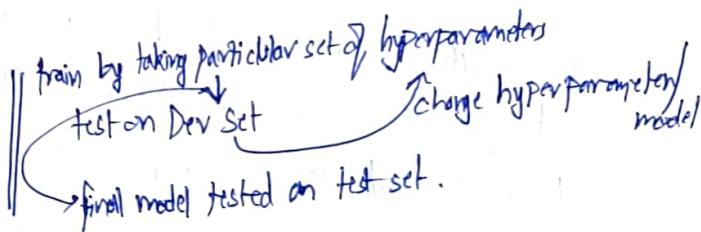
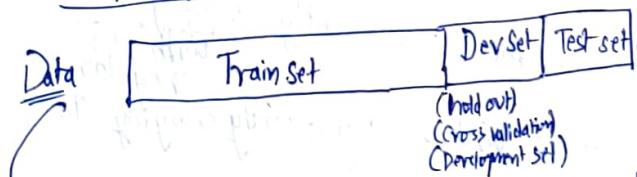
$$\text{eg for binary classification } da^{[2]} = -\frac{y}{a} + \frac{(1-y)}{(1-a)}$$

Hyperparameters: {learning rate α , # iterations, # hidden layers L , # hidden units $n^{[0]}, n^{[1]}, \dots$ }

Choice of activation function

These are tunable parameters which control the values of Parameters ($W^{[2]}, b^{[2]}$)

Test / Dev / Train Sets



for smaller dataset: good split $\rightarrow 70/30$ or $60/20/20$

for big data: dev set and test set \rightarrow going smaller.

Make sure that the dev and test sets come from the same distribution as train set

Not having a test set might be okay \rightarrow train/dev "test"

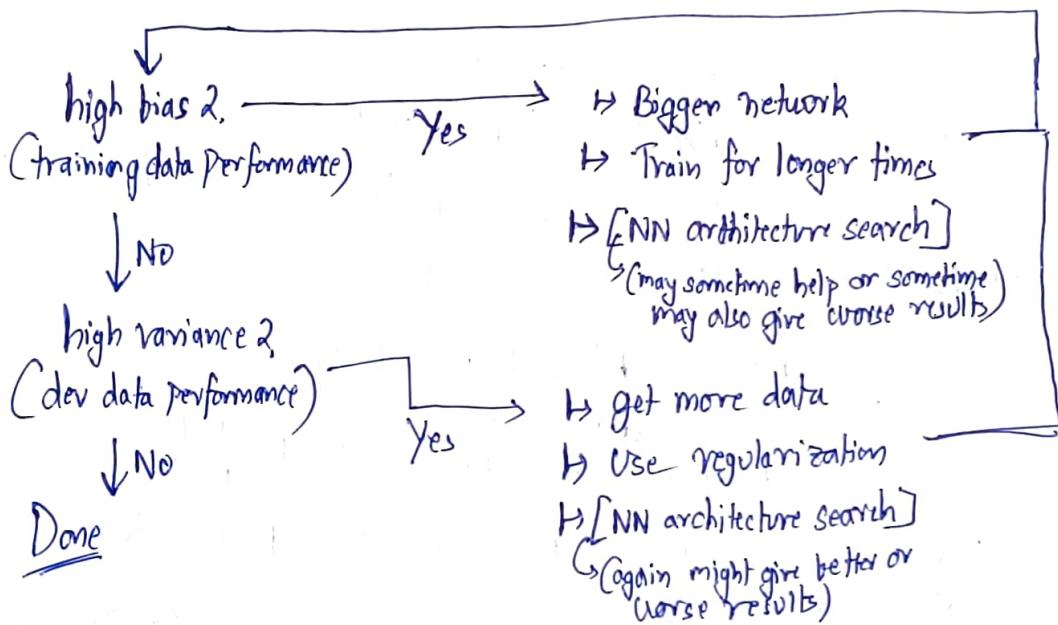
okay if you don't want a completely unbiased performance of the algorithm as this might be overfitting on dev set.

Bias and Variance in Deep learning

(22)

- ↳ bias \Rightarrow underfitting \Leftrightarrow oversimplification of model
 - ↳ Variance \Rightarrow overfitting \Leftrightarrow overcomplexity of model
- high train and dev set error \Rightarrow high bias
 low train and high dev set err \Rightarrow dev high variance
- Under the assumption \rightarrow optimal (Bayes) error $\approx 0\%$
- train and dev set from same distribution.
- (Some classifier, e.g. humans can perfectly classify it, i.e. for eg. noisy and ~~diff~~ undistinguishable data like blurry images)

Basic Recipe for Machine Learning



* Bias-variance tradeoff \Rightarrow earlier due to lack of tools, reducing one \Rightarrow increasing the other.

- ↳ but with availability of larger and larger data and ability to compute bigger and bigger networks \Rightarrow no longer a trade-off
- ↳ bigger network almost always reduces bias without significantly changing variance (with proper regularization)
- ↳ bigger dataset almost always reduces variance without significantly changing bias

↳ Regularization in logistic regression (Parameters: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$)

$$\text{Cost function modified} := J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} (\text{Reg-factor})$$

$$\hookrightarrow L_2 \text{ regularization} \Rightarrow \text{Reg-factor} = \|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$$

$$\hookrightarrow L_1 \text{ regularization} \Rightarrow \text{Reg-factor} = \|w\|_1 = \sum_{i=1}^{n_x} |w_i| \rightarrow w \text{ will be sparse here}$$

(dimensionality reduction
probabilistic
(see earlier))

$$\hookrightarrow \text{not considering } b \text{ as } \text{len}(b) = 1 \ll \text{len}(w)$$

→ Regularization in Neural Network

23

$$\hookrightarrow \text{Cost function} \Rightarrow J(w^{(0)}, b^{(0)}, \dots, w^{(L)}, b^{(L)}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + (\text{reg-term})$$

$$\hookrightarrow \ell_2\text{-norm known as } \underline{\text{frobenius norm}} \hookrightarrow \boxed{\|W^{[2]}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^{m^{[2-1]}} (W_{ij}^{[2]})^2}$$

$$\text{reg-term} = \frac{1}{2m} \sum_{i=1}^L \|w^{(i)}\|_F^2$$

Change in the update term in gradient descent:

$$dW^{[2]} = (\text{from backprop}) + \frac{\lambda}{m} W^{[2]} \rightarrow W^{[2]} := W^{[2]} - \alpha dW^{[2]}$$

↳ hence Frobenius regularization also called weight decay

↳ How does regularization help? \rightarrow penalty for larger $W^{(2)}$ values \Rightarrow $W^{(2)}$ values kept small
 \Rightarrow dimensionality reduction of the parameter space
 \Rightarrow reduced overfitting.

→ Dropout Regularization

↳ overview: choose randomly some neural nodes and drop them out \Rightarrow simpler NN.

Inverted Drop-out of let suppose for layer $L=3$, def keep-prob = 0.8. }

$d_3 = \text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1])$ { $a_3^{[3]} \in \mathbb{R}^{(n^3, m)}$ }
 ↘ (n^3, m) boolean matrix ↘ for each training example, different nodes
 shut-off.

$$a^{[3]} = a^{[z]} * d3$$

$a^{[3]}_1 = \left(a^{[3]}/\text{keep-prob}\right)$ because $z^{[4]} = W^{[4]}a^{[3]} + b^{[4]}$
 otherwise will get reduced by 0.2

When making predictions at test time

- When making predictions at test time
- no nodes are dropped out and \hat{y} calculated as usual using all nodes of all layers
- because random dropout of nodes will be incoherent and random outputs increasing noise in the predictions.

- because random dropout of noise in the predictions.
- also due to the "inverted" dropout step (third step) no need of any normalizing term while making predictions.

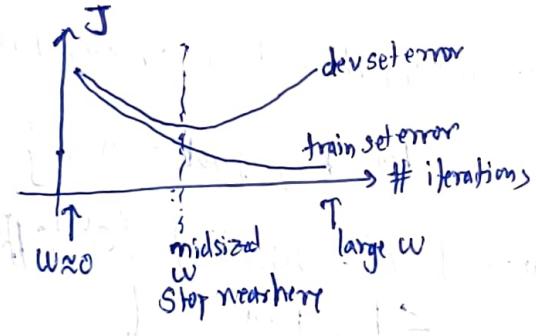
↳ Why does drop-out work? \Rightarrow Each node can't rely on any one pre-v-node / feature
 So have to spread out the weights \Rightarrow shrink weights \Rightarrow similar to L2 regularization.

↳ each layer can be given different "keep-probs" values based on the number of nodes in them. ↳ may be smaller layers given larger keep-probs ↳ but increases # hyperparameters.

↳ Downside of Dropout: ↳ cost function J is no longer well defined ↳ can't get smooth J vs # iterations curve to troubleshoot regularization ↳ $\begin{cases} J \\ \# \text{ of iterations} \end{cases}$

→ Data Augmentation: ↳ generating additional (pseudo) dataset from existing data to simulate having a large dataset ↳ eg: ↳ flipping cat pictures, adding some kind of noise and rotations in OCR data

Early stopping



One-downside: ↳ couples the two tasks of
① optimizing cost function J and to
② not overfitting too

Want to orthogonalize these tasks so as to better find the hyperparameters.

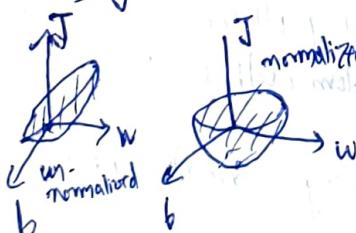
One-upside: ↳ less computation required compared to L_2 regularization where it is needed to train again and again for different λ

Normalizing training sets

$$X := \frac{x - \mu}{\sigma} \quad || \quad \mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}; \quad \sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$

↳ Note: ↳ use the same μ and σ for both training and test set

↳ Why normalize? ↳ data features otherwise have different ranges \Rightarrow W and b parameters have significantly different ranges \Rightarrow elongated cost-function ↳ longer time for gradient descent



Vanishing / Exploding gradients

$W^{[L]} > 1 \rightarrow$ very deep NN \rightarrow Exploding gradients
 $W^{[L]} < 1 \rightarrow$ very deep NN \rightarrow Vanishing gradients

$$\text{as } \hat{y} = W^{[L]} W^{[L-1]} \dots W^{[0]} x \Rightarrow \{\hat{y} \propto (W^{[0]})^L\}$$

(for e.g. for $g(z) = z$ and $b^{[0]} = 0$)

and all $W^{[0]}$ same

Similar scaling in other scenarios with very large deep Neural Nets.

→ One approach to reduce the problem of exploding/vanishing gradients

for a single neuron

$$z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b \xrightarrow{\text{ignore}} \hat{y} = g(z)$$

to keep gradients stable, large $n \rightarrow$ smaller w_i

$$\text{Var}(w_i) \propto \frac{1}{n}$$

So, initialize weights (for deeper nets) as

$$W^{[2]} = \text{np. random. random}(\text{shape}(W^{[2]}) \times \sqrt{\frac{2}{n^{[2-1]}}}) \xrightarrow{\text{Variance of a normal dist}}$$

$$\sqrt{\frac{2}{n^{[2-1]} + n^{[2]}}} \xrightarrow{\text{another helpful}}$$

$$\xrightarrow{\text{for ReLU}}$$

$$\sqrt{\frac{1}{n^{[2-1]}}} \xrightarrow{\text{for tanh}}$$

Xavier initialisation

→ Verifying backprop steps through gradient checking using numerical approximation

$$\hookrightarrow \text{two point formula to approx gradient: } f'(\theta) \approx \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \quad || \text{ error} = O(\epsilon^2)$$

first take $W^{[1]}, W^{[2]}, b^{[1]}, W^{[2]}, b^{[2]}$ and reshape it into a big vector θ

take $dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}$ and reshape to vector $d\theta$

$$\hookrightarrow \text{now, } d\theta_{\text{approx}}^{(i)} = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots)}{2\epsilon} \approx \frac{d\theta^{(i)}}{d\theta_i}$$

Check: $\epsilon = 10^{-7}$

$$\text{if } \frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2} \approx 10^{-7} \text{ — great!}$$

$$10^{-5} \text{ — double check}$$

$$10^{-3} \text{ — wrong.}$$

→ Grad-check only to be used in debug mode not in training mode

When regularization is on, take the ~~reg term~~ in your cost function definition.

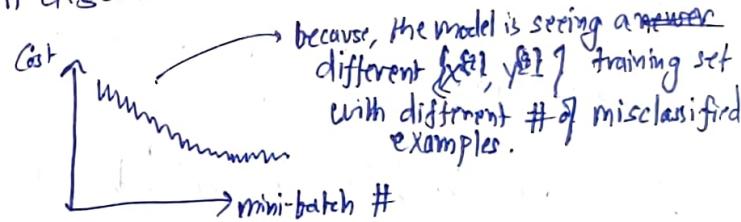
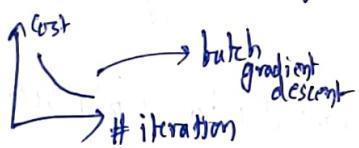
When regularization is off, switch-off dropout to check grads

Grad-check don't work with dropout so switch-off dropout to check grads

26

- Batch gradient descent : in which the whole training example is used at once in a single gradient descent step
- Mini-batch gradient descent : the whole training set is divided into mini-batches to be each used in a single gradient descent step.
- When # training example is very large \Rightarrow long time for batch gradient descent
- \Rightarrow mini-batch found to be a faster approach
- $X = [x^{(1)}, x^{(2)}, \dots, x^{(m)} | x^{(m+1)}, \dots, x^{(2m)}]$
- minibatch $t \Rightarrow x^{t1}, x^{t2}, \dots, x^{tT}, y^{tT}$
- After doing forward prop and backprop on each minibatch t , the a single step of gradient descent is taken.
- Epoch : After the whole training dataset has passed through forward-backprop once
- \hookrightarrow single step of grad descent in batch grad descent whereas multiple steps for minibatch

- Cost function with mini-batch gradient descent :



- Choosing the minibatch size

stochastic grad descent

lose the speed gained from vectorization
the movement towards minima highly noisy and never converge

minibatch

faster as vectorized make progress without waiting for whole training set to be processed

batch

too long per iteration

Size = $m \Rightarrow$ batch grad descent
Size = 1 \Rightarrow stochastic gradient descent

Tips :

- if small training set \Rightarrow batch gradient descent ($m \leq 2000$)
- if bigger then typical mini-batch sizes \Rightarrow 64, 128, 256, 512
- usually in powers of 2 as it is faster b/c the memory is there in batch

make sure minibatch size \leq CPU/GPU memory

- Exponentially Weighted Moving Averages

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

here V_t is approximately averaging over the θ values of $(1-\beta)$ previous steps

$\hookrightarrow V_t = (1-\beta) \theta_t + (1-\beta)\beta \theta_{t-1} + (\beta^2) \theta_{t-2} + (\beta^3) \theta_{t-3} + \dots$

\hookrightarrow "exponentially decaying weights"

27

$V_\theta := \beta V + (1-\beta)\Theta$

Bias Correction \rightarrow In exponentially weighted moving average, we have a problem in getting the proper values of initial V_t as they are given very small contribution from data points as there are less number of previous points.

$$\rightarrow V_0 = 0 \parallel V_1 = \underbrace{(1-\beta)\theta_1}_{\text{small value}} \parallel V_2 = \underbrace{\beta(1-\alpha)\theta_1 + (1-\beta)\theta_2} \parallel \dots$$

hence add a bias component weight. E.g. -

ence add a bias correcting weight: $\left[\frac{V_t}{1 - A_t^2} \right]$

→ Disadvantage of gradient descent:



→ too large a learning rate \Rightarrow the gradient descent oscillations will diverge in a zig-zag direction

- too large a learning rate will diverge
- too small a learning rate: although will go in the right direction (overall) but due to many oscillatory steps will take a large number of steps

Gradient descent with momentum

On iteration 7 of gradient descent:

dW, db from batch / mini-batch

$$\begin{cases} V_{dw} := \beta \cdot V_{dw} + (1-\beta) dW \\ V_{db} := \beta \cdot V_{db} + (1-\beta) db \end{cases}$$

$$V_{db} := \beta V_{ab} + (1-\beta) d_b$$

$$\begin{aligned} w &:= w - \alpha V_{dw} \\ b &:= b - \alpha V_{db} \end{aligned}$$

using the exponentially weighted avg of gradient helps in reducing the grads in the oscillating direction and enforces the grads towards the global minima :-

$\nearrow \nearrow \nearrow \nearrow \nearrow$ I grads gets cancelled out
grads get enforced in this dir by avg
this dir by avg

∇ Almost always works better than simple grad descent
This dir by avg

\rightarrow hyperparameter β := most usual and best value = 0.9 (avg over 10 grids)

RMSProp

Iteration of grad descent:

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dW^2$$

$$S_{db} := \beta_2 S_{db} + (1-\beta) db^2$$

$$W := W - \alpha \frac{dW}{\sqrt{S_{dW}} + \epsilon} \quad \left| \begin{array}{l} b: b - \alpha \frac{db}{\sqrt{S_{db}}} \end{array} \right.$$

Logic: → The parameter directions along which the oscillations are happening could have large RMS gradient(s) and hence the parameters in that direction would be updated by a smaller factor → rate of factor convergence.

→ can take larger learning rate for faster convergence

→ to avoid division by zero error $\epsilon \approx 10^{-8}$ used

(28)

→ Adam Optimization Algorithm (RMSprop + momentum)

Adam → Adaptive moment estimation

$$\hookrightarrow \text{initialize} : \quad V_{dw} = S_{dw} = V_{db} = S_{db} = 0$$

iteration $t \Rightarrow$ get dW, db using mini-batch

$$\hookrightarrow V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dW \quad || \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) db$$

$$\hookrightarrow S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dW^2 \quad || \quad S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2$$

$$\hookrightarrow V_{dw}^{\text{corrected}} := V_{dw}/(1-\beta_1^t) \quad || \quad V_{db}^{\text{corrected}} := V_{db}/(1-\beta_1^t) \quad || \quad S_{dw}^{\text{corrected}} = \frac{S_{dw}}{(1-\beta_2^t)} \quad || \quad S_{db}^{\text{corrected}} = \frac{S_{db}}{(1-\beta_2^t)}$$

$$\hookrightarrow W := W - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}}} + \epsilon} \quad || \quad b := b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}}} + \epsilon}$$

$$\rightarrow \text{Hyperparameter choice} : \quad \alpha: \text{need to tune} \quad || \quad \beta_1 = 0.9 \quad || \quad \beta_2 = 0.999 \quad || \quad \epsilon = 10^{-8} \quad \left\{ \begin{array}{l} \text{Recommended} \end{array} \right.$$

→ Learning Rate Decay → learning rate is reduced with increasing # iterations because as we move closer to the global minima, we want to oscillate in a smaller and smaller neighbourhood of the minima.

$$\hookrightarrow \text{decay-rate} : \quad \alpha = \frac{\alpha_0}{1 + (\text{decay-rate}) * (\text{epoch-num})} \quad || \quad \text{hyper-parameters} : \quad \alpha_0, \text{decay-rate}$$

$$\begin{cases} \alpha \downarrow \\ \text{discrete decay} \end{cases} \quad \rightarrow \quad \alpha = (0.95)^{\text{epoch}\#} \alpha_0 \quad \text{- exponentially decay} \quad \left\{ \begin{array}{l} \text{large training models} \\ \text{sometimes have manual} \\ \text{changes in decay rate} \end{array} \right. \quad \left\{ \begin{array}{l} \text{taking days to complete} \end{array} \right.$$

→ The problem of local optima : is very-very rare to get a local optima in which all the directions have inflection.

$$\hookrightarrow \text{Probability}(\alpha \geq 2^{-\# \text{ of parameters}}) \quad \rightarrow \quad \text{usually get saddle points}$$

→ The problem of plateaux → can make learning slow as gradient at those places are nearly zero.
algorithms like Adam consider- ∇P the rate.

→ hyperparameters tuning preferences :-

← highly likely to tune

highly unlikely to tune →

(29)

$\alpha \xrightarrow{\text{learning rate}} \beta, \# \text{hidden units}, \text{mini-batch size}, \xrightarrow{\text{# layers, learning rate decay}} \text{Adam} \beta_1, \beta_2, \epsilon$

* Trying out hyperparameters values : select points in the hyperparameter space in random

→ don't select hyperparameters using grid search

* Try to implement Coarse to fine

approach : → Search for hps in a broad area first and then search densely in a promising smaller area.

because, random way of using hyperparameter search gives the advantage of having an unique value for each hyperparameter in each trial run.

↳ hyperparameter search scales : → Uniform random search in linear scale : → # layers, # hidden units uniform random search in log scale : → α, β

→ Two approaches for hp tuning → baby sitting one model : → When low computational resources

↳ Training many models in parallel : → get the hps which have best and fastest convergence.

→ Keep on looking the performance of ongoing model training and update the hyperparameters accordingly.

→ Normalising inputs to speed-up learning

↳ just like in logistic regression, normalising inputs helps to train w, b faster,

↳ in DL for layer 2 : → normalizing $z^{[2]}$ can help to train $w^{[2]}, b^{[2]}$ faster

Batch Normalization

for $z^{[2]} : \mu = \frac{1}{m} \sum z_i^{(i)} // \sigma^2 = \frac{1}{m} \sum (z_i^{(i)} - \mu)^2$

here we use different γ and β which may or may not be 1 and 0

as we may want to have the mean and std of the dataset to be something more than just the middle portion of the activation function

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$z^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

{ γ, β are learnable parameters }
{ where γ is new std and β is new mean }

use $\tilde{z}^{[2](i)}$ instead of $z^{[2](i)}$

Adds additional set of parameters to the model : -

$\beta^{[2]}, \gamma^{[2]}, f^{[2]}, \gamma^{[2]}, \dots, \beta^{[L]}, \gamma^{[L]}$

$\{\beta^{[2]} = \beta^{[2]} - \alpha d\beta^{[2]}\}$ → parameter update

but we no longer have the set of parameters $b^{[2]}$

as they get cancelled out while taking mean of $z^{[2]}$ to 0 and then to $\beta^{[2]}$ as $f^{[2]} = w^{[2]} z^{[2]} + b^{[2]}$

Batch normalization implemented usually with mini-batches

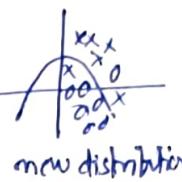
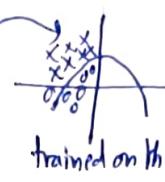
→ Covariate shift: → related to the need to retrain the model if the distribution of the inputs (x) changes even though the ground truth function is same. (30)

Why this a problem
in NN

Cg-1

want to identify white-cats

when the model to identify cats
are trained on darker-colored cats.



new distribution

→ the input distribution

of the m^{th} hidden layer
gets Covariate shifted when
 W, b of earlier layers are
changed during training

→ batch normalization (BN)

helps in stabilizing the
inputs to the m^{th} hidden
layer by keeping their
distribution almost the same

→ This helps in weakly,
de-coupling the hidden
layers' training as
each hidden layer has
its own input distribution
which is full-proof

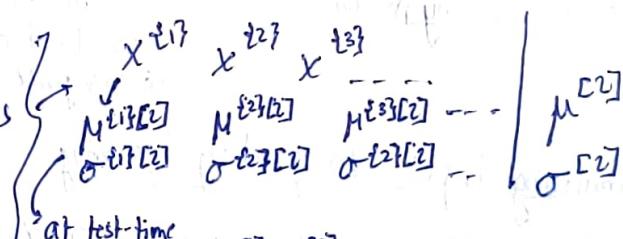
→ Batch Norm at test time

→ BN uses μ and σ^2 of the given
minibatch to normalize.

This weak
decoupling helps in
getting faster training

but how to get those values at test time?

→ we use exponentially weighted averages
over mini-batches.



$$z_{\text{norm}}^{[2]} = \frac{z^{[2]} - \mu^{[2]}}{\sqrt{\sigma^{[2]}} + \epsilon} \rightarrow z_{\text{norm}}^{[2]} = \gamma z_{\text{norm}}^{[1]} + \beta$$

→ Multiclass-Classification in Neural Networks

→ $C = \# \text{ of classes} \Rightarrow n^{[L]} = \# \text{ nodes in output layer} = C$ // for e.g. $a^{[L]} = \hat{y} \in \mathbb{R}^{C \times 1}$

$$\begin{aligned} \hat{y} &= P(C=0|x) \\ &= P(C=1|x) \\ &= P(C=2|x) \end{aligned}$$

→ uses a softmax activation function in
the final softmax layer

$$z^{[L]} = W^{[L]} a^{[L-1]} + b^{[L]} \quad (C, 1)$$

$$\text{softmax activation } \hat{y} = e^{(z^{[L]})} \quad (C, 1)$$

$$a^{[L]} = \frac{e^{z^{[L]}}}{\sum_i e^{z^{[L]}}} \quad (C, 1)$$

$$\text{for e.g. } z^{[L]} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \Rightarrow a^{[L]} = \frac{e^{z^{[L]}}}{\sum_i e^{z^{[L]}}} = \begin{bmatrix} e^{z^{[1]}} \\ e^{z^{[2]}} \\ e^{z^{[3]}} \end{bmatrix} / \begin{bmatrix} e^{z^{[1]}} + e^{z^{[2]}} + e^{z^{[3]}} \end{bmatrix}$$

* The softmax classifier reduces to
logistic regression for the case $C=2$.

$$\hookrightarrow \text{Loss function: } \mathcal{L}(\hat{y}, y) = - \sum_{j=1}^C y_j \log(\hat{y}_j) \quad || \text{ Cost function } J(W^{[L]}, b^{[L]}, \dots, W^{[1]}, b^{[1]}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(y^{(i)}, \hat{y}^{(i)})$$

$$\hookrightarrow \text{Backprop: initialization: } dz^{[L]} = \hat{y} - y \quad (C, 1)$$

Orthogonalization of controls in ML

→ we want to orthogonalize the cause and effect of tuning the various hyperparameters/strategies so as to narrow down our search for better models

Chain of Assumptions in ML

- ① fit training set well on cost function: → tuned by → bigger network, changing optimizer, etc
 - ② fit dev set well on cost function: → tuned by → regularization, bigger train set
 - ③ fit dev/test set well on cost function: → tuned by → bigger dev set
 - ④ performs well in real world: → tuned by → change dev set / change cost function.
- We tend to avoid strategies like early stopping which can tune more than one assumption above (e.g. early stopping tunes simultaneously both ① and ②)

* It is always useful to have a single number evaluation metric which would help to objectively differentiate between different models.
(e.g. use F1 score instead of using both of the precision and recall.)

Optimising and Satisficing metric

the metric which we want to maximize

the metric which we just want to be above/below certain values

e.g. maximizing Accuracy keeping running time < 100ms

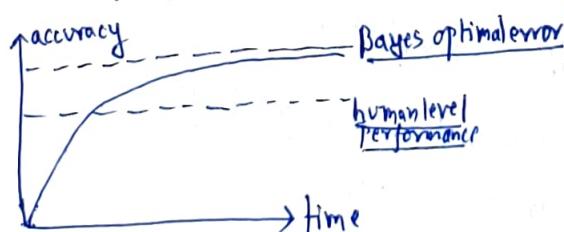
When we care about N number of metrics then

take one (most-)important one as an optimizing metric
all other N-1 as satisficing metric

* Choose a dev set and test set to reflect data you expect to get in the future and consider important to do well on

* choosing the dev set and the evaluation metric determines what target you want to achieve

→ Comparing to human-level performance:

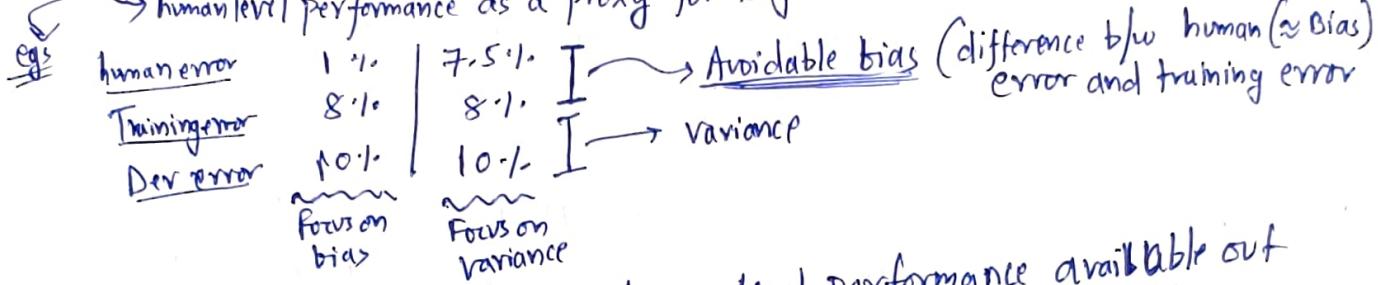


→ Bayes optimal error: → maximum accuracy theoretically any function can achieve from mapping $X \rightarrow Y$ limited by noise in the input

Why compare with human level performance?

- Because, so long as ML is performing worse than Human, the model accuracy can be tried to be improved by:
 - getting labeled data from humans
 - better analysis of bias/variance
- gain insights from manual error analysis: why person got it right? and machine got it wrong?

→ human level performance as a proxy for Bayes error



- You should always consider the best human level performance available out there as a proxy to decide upon Bayes error and whether to reduce bias or variance. (choose the lower error) from team of doctors vs higher prior for one doctor as human level proxy (\approx Bayes error) → (example)

Carrying out Error analysis

- Check the mislabelled data → Categorise the mislabelled based on what went wrong. (for eg: blurry image, mislabelling on a special kind of images, etc)
 - try to find-out
 - the major issues
 - (which issue causing the most mislabelling)
- Strategise to figure-out those major issues first.

* Deep-learning Algorithms are quite robust to random errors in the training set. But they are not robust to systematic errors.

* If there is a significant amt of incorrectly-labelled data in the dev set which is causing a significant % of total mislabelled cases by the model, then it is worthwhile to go and correct out the incorrectly-labelled dev data (by the humans).

* Build your first system quickly and then iterate:

- Set up dev/test set and metric → build initial system quickly
- use bias/variance analysis and error analysis to prioritize next steps.

→ Scenario : Different training and testing distribution

↳ Suppose you have a large dataset from one distribution (e.g. cub scrap) which is different from the distribution where the model is going to be used (e.g. user uploaded images). Two approaches here → (very little data from the target distribution)

- ① shuffle the data from both large dataset and target dataset and make the train/dev/test split (all from same distribution)
 ↳ Problem: → the dev/test distribution will have a very small fraction of the target dataset and hence model will get more optimised for the large dataset.

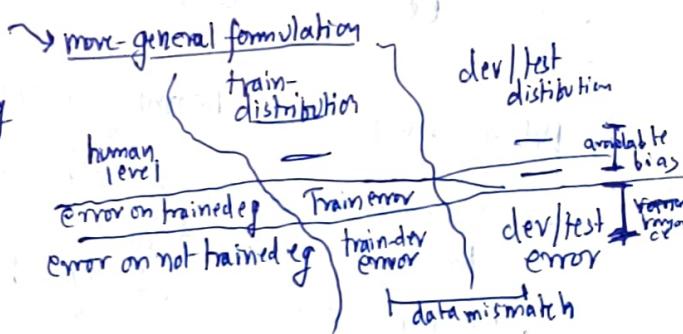
- ② make the test/dev from the target distribution and keep the remaining part of target distribution along with the large one as the training set
 ↳ Train distribution different than test/dev
 ↳ but optimization more towards target distribution

↳ Bias, variance on different training and testing distribution

↳ Train-dev-set: same distribution as training set, but not used for training:

Some distribution	Bayes error	available bias
	Training error	variance
	Train-dev error	data mismatch
Some distribution	dev error	degree of overfitting
	test error	to the dev set

* Note: → when the dev/test distribution is "easier" then it might sometimes happen that train error > dev error.



→ Addressing data Mismatch

↳ carry out manual error analysis to try to understand difference b/w training and dev/test sets

↳ make training data more similar, or collect more data similar to dev/test sets

Artificial data synthesis

↳ eg: → adding noise to data

↳ eg: → generating images using computer graphics

↳ Risk: → might overfit to noise/graphical objects.

↳ Transfer learning

- ↳ taking a model trained on one dataset to use it for a different application with fewer data. e.g. → model trained on image recognition transferred to learn radiology diagnosis.
- ↳ Makes sense only when, the ^{transfer} learning is happening from a model with a large amount of data to an application with limited amount of data.
As the initial basic layers would stay similar in transfer learning, priority is given to train the final layer(s) (retraining the $W^{[L]} b^{[L]}$ or adding more layers at the end, etc.) based on the amount of data available for finetuning
- ↳ qA) Pre-trained model → transfer learning → Fine-tuning f.B?
e.g.: retraining output layer
adding additional layer at the end
(Can large dataset)
- ↳ Task A and B must have same input
- ↳ low level features from A should be helpful for learning B

→ Multi-task learning :> using the same NN to do multiple things
(e.g. classify multiple different things) → presence of car, traffic light, sign, etc

This case is different from softmax regression → as one image has multiple labels here

When multi-task learning makes sense

- ↳ Training on a set of tasks that could benefit from having shared lower-level features
- ↳ usually: amount of data you have for each task is quite similar
- ↳ can train a big enough neural network to do well on all the tasks

→ End-to-End Deep Learning

- ↳ Pros: → let the data speak
- ↳ less hand-designing of components needed

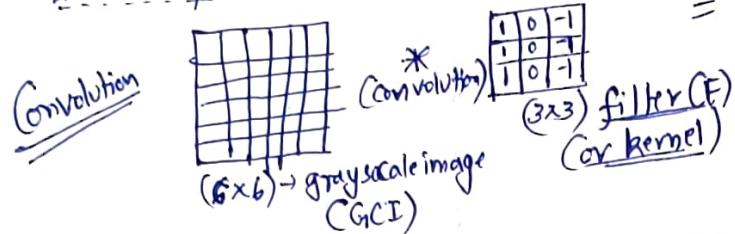
- ↳ Cons:
- ↳ need large amount of data
- ↳ excludes potentially useful hand-designed components in case of less data.

↳ end-to-end deep learning only when there is enough data to learn the complexity of the problem in hand.

Computer Vision

↳ one-problem: + Very large input parameter and hence input W and b
 (for a 1 Megapixel image \rightarrow input vector dim = $1000 \times 1000 \times 1000 \times 3$)

↳ An-example: \rightarrow Vertical edge detection:



$$\begin{aligned} &\rightarrow \text{sum}(\text{sum}(GCI[0:3, 0:3] * F[:, :])) \\ &\rightarrow \text{sum}(\text{sum}(GCI[0:3, 1:4] * F[:, :])) \end{aligned}$$

↳ Other vertical edge filters: \rightarrow

$$\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$$

Sobel filter

$$\begin{bmatrix} 3 & 0 & -3 \\ 10 & 0 & -10 \\ 3 & 0 & -3 \end{bmatrix}$$

Scharr filter

↳ learning the filter: \rightarrow

$$\begin{bmatrix} \text{img} \end{bmatrix} * \begin{bmatrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \\ w_7 & w_8 & w_9 \end{bmatrix} = \boxed{\quad}$$

learn them as "weights"
 using NN these filters to detect low-level features in images can be learnt easily.

$$\hookrightarrow (n \times n) * (f, f) = (n-f+1, n-f+1)$$

every time we apply convolution, the image shrinks,
 also, the pixels at the edges are used much less (~~info thrown away~~)

↳ Solution: \rightarrow Pad the image: $\rightarrow p = \text{padding} = 1 \rightarrow$

$$\begin{bmatrix} \text{img} \end{bmatrix} \rightarrow (n+2, n+2)$$

$$(n+2, n+2) * (3, 3) = (n, n)$$

Two types

"Valid": no-padding

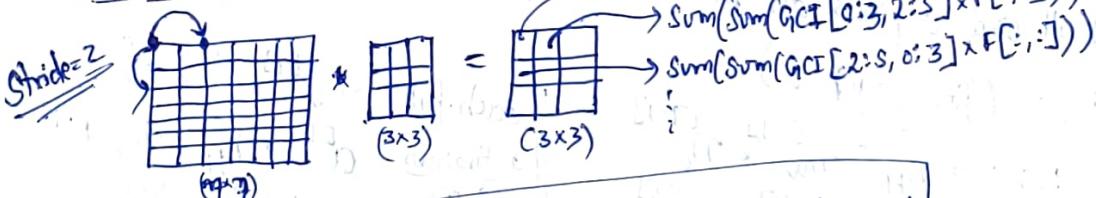
$$(n \times n) * (f, f) = (n-f+1, n-f+1)$$

"Same": same input and output size

$$(n+2p, n+2p) * (f, f) = (n, n)$$

by convention: in Computer vision
 $f = \text{odd}$

Strided Convolution



$$\rightarrow \text{sum}(\text{sum}(GCI[0:3, 0:3] * F[:, :]))$$

$$\rightarrow \text{sum}(\text{sum}(GCI[0:3, 2:5] * F[:, :]))$$

$$\rightarrow \text{sum}(\text{sum}(GCI[2:5, 0:3] * F[:, :]))$$

(general)
 (Stride=s) \rightarrow $(n \times n) * (f \times f) = \left\lfloor \frac{n+2p-f+1}{s} \right\rfloor \times \left\lfloor \frac{n+2p-f+1}{s} \right\rfloor$

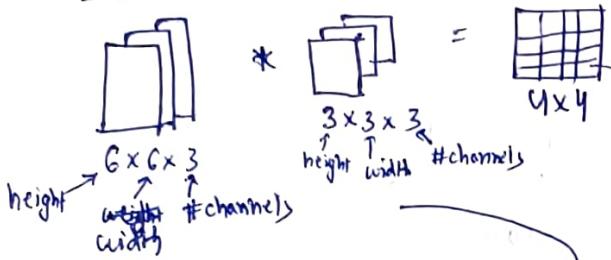
$\lfloor z \rfloor = \text{floor of } z$

Technical note on Cross-correlation and convolution

→ In reality the "convolution" done in dL is actually cross-correlation, but by convention it is being called convolution

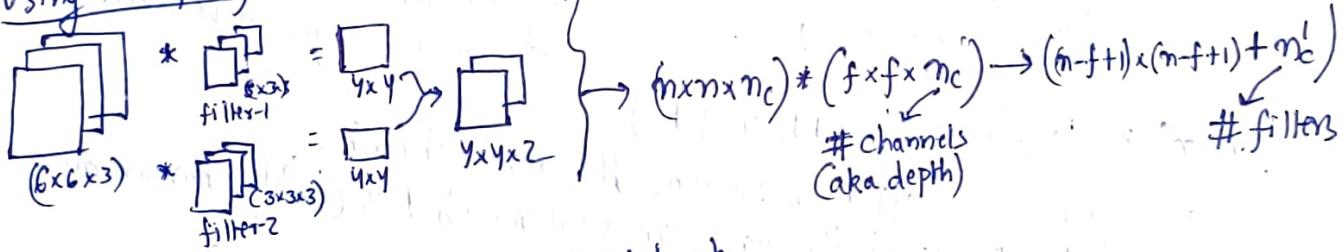
→ In actual convolution, an additional step of flipping the filter both horizontally and vertically is done before multiplying to preserve associativity ($(a * b) * c = a * (b * c)$)

→ Convolutions on RGB images :-



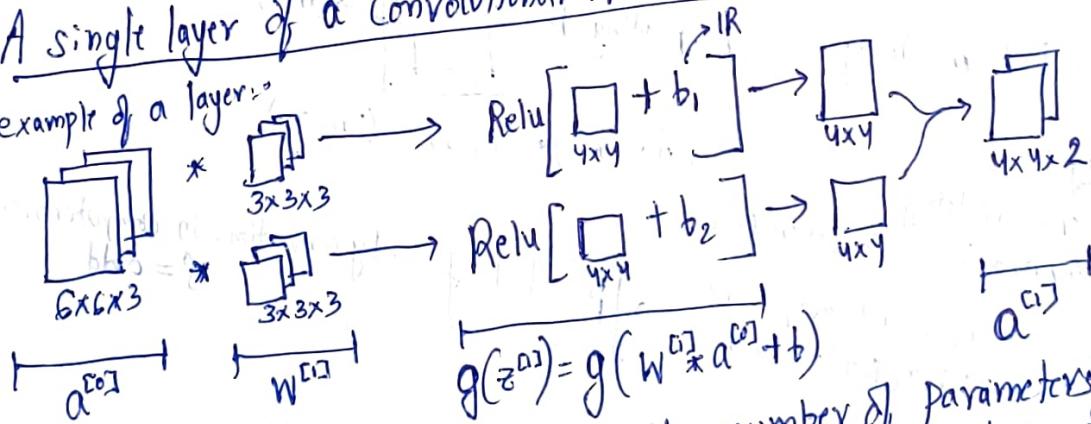
- Note:-
- (1) # channels of both the input and the filter must be same
 - (2) the output is a single layer always
 - (3) each number is the summation of the 27 multiplications (element wise) by sliding the cube over the whole image cube
 - (4) each channel of the filter detects features in corresponding each channel of input

→ Using multiple-filters at once



→ A single layer of a Convolutional Network

Example of a layer:-



* Note:- no matter how big the image is, the number of parameters is pretty small and depends on the number of filters and size of filters.

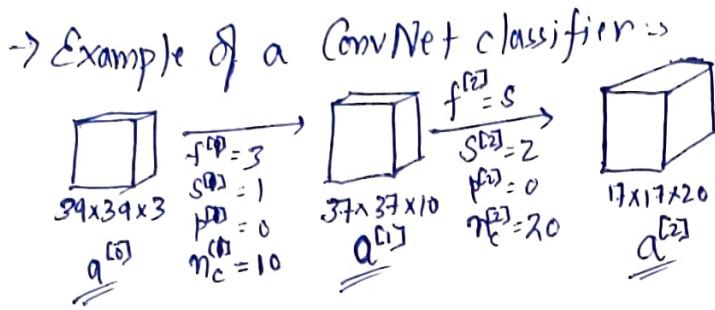
→ Summary of notation:- (for layer 2)

$f^{[2]}$	filter size
$p^{[2]}$	padding
$s^{[2]}$	stride

$n_H^{[2]}$	$n_W^{[2]}$	$n_C^{[2]}$
$m_H^{[2]}$	$m_W^{[2]}$	$m_C^{[2]}$
$n^{[2]}$	$m^{[2]}$	$n^{[2]}$

$$n^{[2]} = \left\lceil \frac{m^{[2]} + 2p^{[2]} - f^{[2]}}{s^{[2]}} + 1 \right\rceil$$

each filter is: $f^{[2]} \times f^{[2]} \times n_c^{[2]}$
activations: $a^{[2]} \rightarrow m_H^{[2]} \times m_W^{[2]} \times m_C^{[2]}$
 $(A^{[2]} \rightarrow m \times m_H^{[2]} \times m_W^{[2]} \times m_C^{[2]})$
weights: $f^{[2]} \times f^{[2]} \times m_C^{[2]} \times n_C^{[2]}$
bias: $n_C^{[2]} \rightarrow (1, 1, 1, n_C^{[2]})$



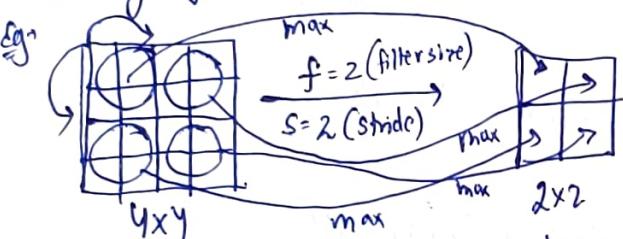
→ Types of layers in a convolutional network: →

↪ Convolution (CONV)

↪ Pooling (POOL)

↪ Fully connected (FC)

↪ Pooling layer: → Max pooling



take the max value from the sub-matrix on which the filter is currently situated

→ Intuition: → max pooling helps to preserve the info on the presence/absence of a feature at a particular region of the filter.

→ Input size: $n \times n \times n_c$

→ Output size: $\left\lfloor \frac{n+2p-f+1}{s} \right\rfloor \times \left\lfloor \frac{n+2p-f+1}{s} \right\rfloor \times n_c$

→ f, s are the hyperparameters here

→ No parameters to learn here

→ if multiple channels are present, the same computation is done over all the channels independently

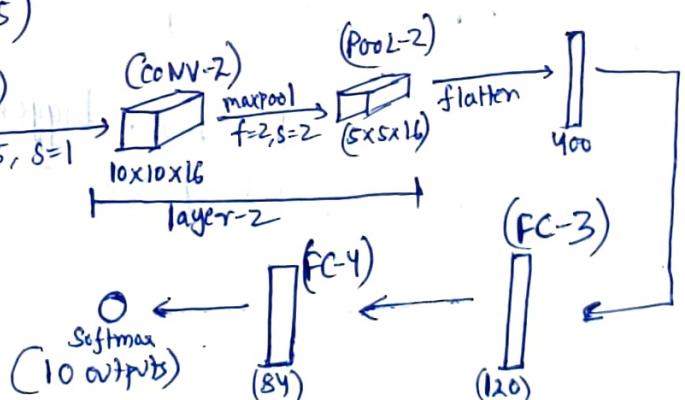
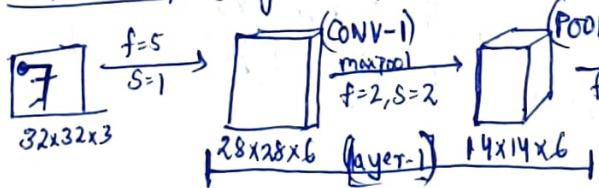
↪ Pooling layer: → Average pooling (Not used very often) → instead of taking max, average is taken.

↪ Sometimes used very deep into the neural network to collapse a larger representation (image) to a smaller one.

↪ Usual hyperparameters for Pooling: → $f=2, s=2$ or $f=3, s=2, p=0$, max pooling

↪ Note: → # of layers in a CNN is usually the # of layers with trainable parameters. So, the pooling layer is considered with the conv layer as a single layer.

↪ CNN example (digit recognition) (LeNet-5)



(38) → Fully Connected layers (FC): These are the visual ANN layers in which all input and output nodes are fully interconnected and they use the visual $W^{[2]}, b^{[2]}$ parameters.

↳ Note: In a visual CNN, as the network progresses, the $n_h, n_w \downarrow$ and $n_c \uparrow$

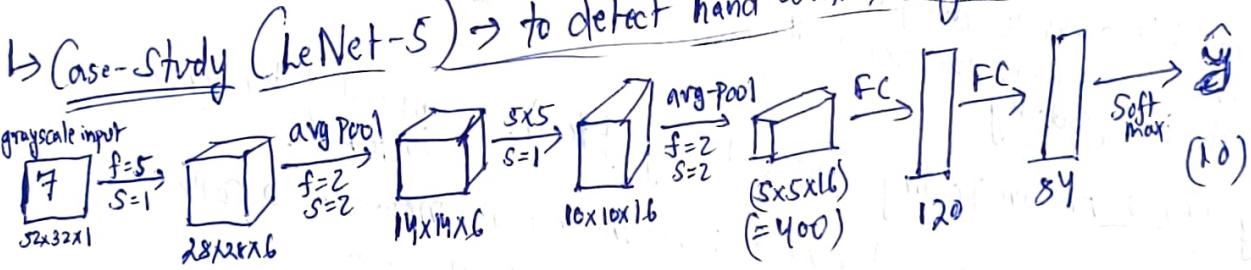
→ Why CNN can get away with smaller number of parameters compared to ANN

↳ Parameter sharing: A feature detector that's useful in one part of the image is probably useful in another part of the image.

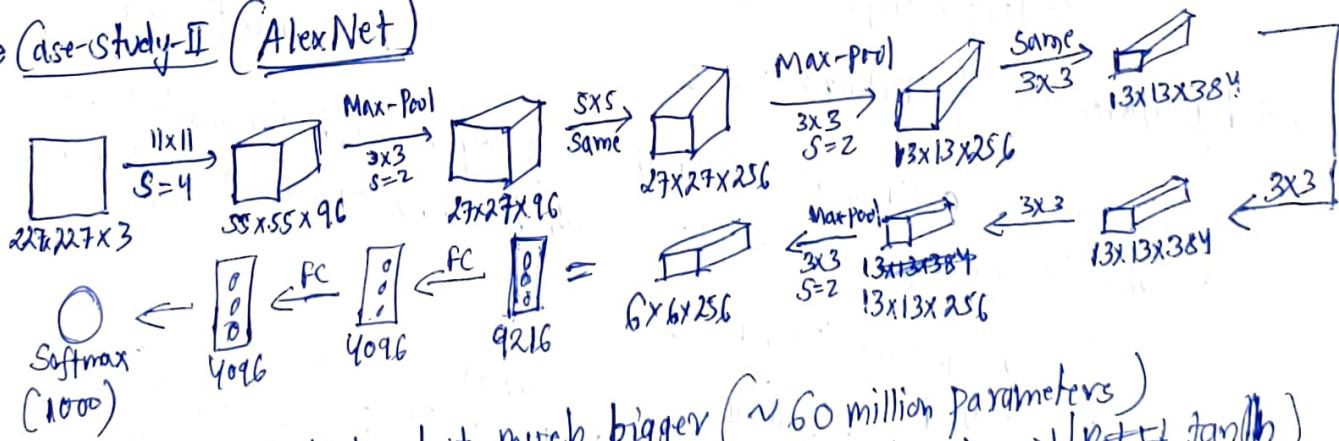
↳ Sparcity of connections: in each layer, each output value depends only on a small number of inputs.

↳ CNNs are good at working with translational invariance

↳ Case-Study (LeNet-5) → to detect hand written digits.



↳ Case-Study-II (AlexNet)



→ Similar to LeNet, but much bigger (~ 60 million parameters)

→ Uses ReLU compared to LeNet-5 which used sigmoid/ReLU tanh

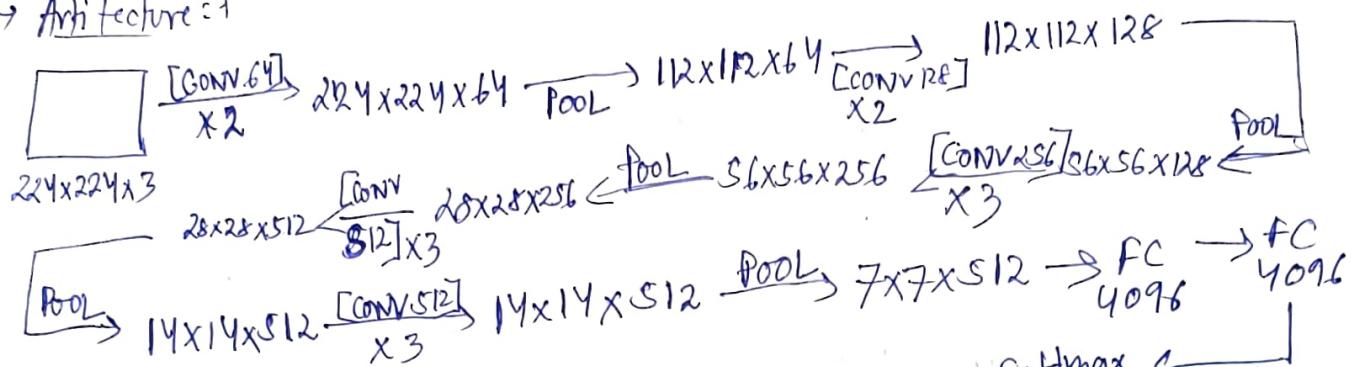
→ VGG1-16 (Case-study)

↳ uses the same CONV {3x3 filter, s=1, p=same} and Max-pool {2x2 filter, s=2}

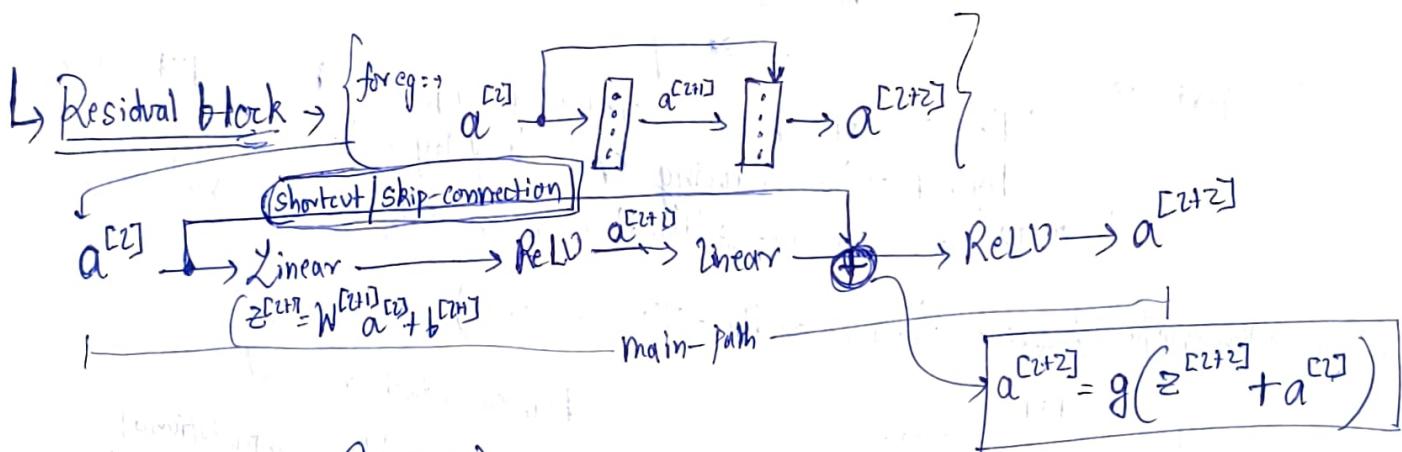
↳ layers throughout

↳ eg. $[CONV\ 64] \times 2 \equiv$ two conv layers with 64 filters each

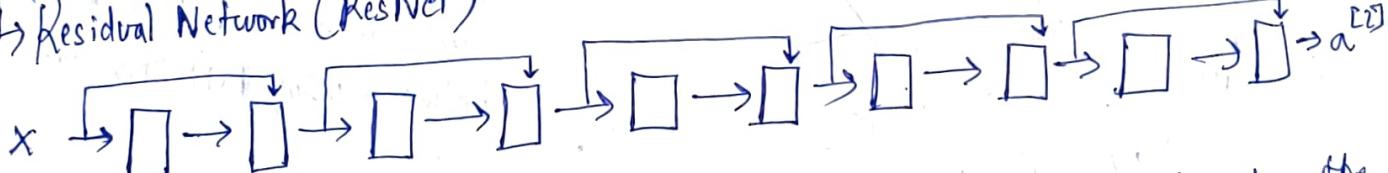
↳ Architecture:



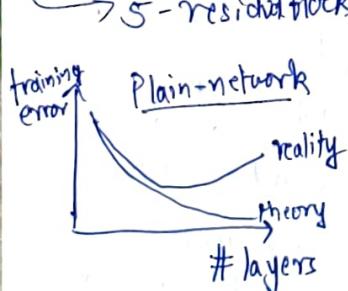
↳ has 128 million parameters in 16 layers with trainable weights.
↳ but is simplistic.



↳ Residual Network (ResNet)



↳ 5-residual blocks



→ ResNet very effectively solves the problem of exploding/vanishing gradients giving the option to have very deep Neural networks.

↳ As the # layers increases, we get the problem of exploding/vanishing gradients

→ Why do residual networks work?

↪ the residual blocks can easily learn the identity function in case of small W and b values so that the NN does not "stop" giving non-zero values.

$$a^{[L+2]} = g(C^L W^{[L+2]} a^{[L+1]} + b^{[L+2]} + a^{[L]}) \stackrel{\text{ReLU}}{=} g(a^{[L]}) = a^{[L]}$$

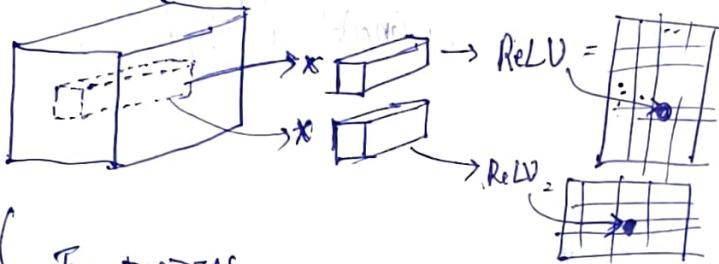
↪ So, this residual block can be effectively "skipped" in these scenarios without decaying the performance of the whole NN.

NT → if $\dim(a^{[L+2]}) \neq \dim(a^{[L]})$ in residual layer, then,

$$a^{[L+2]} = g(C^L z^{[L+2]} + W_s a^{[L]}) \quad || \quad W_s \in \mathbb{R}^{\dim(a^{[L+2]}) \times \dim(a^{[L]})}$$

either learnable or
to make zero-padding

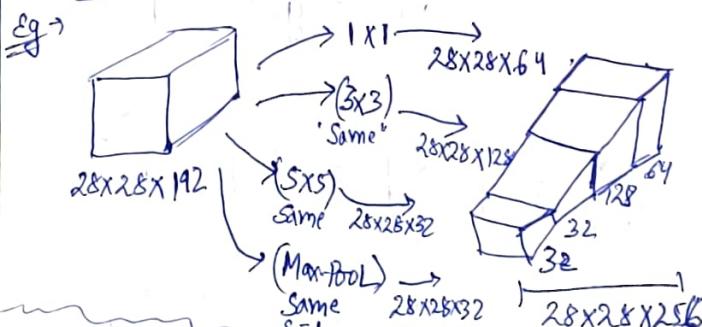
→ Network in Networks (1x1 convolutions)



≡ a 1×1 convolution acts like a neural network.

- ↪ Two purposes
 - like pool-layers helps in reducing n_H and n_W . 1×1 CONV layer helps in reducing n_C (# of layers)
 - It does also add additional non-linearity to the network

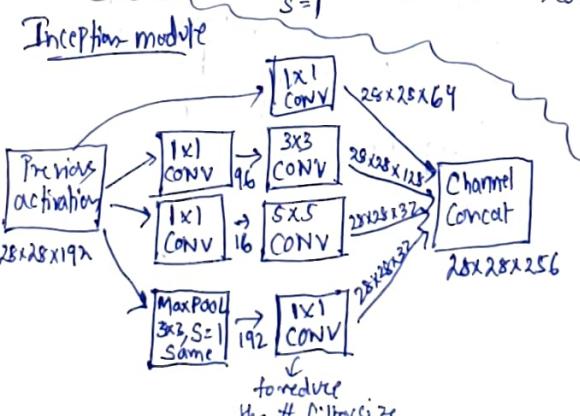
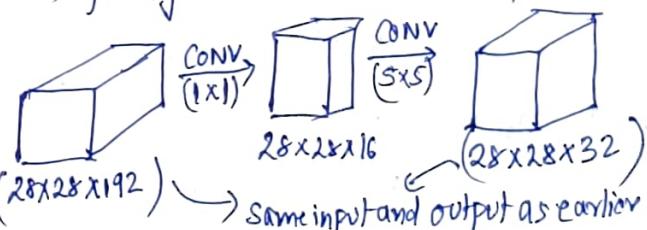
→ Inception network block (trying out all the filters at once)



↪ high computational cost

↪ to reduce the computational cost, we introduce a bottleneck layer in between using (1x1) CONV layer

for eg: →



Inception Network → made by stacking up inception modules

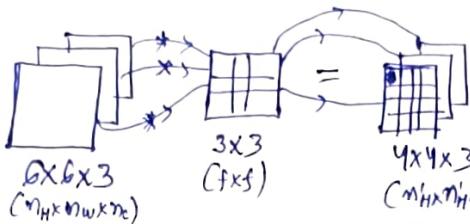
↪ later inception modules in this network has an additional branch ending to a softmax to check output from the hidden layers also called (GoogLeNet). { "we need to go deeper" } { - from inception movie }

→ Computational cost of a normal convolution:

$$\text{Ex: } \begin{matrix} 6 \times 6 \times 3 \\ \text{Input} \end{matrix} * \begin{matrix} 3 \times 3 \times 5 \\ \text{Filter} \end{matrix} = \begin{matrix} 4 \times 4 \times 5 \\ \text{Output} \end{matrix}$$

$$\left. \begin{array}{l} \text{Computational cost} = \# \text{filter params} \times \# \text{filter positions} \\ \times \# \text{of filters} \end{array} \right\} \Rightarrow \text{Cost} = (3 \times 3 \times 3) \times (4 \times 4) \times 5 = 2160$$

→ Depthwise-Convolution:



$$\left. \begin{array}{l} \text{Computational cost} = \# \text{filter params} \times \# \text{filter positions} \\ \times \# \text{of filters} \end{array} \right\} 432 = (3 \times 3) \times (4 \times 4) \times 3$$

→ Point-wise-Convolution:

$$\begin{matrix} 6 \times 6 \times 3 \\ \text{Input} \end{matrix} * \begin{cases} 1 \times 1 \times 3 \\ (1 \times 1 \times n_C) \\ m_C = 5 \text{ filters} \end{cases} = \begin{matrix} 4 \times 4 \times 5 \\ \text{Output} \end{matrix}$$

$$\left. \begin{array}{l} \text{Computational cost} = \# \text{filter params} \times \# \text{filter positions} \times \# \text{filters} \end{array} \right\} 240 = 1 \times 1 \times 3 \times 4 \times 4 \times 5$$

→ Depthwise-Separable convolutions:

$$\begin{matrix} m_H \times m_W \times n_C \\ \text{Input} \end{matrix} * \begin{matrix} m_H \times m_W \times n_C \\ \text{depthwise convolution} \end{matrix} * \begin{matrix} (1 \times 1 \times n_C) \\ \times n_C' (\# \text{filters}) \\ \text{Pointwise convolution} \end{matrix} = \begin{matrix} m_H \times m_W \times n_C' \\ \text{Output} \end{matrix}$$

Mobile Nets

↳ using depthwise-separable convolutions

↳ used to reduce computations

$$\frac{\text{Computation cost [MobileNets]}}{\text{Computation cost [Normal convolution]}} = \frac{1}{m_C'} + \frac{1}{f^2} \quad \left. \begin{array}{l} \text{typically} \\ \approx 1/10 \\ \text{as } m_C \text{ very} \\ \text{large and} \\ f = 3 \end{array} \right\}$$

→ MobileNet V1:

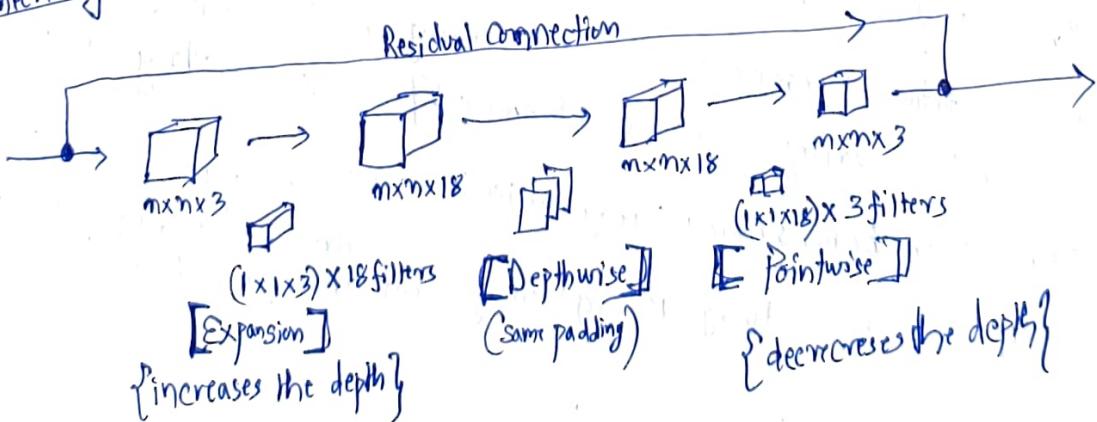
$$\begin{matrix} \text{Input} \end{matrix} \rightarrow \underbrace{\begin{matrix} \text{depthwise-separable} \\ \text{CONVS} \end{matrix}}_{\times 13} \rightarrow \text{POOL, FC, SOFTMAX}$$

→ MobileNet V2:

$$\begin{matrix} \text{Input} \end{matrix} \rightarrow \underbrace{\begin{matrix} \text{expansion} \rightarrow \text{depthwise} \rightarrow \text{projection} \end{matrix}}_{\times 17} \rightarrow \text{POOL, FC, SOFTMAX}$$

Bottleneck layer

Bottleneck layer



↳ Why mobilenet v2 bottleneck

→ it allows to learn a richer (more complex) function as it expands the previous activation, as well as keeping the sizes of the activations small, as it converts back (projects back) to the original size (small) of the activation.

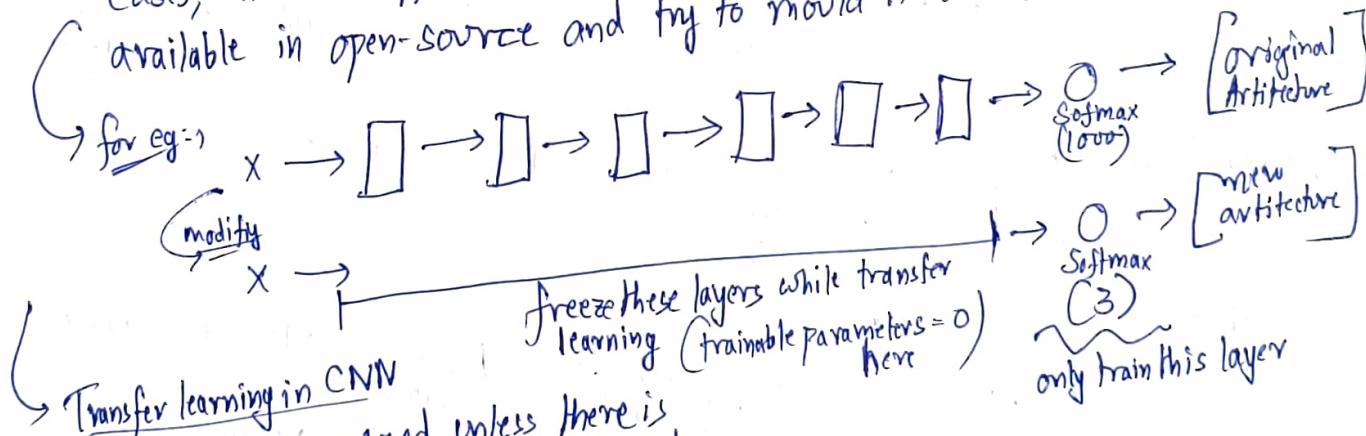
→ EfficientNet :> { how to make a CNN suitable for given computational resources }

- Three things to control :> ① r :> resolution of image ② d :> depth of NN ③ w : width of layers
- see open-source implementations of efficient-Net.

→ Practical Advices on using ConvNets

→ try to find open source implementations of the architecture that you are willing to try in sites like github.

→ Many a times there are very large Neural Networks which would take a lot of time to set up and to do the optimal parameters searches. In these cases, it is helpful to use transfer learning from pre-trained models available in open-source and try to mould it for your application.



Transfer learning in CNN

→ almost always used unless there is a large amount of data and computational resources at your service.

↳ Common data augmentation techniques in CV :> mirroring, random cropping, [rotation, shearing, local warping, -- less used techniques], Color shifting, PCA color augment

↳ Tips for doing well on benchmarks/winning Competitions :> Ensembling ; using different NNs and averaging their outputs

→ Multi-crop at testime :> taking (10-crop) images and averaging their outputs

→ Terminology used:-

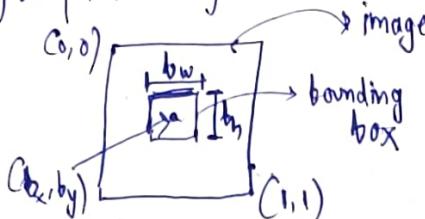
↳ (Image classification) and (Image classification and localization)

↳ both concerning only a single object classification.

(detection) → Concerning multiple objects classification and detection.

→ Classification with localization

↳ (target) output parameters for localization using bounding boxes:



class label

$$\hat{y} = (\text{softmax}, b_x, b_y, b_H, b_w)$$

Defining the target label for a classification like $\{o_1, o_2, \dots, o_m, \text{no-object}\}$
 {ie either one of the object present or no object there}.

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_H \\ b_w \\ c_1 \\ c_2 \\ \vdots \\ c_m \end{bmatrix} \rightarrow \begin{cases} 0: \text{no-object} \\ 1: \text{object there} \end{cases}$$

if $p_c = 0 \Rightarrow y = \begin{bmatrix} 0 \\ 2 \\ 2 \\ 2 \\ 2 \\ \vdots \\ 2 \end{bmatrix}$

if $p_c > 0 \Rightarrow y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$

Softmax for which object present.

One kind of loss function: $\rightarrow L(g_i, y) = \begin{cases} (g_i - y)^2 + -(\hat{y}_i - y_i)^2 & \text{if } y_i = 1 \\ (\hat{y}_i - y_i)^2 & \text{if } y_i = 0 \end{cases}$

→ Landmark detections → Landmarks are basically the x,y coordinates that can be taken as labels which would be representing certain key features in an image
 Eg:- Corner of eyes, etc.

→ Object detection using sliding windows detection

- First build a CNN which is trained on cropped images of single occurrences/non-occurrences of the object(s) of interest which would classify if the image is of that object or not.
- Then on the target actual images, slide a window (from top to bottom) which would take that part of image into CNN and see if that is that object or not.
- After that keep on increasing the window size.

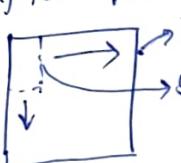


image
sliding window

0	0	1	0
0	1	0	1
1	1	0	1
0	0	0	0

increase
the window
size
and continue

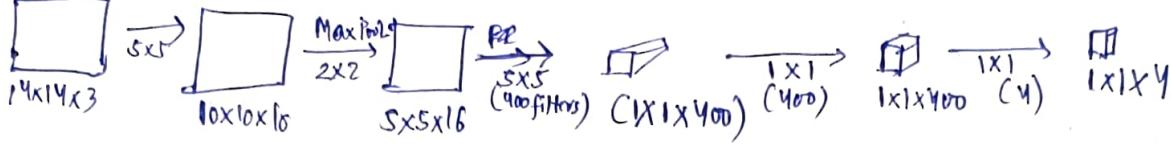
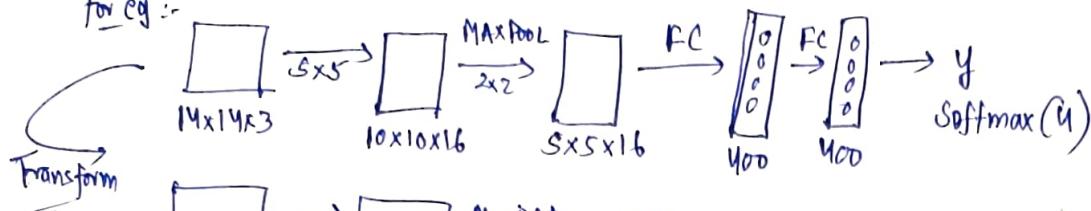
Drawback:-

↳ high computational cost at inference time.

→ Turning FC layers to CONV layers in a CNN

44

For eg:-



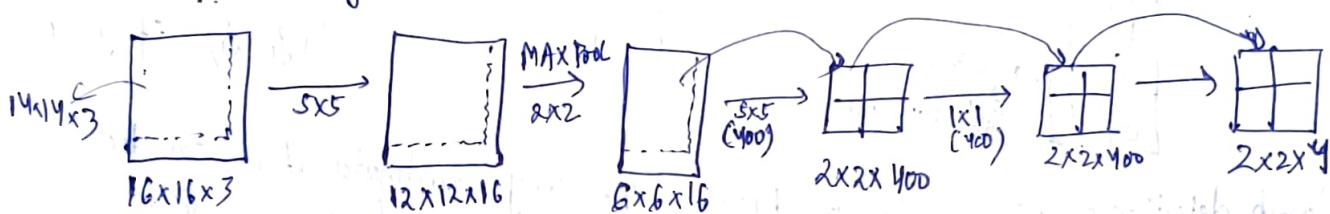
↳ Convolution implementation of Sliding window

→ it is based on the fact that sliding filter calculation shares a lot of similar calculations (as they are being calculated on shared parts of images)

→ So, it can be done by a single forward propagation of the CNN

(using the FC → CONV conversion)

→ let the sliding filter be of size 14x14x3 on an image of size 16x16x3 (can be of any greater size)

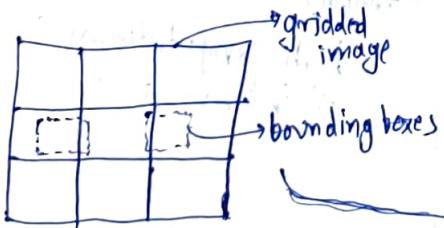


↳ Drawback: → due to fixed shape and position of the bounding boxes, the algorithm is not so accurate in getting the correct (nearly) positions of them.

↳ YOLO [You Only Look Once] Algorithm

→ follows the CNN way of splitting the image into 'grids' and each grid is associated with a target vector $\hat{y} = \begin{pmatrix} p_c, b_x, b_y, b_h, b_w \\ c_1, c_2, \dots, c_n \end{pmatrix}$ → as seen in classification with localization

→ Note that it is done by single CNN network which has out layer of $(3 \times 3 \times 8)$ → {for eg for a grid of 3x3 on the image having $\text{len}(\hat{y}) = 8\}$



↳ Assumptions: → ① The object is associated to the grid box where its center is lying
② The grid sizes are small enough so that no two objects are in the same grid.

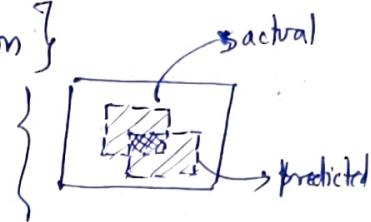
$\rightarrow x \rightarrow \boxed{\quad} \rightarrow \text{CONV} \rightarrow \text{MaxPool} \rightarrow \dots \rightarrow \boxed{3 \times 3 \times 8}$

So in YOLO, the bounding boxes are given by the (bx, by, bh, bw) labels of each grid box

Note that (bx, by, bh, bw) are relative to a particular grid box and they must be transformed according to the whole image.

Evaluating Object Localization of Intersection over Union

$$IoU = \frac{\text{Intersection of ground truth and Predicted bounding box}}{\text{Union of ground truth and predicted bounding box}}$$



Usually $IoU \geq 0.5 \Rightarrow$ Correct detection (Convention)

Non-max Suppression Algorithm

One problem with YOLO: multiple grid-boxes might predict multiple different bounding boxes for the same object

Solution: non-max suppression algorithm (for one object detection)

Each output prediction from each grid boxes \rightarrow

$$\begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$$

discard all boxes with $P_c \leq 0.6$

for remaining boxes \rightarrow pick the box with the largest P_c

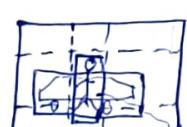
discard any remaining box with $IoU \geq 0.5$ with the box output in the previous step

Note

For multi-class
Object detection

run the algorithm independently for each class.

→ Anchor boxes → used when there are more than one objects in the same grid



pre-define shape of anchor boxes for each grid {eg: \square and \square }

the output label size gets multiplied by the number of anchor boxes.

With two Anchor boxes: → Each object in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IoU

P_c	related to anchor box
b_x	
b_y	
b_h	
b_w	
C_1	
C_2	
b_x	
b_y	
b_h	
b_w	
E_1	
E_2	

$$Q = \begin{bmatrix} 1, b_x, b_y, b_h, b_w, 1, 0, \\ 1, b_x, b_y, b_h, b_w, 0, 1 \end{bmatrix}$$

Limitations: → Can't work if # of objects in single grid $>$ # anchor boxes

Can't work if the objects in the same grid have the similar anchor box shape.

R-CNN (Region Proposal CNN)

- instead of running CNN over the whole image, R-CNN runs CNN only on particular regions of the image where there is high chance of finding the object using segmentation algorithm.
- R-CNN → Propose regions, classify proposed regions one at a time. output label + bounding box
- Fast R-CNN → propose regions, use convolution implementation of sliding windows to classify all the proposed regions.
- Faster R-CNN → use convolutional network to propose regions.

→ Semantic Segmentation → in this each pixel is classified to which class it belongs to. This helps in outlining the objects properly.

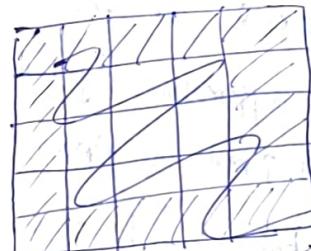
→ The input and output size is to be same.
→ as the CNN goes deeper, the size of the image block needs to be increased and filter size decreased.

Transpose Convolutions

an example

2	1
3	2

$$\begin{array}{|c|c|c|} \hline 1 & 2 & 1 \\ \hline 2 & 0 & \\ \hline 0 & 2 & 1 \\ \hline \end{array} \quad f \times f = 3 \times 3, p = 1, s = 2$$

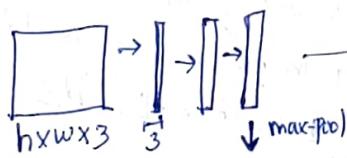


1	1	1	1	1
0	2+2	0	1	
4+6	2+0	2+4	1+2	
1+3	2+2	0	1	
0	3+4	0	2	
6	3+0	4	2	

{ in case of overlapped boxes in the output, the values are just added }

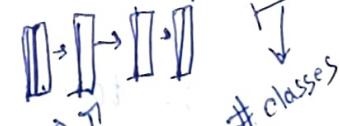
{ keep multiplying each value in the input with the filter values and placing the multiplied values on top of the output }

U-Net (for semantic segmentation)

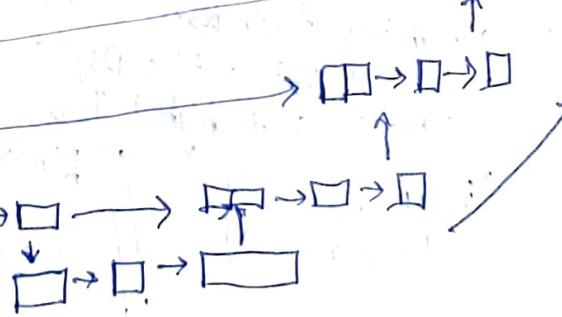


(high resolution, low level feature into)

Skip connection



Trans-Conv



{ NT skip connections used to deliver high resolution info from earlier layers along with high level feature info from previous layer }

Face verification

- ↳ input image, name/ID
- ↳ output whether the input image is that of the claimed person

Face recognition

- ↳ has a database of K persons
- ↳ get an image as input
- ↳ output id if the image is any of the K persons.

One-shot learning :> learning from one example

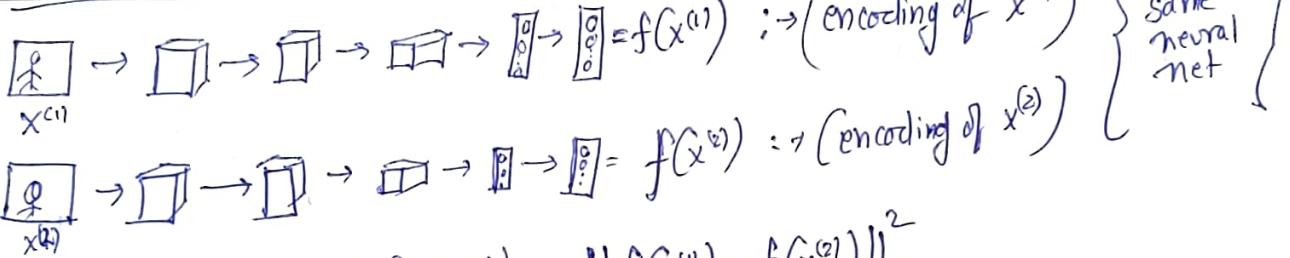
→ Problem with face recognition: it has to be one-shot, also can't use softmax max layer as output as each time a new person is added to the database, the model have to be retrained.

Learning a similarity function

$d(\text{img1}, \text{img2}) = \text{degree of difference b/w images}$

↳ if $d(\text{img1}, \text{img2}) \leq \tau \Rightarrow \text{"same-person"}$
 ↳ if $d(\text{img1}, \text{img2}) > \tau \Rightarrow \text{"different-person"}$

Siamese network



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

↳ Goal of learning in Siamese network: learn parameters so that:
 ↳ Parameters of NN here define an encoding $f(x^{(i)})$

$d(x^{(1)}, x^{(1)})$ is small for same person's images
 $d(x^{(1)}, x^{(2)})$ is large for different person's images

→ Triplet loss: three images compared at a time:
 ↳ Anchor image (A) → Positive image (P) { same person example }
 ↳ Anchor image (A) → Negative image (N) { different person example }

$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha \leq 0$$

→ margin { to ensure a distance b/w positive diff neg diff }

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

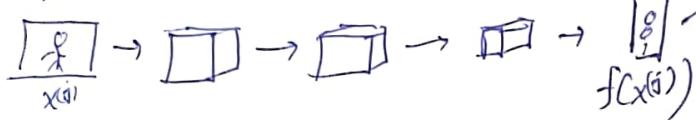
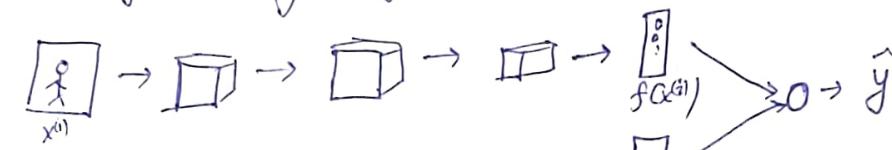
$$J = \sum_{i=1}^m \mathcal{L}(A^{(i)}, P^{(i)}, N^{(i)})$$

Note: for training with triple loss, we require multiple images of the same person. once the training of the encoder is done, we can try out the one-shot learning.

(48)

↳ Choosing the triplets $A, P, N \Rightarrow$ if A, P, N are chosen randomly then $d(A, P) + \alpha \leq d(A, N)$ is easily satisfied.
 choose triplets that are hard to train.
 (take triplets where $d(A, P) \approx d(A, N)$)

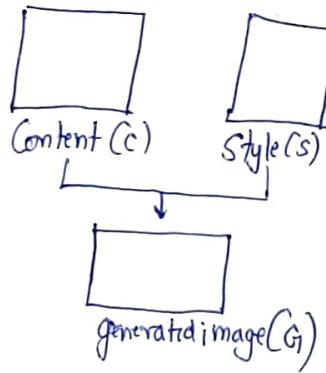
→ Face recognition using binary classification



$$\hat{y} = \sigma \left[\sum_{k=1}^{128} w_k |f(x^{(i)})_k - f(x^{(j)})_k| + b \right]$$

Other variations: $\frac{(f(x^{(i)})_k - f(x^{(j)})_k)^2}{f(x^{(i)})_k + f(x^{(j)})_k}$

→ Neural Style Transfer



Cost function

$$J(G_i) = \alpha J_{\text{content}}(C, G_i) + \beta J_{\text{style}}(S, G_i)$$

↳ finding the generated image (G_i)

↳ (1) initialize the G_i randomly

$$G_i: 100 \times 100 \times 3$$

↳ (2) use gradient descent to minimize $J(G_i)$

→ Content Cost function

↳ say you use hidden layer 2 to compute content cost. (visually a layer near the middle)

↳ use pretrained ConvNet to pass the content (C) and generated (G_i) images.

↳ let $a^{[2]}(C)$ and $a^{[2]}(G_i)$ be activations from layer 2 on C and G_i images

$$J_{\text{content}}(C, G_i) = \frac{1}{2} \|a^{[2]}(C) - a^{[2]}(G_i)\|^2$$

(99)

→ Style cost function :-

↳ measure of style for a particular layer l of CNN :- how much correlations there is present between the different channels (having different features) Susceptibility

↳ ie which feature appears along with which other features -
not appears

Style matrix

$$G^{[l]} = \sum_{i=1}^{n_H^{[l]}} \sum_{j=1}^{n_W^{[l]}} a_{ijk}^{[l]} a_{ijk'}^{[l]} \quad \left| \begin{array}{l} a_{ijk}^{[l]} = \text{activation at } (i, j, k) \\ G^{[l]} \in \mathbb{R}^{n_C^{[l]} \times n_C^{[l]}} \end{array} \right.$$

height
channel
width

$G_{kk'}^{[l](s)}$:- for style image // $G_{kk'}^{[l](g)}$:- for generated image

↳ Called "gram" matrix in linear algebra.

↳ Style Cost function :- $J_{\text{style}}^{[l]}(S, G) = \|G^{[l](s)} - G^{[l](g)}\|^2$

$$= \frac{1}{(2n_H^{[l]}n_W^{[l]}n_C^{[l]})^2} \sum_k \sum_{k'} (G_{kk'}^{[l](s)} - G_{kk'}^{[l](g)})^2$$

↳ Gives better results if we use different layers for style cost function :-

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

* Note, 2D-CONV Nets can be easily generalized to 1D-CONV Nets and 3D-CONV Nets by taking 1-dimensional and 3-dimensional filters respectively.

(50)

Notations in RNN

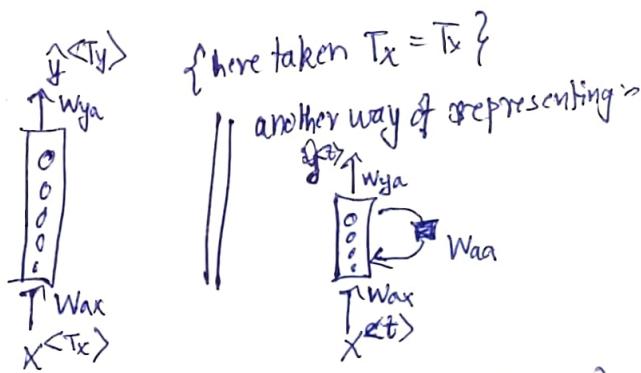
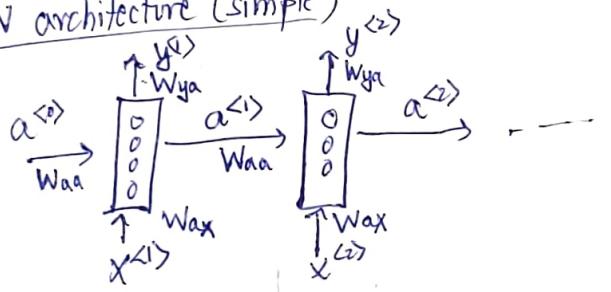
input sequence: $x = \text{Harry Potter and Hermione Granger invented a new spell.}$
 $(T_x = \text{total length of the } x)$ $x^{(1)} \quad x^{(2)} \quad \dots \quad x^{(T_x)}$

output sequence: $y = [1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 0 \ 0]$
 $y^{(1)} \quad y^{(2)} \quad y^{(3)} \quad \dots \quad y^{(T_x)}$

↳ the words are given as input as a one-hot encoding of a vocabulary of $> 10,000$ words or so. This means each input word is a $> 10,000$ dimensional vector with ~~1~~ at only one place and 0 everywhere else.

↳ Why not a standard ANN:
 ↳ (1) Input and outputs can be of different lengths
 ↳ (2) does not share features learned across different positions of text

RNN architecture (simple)



$$a^{(0)} = \vec{0} \quad \| \quad a^{(1)} = g_1(W_{aa} a^{(0)} + W_{ax} x^{(1)} + b_a) \quad \| \quad a^{(t)} = g_1(W_{aa} a^{(t-1)} + W_{ax} x^{(t)} + b_a)$$

$$\| \quad y^{(1)} = g_2(W_{ya} a^{(1)} + b_y) \quad \| \quad \hat{y}^{(t)} = g_2(W_{ya} a^{(t)} + b_y)$$

↳ common choice of activation here: tanh / relu

simplified RNN notation

$$a^{(t)} = g(W_a [a^{(t-1)}, x^{(t)}]_x + b_a)$$

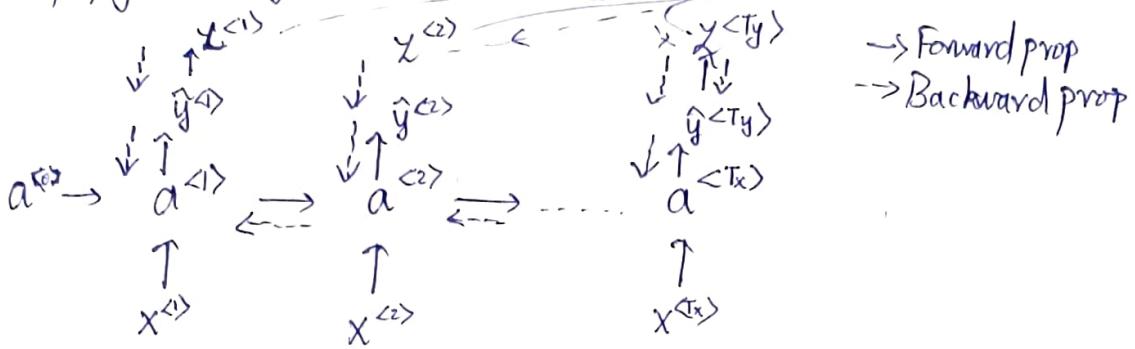
where, $W_a = [W_{aa} \ | \ W_{ax}]$

$$\rightarrow \text{Loss at a particular time } t \\ \mathcal{L}^{(t)}(\hat{y}^{(t)}, y^{(t)}) = -\hat{y}^{(t)} \log \hat{y}^{(t)} - (1 - \hat{y}^{(t)}) \log (1 - \hat{y}^{(t)})$$

$$\text{total loss } \mathcal{L} = \sum_{t=1}^{T_y} \mathcal{L}^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

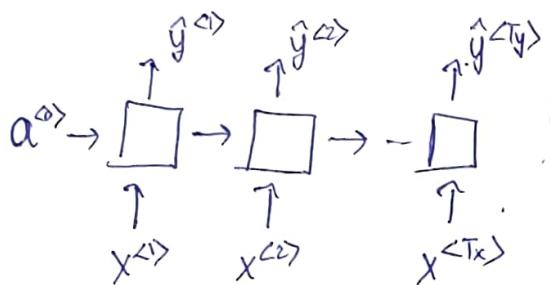
(SI)

→ Backpropagation through time for RNN

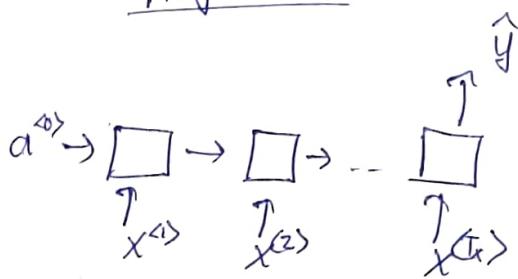


→ Examples of RNN architecture:

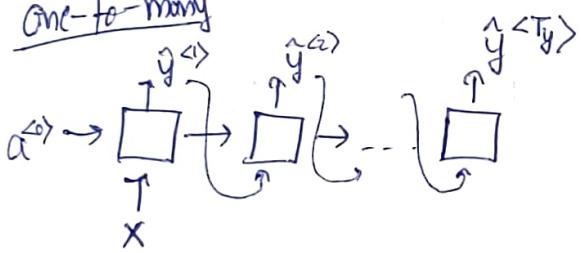
Many-to-many ($T_x = T_y$)



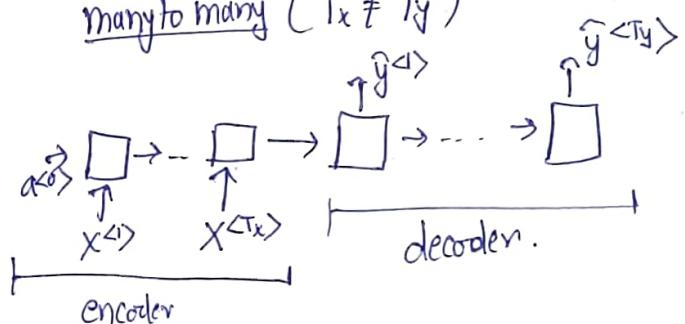
Many-to-one



One-to-many



many-to-many ($T_x \neq T_y$)



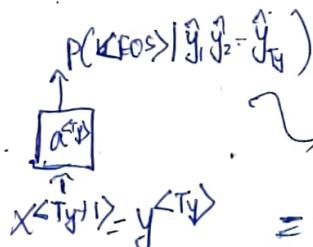
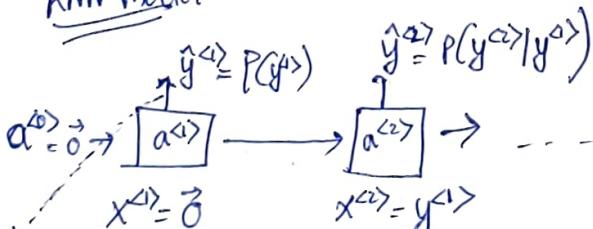
→ Language Modelling → given a sentence, predict its probability of occurrence,
 $P(\text{sentence}) = P(y^{<1>} | y^{<2>} | \dots | y^{<Ty>})$

Before training, first need to tokenize the text.

for eg.: $\text{cats averages } 15 \text{ hours of sleep a day. } \langle \text{EOS} \rangle$

→ taken from a dictionary
 if word not in dictionary
 ⇒ use <unk> (unknown) tag

RNN model



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

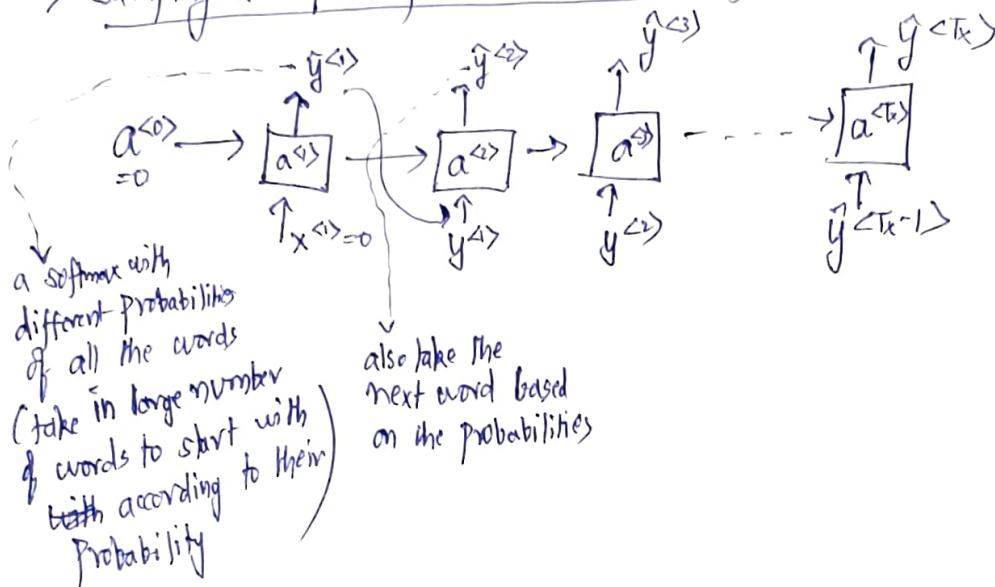
$$= P(y^{<1>}) P(y^{<2>} | y^{<1>}) \dots P(y^{<Ty>} | y^{<Ty-1>})$$

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

softmax layer
of dimension of dictionary + 1

(52)

→ Sampling a sequence from a trained RNN for language model



→ Problem of vanishing gradients

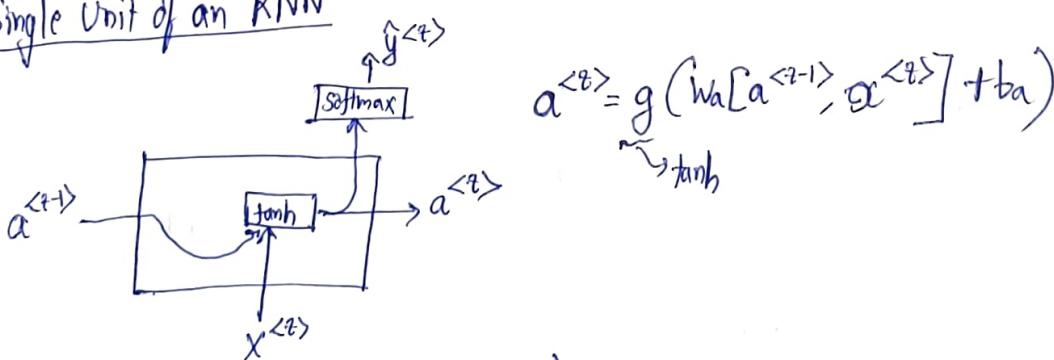
for e.g. → The fat, which were --, were full } long distance b/w the word
The cat, which was --, was full } (at+3) which directly affects word was/were

Because long sequences in RNN ⇒ vanishing gradients can't be able to affect these far words relationships

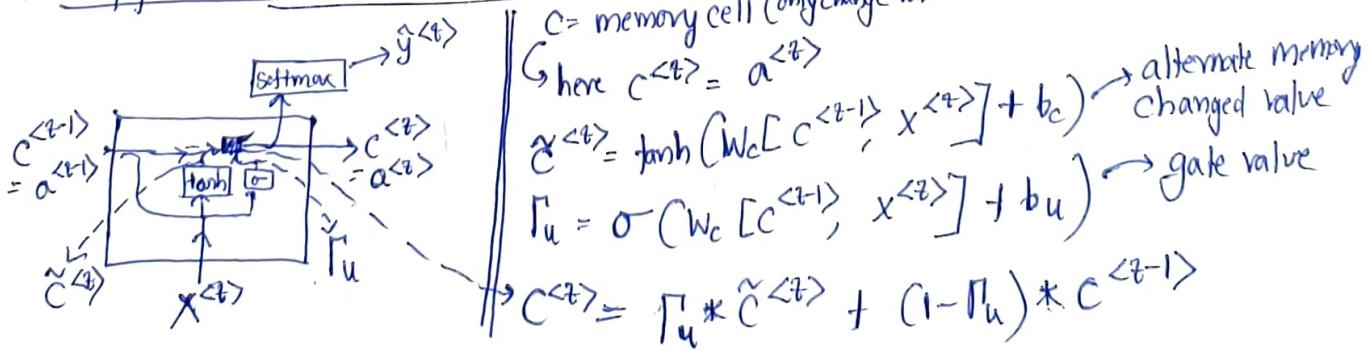
NA → Problem of Exploding gradients: easier to handle, can be detected by NaN values of the gradients and can be handled by gradient clipping

Upper limit on gradient value

→ Single Unit of an RNN



→ Simplified GRU (Gated Recurrent Unit)



↳ So, if Γ_u is kept to be zero (or near to zero), then the value of (53)
 $C^{<t>}$ is kept across many long distance (in time) to influence the
values of further outputs \rightarrow helps to address the vanishing gradient problem
Can learn long distance relationships.

↳ Full GRU

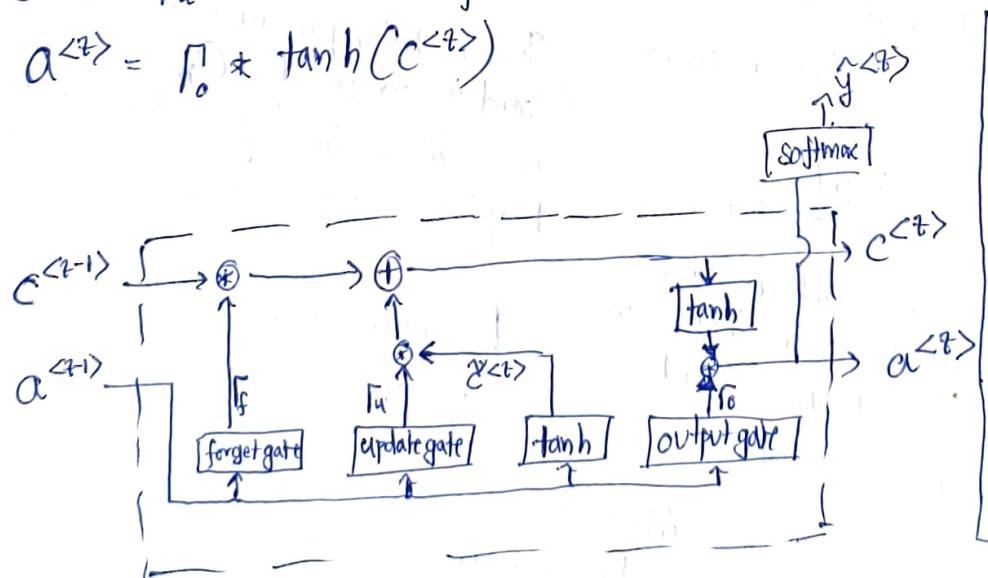
$$\begin{aligned}\tilde{C}^{<t>} &= \tanh(W_c[\Gamma_r * C^{<t-1>}, x^{<t>}] + b_c) \\ \Gamma_u &= \sigma(W_c[C^{<t-1>}, x^{<t>}] + b_u) \\ \Gamma_r &= \sigma(W_r[C^{<t-1>}, x^{<t>}] + b_r) \\ C^{<t>} &= \Gamma_u * \tilde{C}^{<t>} + (1 - \Gamma_u) * C^{<t-1>}\end{aligned}$$

↗ relevance gate
Why Γ_u ?
↳ to map the relevance of
 $C^{<t-1>}$ in calculating
 $C^{<t>}$

→ LSTM (Long-Short Term Memory)

$$\begin{aligned}\tilde{C}^{<t>} &= \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c) \\ \Gamma_u &= \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u) \\ \Gamma_f &= \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f) \\ \Gamma_o &= \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o) \\ C^{<t>} &= \Gamma_u * \tilde{C}^{<t>} + \Gamma_f * C^{<t-1>} \\ a^{<t>} &= \Gamma_o * \tanh(C^{<t>})\end{aligned}$$

↗ here $C^{<t>} \neq a^{<t>}$
↗ update gate
↗ forget gate
↗ output gate

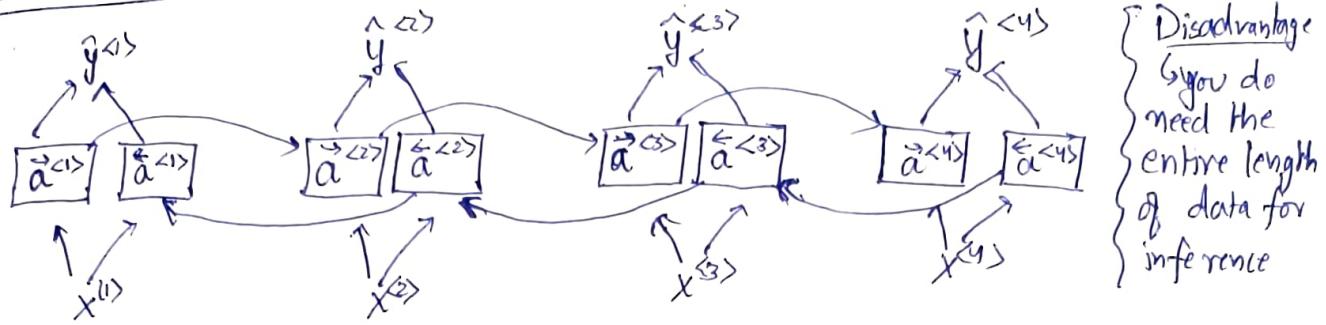


- NT → GRU's are more recent invention than LSTM
- NT → GRU's simpler unit can allow to have larger RNNs
- NT → but LSTM's have more power due to model complexity.

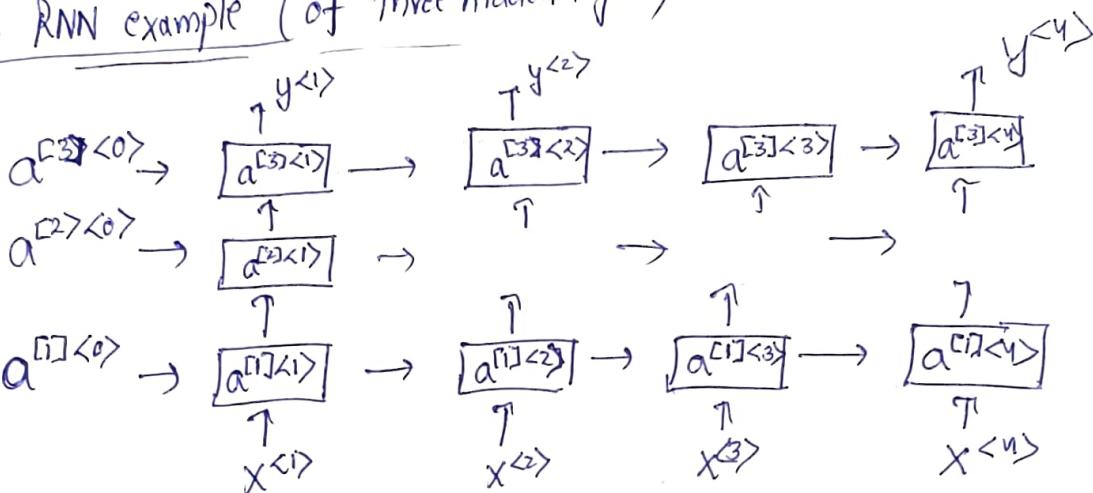
Bidirectional RNN

$$\hat{y}^{<2>} = g(W_y [\vec{a}^{<2>} \ \ \check{a}^{<2>}] + b_y)$$

(54)



Deep RNN example (of three hidden layers)



for eg. $a^{[2]}(3) = g[W_a^{[2]} [a^{[2]}(2), a^{[1]}(3)] + b_a^{[2]}]$

NT: We don't see as many deep recurrent layers, but the output of the recurrent layers can be fed to deep ANN layer (not connected temporally)

↳ In the one-hot representation of the words, all the words are having zero inner product with each other. That means, we don't have any relations between the individual words.

{ In One-hot representation, each word is denoted by On where n is the position of vector where 1 appears (zero everywhere else)}

Factorized Representation (Word Embedding)

↳ each word is associated with a feature vector which gives values for different features related to it. The word for eg. gender, age, cost, etc.

{ Visualize word embeddings - using t-SNE}

{ Approach: Use a large unlabelled text to figure out the word embedding and use a smaller data corpus for the job in hand}

55

→ Transfer learning and word embeddings

- ↳ learn word embedding from large text corpus. (1-10B words)
- ↳ transfer embedding to new task with smaller training set
- ↳ Optional: Continue to finetune the word embeddings with new data.

→ Analogies using word vectors

- ↳ to find out analogical word to the given word, given a prior analogy.
- ↳ e.g.: Man is to Woman as King is to ? → Ans: Queen.
- ↳ So, here the differences in their corresponding word vectors must be similar to reflect same analogy. i.e.
 $e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\text{?}}$
- ↳ find word w such that $\underset{w}{\operatorname{argmax}} \left[\text{Sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}}) \right]$
- ↳ similarity function.
- ↳ Cosine Similarity

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

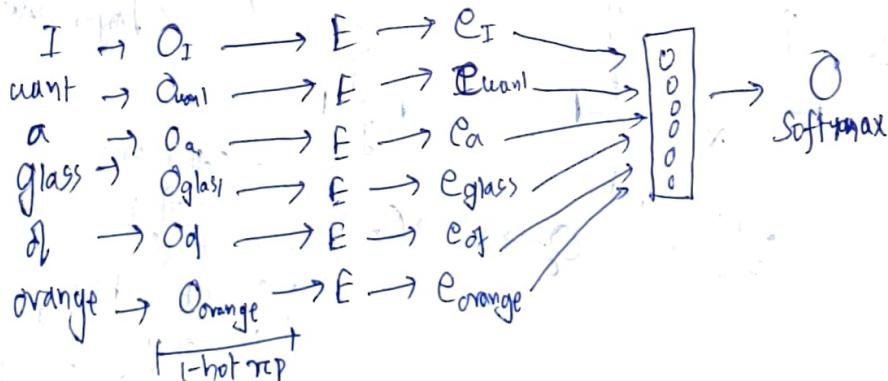
→ Embedding matrix (E)

$$\begin{matrix} \# \text{ of words} \\ \# \text{ of features} \end{matrix} \xrightarrow{\text{T}} E \times \begin{matrix} \text{one-hot vector} \\ \downarrow \end{matrix} = \begin{matrix} e \\ \downarrow \text{feature vector} \end{matrix}$$

To learn to get word embeddings

↳ Neural language model to learn word embeddings

- ↳ learn word embedding by training a neural net to predict the next word
- Given n - # of prior words.
- e.g.: I want a glass of orange juice
- ↳ here E is the embedding matrix where parameters need to be learned.



→ Other context/target pairs

e.g.: I want a glass of orange juice
Context target

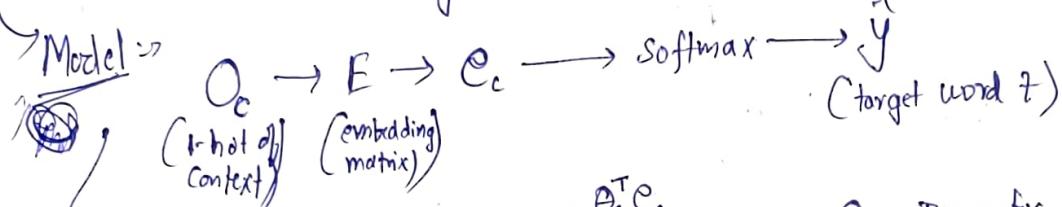
Other context

- ↳ four words on left and right
- ↳ last one word
- ↳ nearly one word.

(Word2Vec)

→ Skip-grams : { given a Context word in a sentence, the problem statement is to predict a randomly chosen Target word in the vicinity (say ± 3 words) of the context word }

Model :-



$$\text{Softmax} : P(z|c) = \frac{e^{\theta_z^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}}$$

θ_z : parameter associated with output z

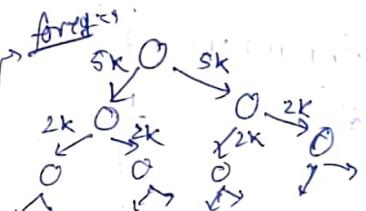
loss-function :-

$$\mathcal{L}(\hat{y}, y) = -\sum_{i=1}^{10000} y_i \log(\hat{y}_i)$$

↳ Problem with Softmax Classification

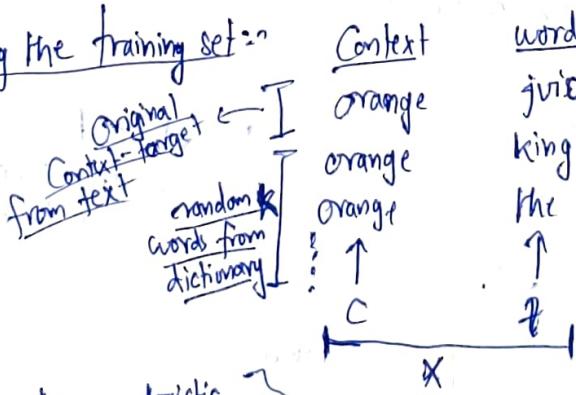
↳ large computational cost to calculate softmax

↳ use hierarchical softmax classifier



↳ Negative Sampling { Problem statement :- given two words, need to determine if they are context-target words }

↳ building the training set :-



target 2

↳ Choosing-k

$$K = 5-20$$

(smaller dataset)

$$K = 2-5$$

(larger dataset)

We will use binary logistic regression here

$$\text{i.e. } P(y=1|c, t) = \sigma(\theta_t^T e_c)$$

↳ neural net structure

$$O_c \rightarrow E \rightarrow E_c$$

$$\sigma(\theta_t^T e_c)$$

↳ Number of logistic regressors = # of words in document
↳ but here we will only use $K+1$ number of regressor

→ Hence, in negative sampling, we are reducing the softmax having 10000 (# of words in dictionary) different classes to $k+1$ different binary logistic regressors out of which we only need to train $k+1$ of them.

Selecting negative samples: two extremes

- uniform random selection from the text corpus.
- cover representation of words like, 'the', 'of'
- uniform random selection from the dictionary (can't have distribution reflective of that in the real world)

a way out (heuristic)

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10000} f(w_j)^{3/4}}$$

f(w_i): frequency of word w_i in the text corpus

→ GloVe (global vectors for word representation)

↳ define a way to get context and target pairs (c, t)

↳ define $X_{ij} = \# \text{ times } j \text{ appears in context of } i$

↳ GloVe: take def of (c, t) as based on proximity (target word is in the context of c if it is in its proximity)

$$\Rightarrow X_{ij} = X_{ji}$$

Model: → need to minimize the function: \rightarrow (let # of words in dictionary = $n = 10000$)

$$\sum_{i=1}^{10000} \sum_{j=1}^{10000} f(X_{ij}) \left[\underbrace{\Theta_i^\top c_j}_\text{dot product of target and context word vectors} + \underbrace{b_i}_t + \underbrace{b_j}_c - \log(X_{ij}) \right]^2$$

weighting term

↳ $f(X_{ij}) = 0$ if $X_{ij} = 0$ (as $\log(0) = \text{m.d.}$)

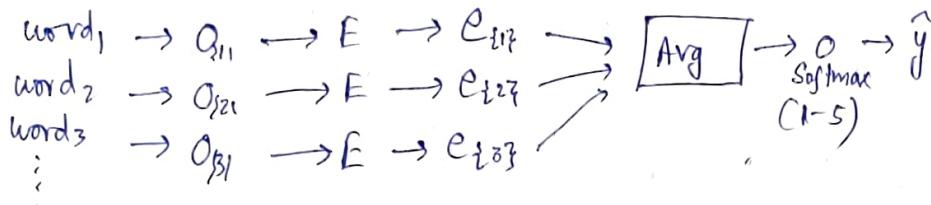
↳ can also be taken (heuristically) such that it is avoided to give too much weightage to frequent words like this, is, a, of, etc and too little weightage to rare words.

→ Sentiment classification Problem

given a comment find the associated rating (5 stars/4 stars/-) }
 → data for training might be less here so word embedding might help

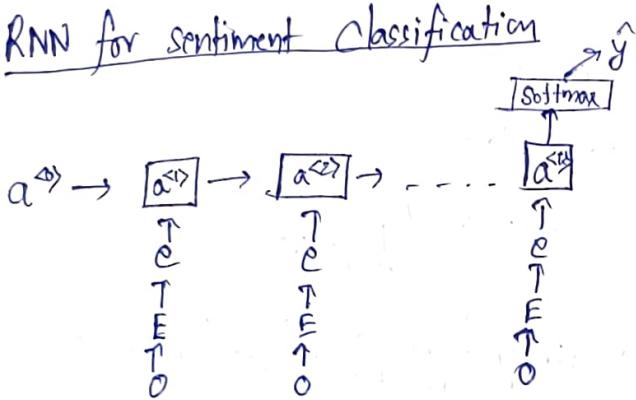
(S8)

↳ Simple sentiment classification model



Limitation: does not preserve word-ordering but preserves only word-meaning. So might not work in cases where ordering is important
 (e.g. "completely not good" +ve)

↳ RNN for sentiment classification



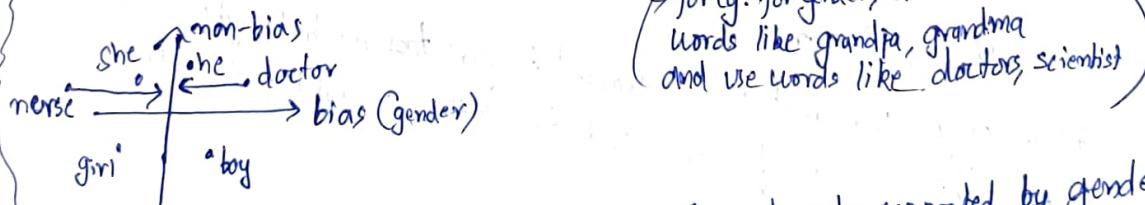
↳ Reducing bias in word embeddings

{ word embeddings can reflect gender, age, etc., }
 { biases of the text used to train the model }

① identify bias direction

Eg. for gender bias take { she - eshe } → [arg] → bias direction
 (definition based on gender) { female - female }
 and are separated by gender only

② Neutralize: for every word that is not definitional, project to get rid of bias



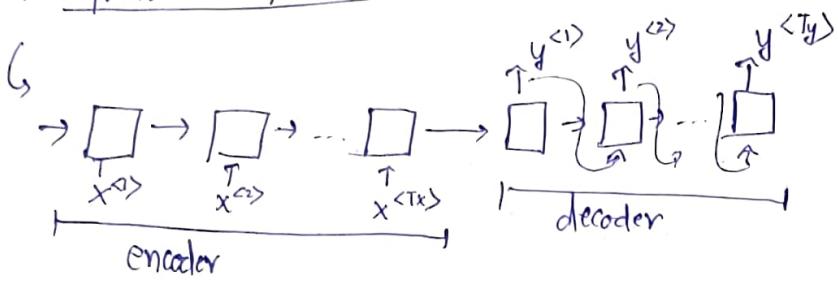
(for e.g. for gender, don't use words like grandpa, grandma and use words like doctors, scientist)

③ Equalize pairs: ie make the pairs of words only separated by gender to be equidistant from non-gender words.

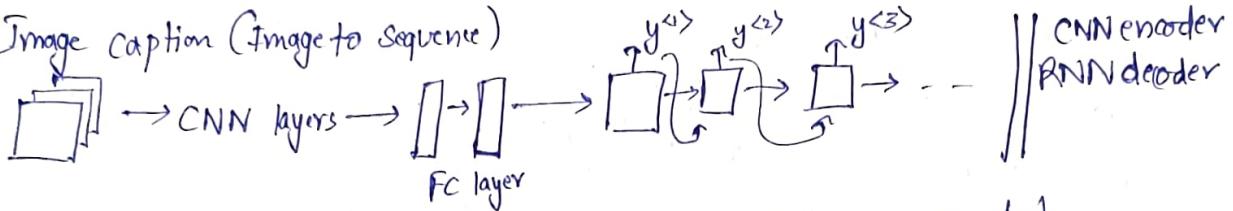
To know which words should be passed through the neutralization step (ie which words are not gender specific) a linear classifier can be used.

Sequence to Sequence architectures

(59)



↳ Image caption (Image to sequence)



↳ Machine learning translation as building a conditional language model

$$P(y^{<1>} \dots y^{<Ty>} | x^{<1>} \dots x^{<Tx>}) \leftarrow \text{as we have to find } y^{<1>} \dots y^{<Ty>} \text{ given } x^{<1>} \dots x^{<Tx>}$$

To find the most likely translation we find: $\arg \max_{y^{<1>} \dots y^{<Ty>}} P(y^{<1>} \dots y^{<Ty>} | x)$

* Why not greedy search?

ie. picking $y^{<i+1>}$ after picking $y^{<i>}$ and not picking all the $y^{<i>}$ simultaneously which maximizes $P(y^{<1>} \dots y^{<Ty>} | x)$

because this does not optimize the probability of having the whole line towards the correct translation.

because here unlike in language model where we get the most-probable sentence based on only the previous words, we want the most-probable language translation (the sentence as a whole)

Difficulty: very large number of probable sentences in the samplespace

$(\# \text{ of words})^{\# \text{ of words in sentence}}$ in the dictionary

Beam Search

* Beam width: # of possibilities to consider to obtain the next word. (let = b)

Choosing 1st word: get the top b # of possibilities based on $P(y^{<1>} | x)$

Choosing 2nd word: get the top b # of possibilities based on $P(y^{<1>} \dots y^{<2>} | x)$
 $(\text{find } P(y^{<1>} \dots y^{<2>} | x) = P(y^{<1>} | x) P(y^{<2>} | x, y^{<1>}))$

Choosing ith word: get the top b # of possibilities based on $P(y^{<1>} \dots y^{<i>} | x)$
 $(\text{find } P(y^{<1>} \dots y^{<i>} | x) = P(y^{<1>} \dots y^{<i-1>} | x) P(y^{<i>} | x, y^{<1>} \dots y^{<i-1>}))$

of probabilities to check = # of words in dictionary

of probabilities to check = (# of words in dictionary) $\times b$

same here.

NT: $b=1 \Rightarrow$ greedy search

→ Normalized log-likelihood objective (length normalization)

↳ Two problems while calculating $\arg \max_y \prod_{i=1}^{Ty} P(y^{<i>} | x, y^{<1>}, \dots, y^{<i-1>})$

① as probabilities are pretty small (large # of vocabulary words) multiplying these probabilities can lead to very small #'s, smaller than floating point precision of computer

↳ Remedy: → take log-likelihood.

② as longer sentence \Rightarrow more # of multiplications \Rightarrow smaller probabilities

\Rightarrow the algorithm supports shorter lengths of output even though they might be wrong.

↳ Remedy: → normalized log-likelihood objective

$$\arg \max_y \frac{1}{Ty^\alpha} \sum_{i=1}^{Ty} \log P(y^{<i>} | x, y^{<1>}, \dots, y^{<i-1>}) \quad || \begin{array}{l} \alpha \rightarrow \text{hyperparameter} \\ \in (0, 1) \\ \text{of soft normalization} \end{array}$$

* for beam search, calculate the normalized log-likelihood objective of all the cases found and take the one with the maximum.

→ Error Analysis for beam search

↳ For a machine translation, given x and its human translation y^* and its machine translation \hat{y} how to know which is causing the error? (RNN or Beam search)

↳ Using the RNN, calculate $P(\hat{y}|x)$ and $P(y^*|x)$ {or calculate the corresponding log-likelihood objective}

↳ Case-1: $P(y^*|x) > P(\hat{y}|x) \Rightarrow$ Beam search is at fault

↳ beam search chose \hat{y} . But y^* attains higher $P(y^*|x)$.

↳ Case-2: $P(y^*|x) < P(\hat{y}|x) \Rightarrow$ RNN is at fault

↳ look into all the mislabelled examples and get the fraction of errors due to beam search. If it is in majority, solve this error by increasing beam width.

→ Evaluating Machine Translation (Bleu Score)

(Bilingual Evaluation Understudy)

↳ Multiple different translations can be a good translation for a given text

→ the bleu score is calculated by comparing the machine translation (MT) output with all the reference outputs (more than one) made by humans for a given input. (61)

bleu score on unigrams (one-word)

e.g.:

Ref-1: The cat is on the mat
Ref-2: There is a cat on the mat
MT: the the the the the the the

$$\text{precision} := \frac{7}{7} = \frac{\# \text{ of words in MT that appear in either ref}}{\# \text{ of words in MT}}$$

(not a good measure)

$$\text{modified precision} := \frac{2}{7} = \frac{\sum_{w \in \text{unique words in MT}} \text{Count}_{\text{clip}}(w)}{\# \text{ of words in MT}}$$

$\text{Count}_{\text{clip}}(w) := \max \text{ count of word } w \text{ in any of the reference}$

bleu score on bigrams (pairs-of-words)

e.g. Ref=1 and Ref=2

MT: the cat the cat on the mat.

bigrams	Count	Count _{clip}	modified precision
the cat	2	1	
cat the	1	0	
cat on	1	1	
on the	1	1	
the mat	1	1	

$$= \frac{\sum \text{Count}_{\text{clip}}}{\sum \text{Count}}$$

$$= \frac{4}{6}$$

Bleu Score on n-grams

$$P_n = \frac{\sum_{n\text{-gram} \in y} \text{Count}_{\text{clip}}(n\text{-gram})}{\sum_{n\text{-gram} \in y} \text{Count}(n\text{-gram})}$$

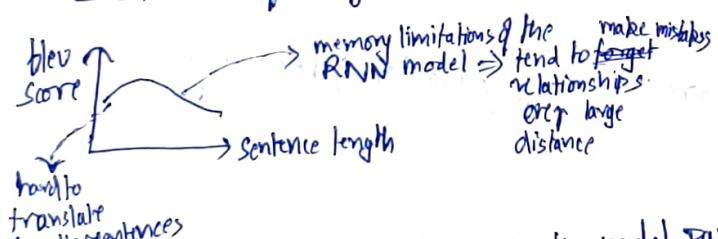
→ Combined Bleu score $\Rightarrow \text{BP} \exp \left[\frac{1}{k} \sum_{n=1}^k P_n \right]$

Brevity Penalty (BP)

$$= \begin{cases} 1, & \text{MT_output_length} \geq \text{ref_output_length} \\ \exp \left(1 - \frac{\text{MT_output_length}}{\text{ref_output_length}} \right), & \text{MT_output_length} < \text{ref_output_length} \end{cases}$$

used to penalize shorter length MT output which may tend to have larger bleu score as they might be having all the words already in reference, but the translation overall might be very wrong.

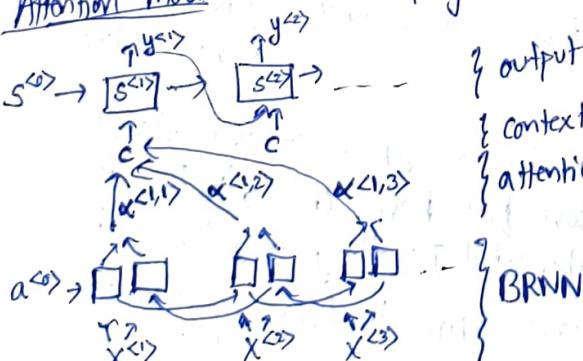
→ The Problem of long sequences



solved by attention model

Attention model

to make the model pay attention to only a part of the input sequence to generate the output.



$\alpha^{(1,1)}, \alpha^{(1,2)}, \alpha^{(1,3)}$ \rightarrow depends on previous words i requires to derive word j in the output

$$\cdot a^{(t)} = (\vec{a}^{(q)}, \vec{a}^{(v)})$$

↳ attention and context: $\sum_{t'} \alpha^{(q,t')} = 1$ (all the attentions are in $[0,1]$)

$$C^{(q)} = \sum_{t'} \alpha^{(q,t')} a^{(t')}$$

(62)

↳ Computing attention $\alpha^{(q,t')} = \text{amount of attention } y^{(t')} \text{ should pay to } a^{(t')}$

$$\alpha^{(q,t')} = \frac{\exp(e^{(q,t')})}{\sum_{t=1}^{T_x} \exp(e^{(q,t')})}$$

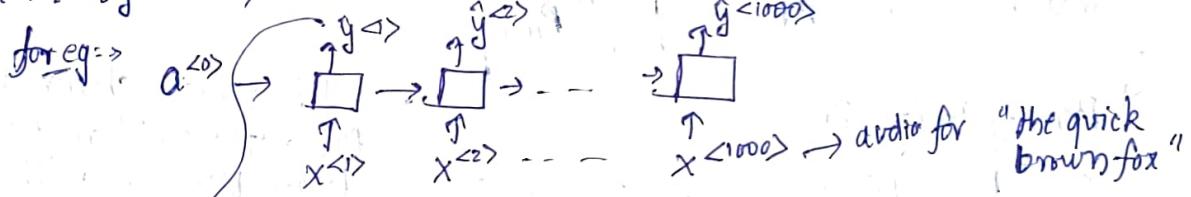
* Softmax used here to ensure $\alpha^{(q,t')} \in [0,1]$

$e^{(q,t')} \rightarrow \begin{matrix} s^{(q,t')} \\ a^{(q,t')} \end{matrix} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow e^{(q,t')}$

NT → Attention models have quadratic time complexity $\rightarrow O(T_x T_y)$

→ Speech Recognition

↳ $T_x \gg T_y$ here, so need to condense the output to the original length



CTC
(Connectionist
Temporal
Classification)

Basic rule: → Collapse repeated character not separated by blank.

→ Transformers

↳ sequence to sequence models which are performing parallel processing.

like CNN models.

→ Self-Attention

$A(q, k, v) = \text{attention-based vector representation of a word.}$

$q \rightarrow \text{query}$
 $k \rightarrow \text{key}$
 $v \rightarrow \text{value}$

Transformer Attention

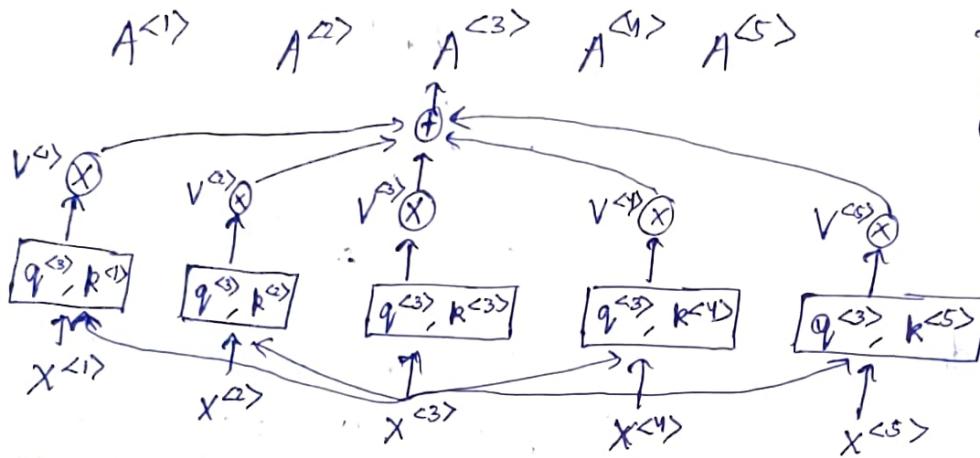
$$A(q, k, v) = \sum_i \frac{\exp(q \cdot k^{(i)})}{\sum_j \exp(q \cdot k^{(j)})} v^{(i)}$$

→ Similar to RNN attention discussed above but here the attention is calculated parallelly (for all the words) whereas RNN attentions are calculated one after another.

→ here, $q^{<i>} = W^Q \cdot x^{<i>}$ → query to enquire about another input
 $k^{<i>} = W^K \cdot x^{<i>}$ → key related to the input vector
 $v^{<i>} = W^V \cdot x^{<i>}$ → value related to the input vector

query, key, value similar to database concepts
here W^Q, W^K, W^V are learnable parameters.

→ Self attention example:-



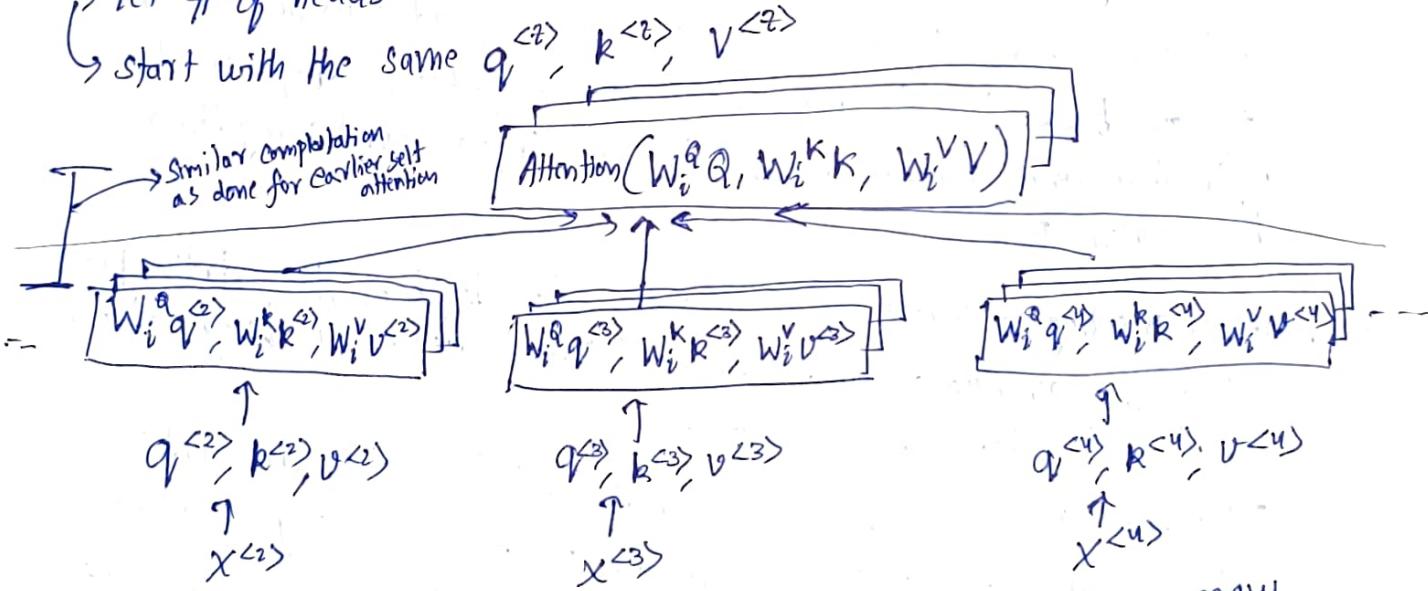
gives a more advanced way of representing every word which also has appropriate attention from nearby words in contrast to word embeddings

→ Multi-head Attention

the self-attention calculated earlier was single-head.

let # of heads = h

start with the same $q^{<2>}, k^{<2>}, v^{<2>}$



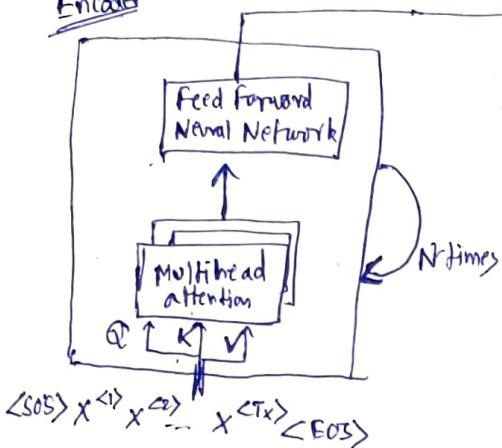
The $\{W_i^Q, W_i^K, W_i^V, i \in [1, h]\}$, each head ($i=i$) try to define a new feature to the self-attention connection { for e.g. if $x^{<3>} \rightarrow$ Africa then $W_1^Q, K, V \rightarrow$ where in Africa
 $W_2^Q, K, V \rightarrow$ when in Africa, etc. }

$$\text{Multihead}(Q, K, V) = \text{Concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h) W_o$$

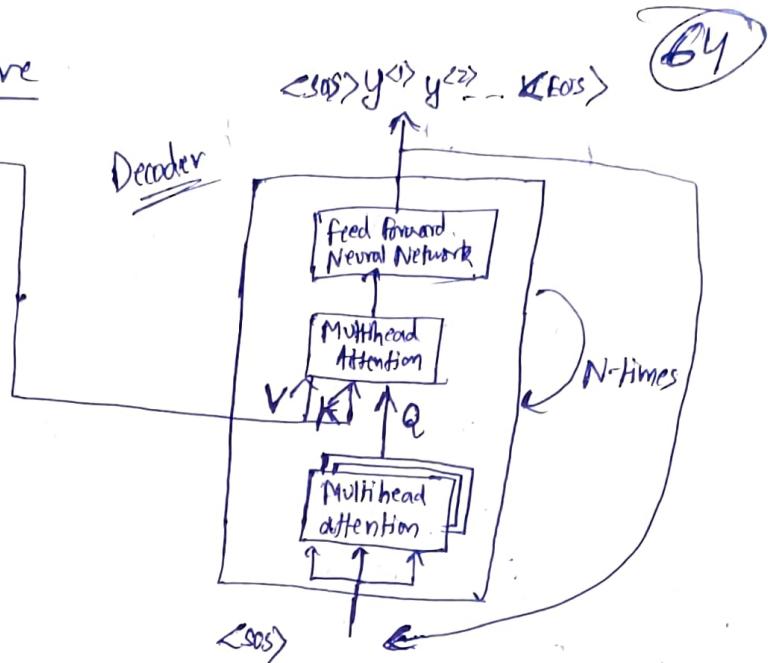
$$\text{head}_i = \text{Attention}(W_i^Q Q, W_i^K K, W_i^V V)$$

→ Transformer Network Architecture

Encoder



Decoder



(64)

NT →

the input to the decoder block starts with $\langle \text{SOS} \rangle$ and then the next input is ~~given~~ taken from the output of the decoder block in the previous run.

NT →

the K and V passed from the encoder contains the info for each input word in context of different features. The Q generated from each of synthesized word till now. So, the 2nd multihead attention in the decoder block tries to query (related to each synthesized output word) information and connections related to each input word using their key, value pair.

NT → Positional encoding :→ while calculating the self attention, the relative position of the queried word to the word whose attention is being calculated, is not taken into consideration.

To take positions of the words (relative) into consideration while calculating self attention, a unique positional encoding is added to each $x^{<i>}$, $y^{<i>}$ given as input to both of the encoder and decoder.

Let $x^{<i>}$ be d-dimensional $x^{<i>} = [\underbrace{-}_{i=0}, \underbrace{-}_{i=1}, \underbrace{-}_{i=2}, \dots, \underbrace{-}_{i=d-1}]$

Corresponding d-dimensional positional encoding is { pos = position } of the word }

$$\text{PE}(\text{pos}, z_i) = \sin \left[\frac{\text{pos}}{(10000)^{2i/d}} \right]$$

$$\text{PE}(\text{pos}, z_{i+1}) = \cos \left[\frac{\text{pos}}{(10000)^{2i/d}} \right]$$