

Deep Learning

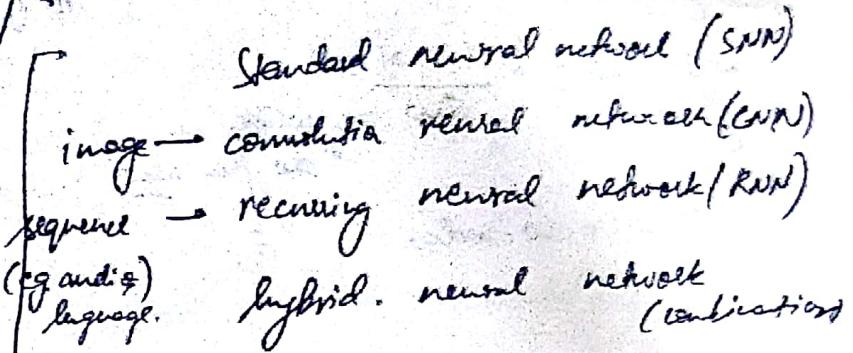
Neural network

info output

(x, y) supervised learning

Structured data → Able, tabular form

Unstructured data → image, audio, text



$m \rightarrow$ training set (labelled data).

Binary Classification $n_x (0 or 1)$

Notation: $(x, y) \quad x \in \mathbb{R}^n \quad y \in \{0, 1\}$ m_{train}, m_{test} .

$$X = \mathbb{R}^{n \times m} \quad Y = \mathbb{R}^{1 \times m} \quad X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}^T \quad Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

Logistic Regression

Given X want $\hat{y} = P(y=1|X)$

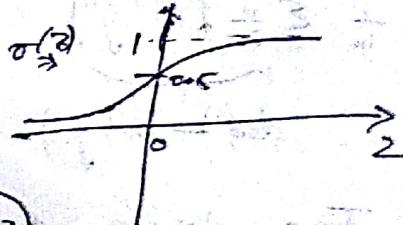
$X \in \mathbb{R}^{n \times m} \Rightarrow$ matrix. $(m \times n)$

Parameters: $w \in \mathbb{R}^{n \times 1}$, $b \in \mathbb{R}$ \Rightarrow matrix $(m \times 1)$

Output.

$$\hat{y} = \sigma(w^T n + b)$$

we want $\hat{y}^{(i)} \approx y^{(i)}$



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

Cost function $J(w, b)$, Loss(error) function: $L(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$

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

$$= \frac{1}{m} \sum_{i=1}^m (y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)}))$$

↑
for single
exmpl.

$$\text{But we use } L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

so $y=1 \rightarrow$ we want \hat{y} large
 $\& y=0 \rightarrow$ we want \hat{y} small

for large (\hat{y} so average)

Gradient descent \Rightarrow find w, b that minimise

(Take any random point, move towards optimum, finally reach.)

Global optimum

$$J(w, b)$$

(Convex so works)

$$\text{Repeat } \left\{ \begin{array}{l} \text{learning rate} \\ w := w - \alpha \frac{\partial J(w)}{\partial w} \end{array} \right. \rightarrow \begin{array}{l} \text{written as} \\ "ds" \text{ from now} \end{array}$$

To $w := w - \alpha \partial w$.

$$\text{to find } \nabla J(w, b). \quad \text{due}$$

$$w := w - \alpha \frac{\partial J(w, b)}{\partial w}$$

$$b := b - \alpha \frac{\partial J(w, b)}{\partial b} \quad \text{due}$$

$$\begin{aligned}
 & z = w_1 n_1 + w_2 n_2 + b \rightarrow a = \sigma(z) \rightarrow f(x, y) \\
 & \frac{dz}{d\theta} = \left(\frac{\partial L}{\partial n_1}, \frac{\partial L}{\partial n_2} \right) \Rightarrow a(\theta) \\
 & \frac{\partial L}{\partial n_1} = \frac{\partial dL}{\partial n_1} = \frac{da}{dn_1} = \frac{1}{a} + \frac{1}{1-a} \\
 & \frac{\partial L}{\partial n_2} = \frac{\partial dL}{\partial n_2} = \frac{da}{dn_2} = \frac{1}{a} - \frac{1}{1-a} \\
 & \frac{\partial L}{\partial w_1} = "dw_1" = n_1 \cdot dz \quad dw_1 = n_1 \cdot dz \quad db = dz \quad w_1 := w_1 - \alpha dw_1 \\
 & \frac{\partial L}{\partial w_2} = "dw_2" = n_2 \cdot dz \quad dw_2 = n_2 \cdot dz \quad w_2 := w_2 - \alpha dw_2 \\
 & b := b - \alpha db
 \end{aligned}$$

Logistic regression on m examples

$$J(\omega, b) = \frac{1}{m} \sum_{i=1}^m \ell(a^{(i)}, y^{(i)}) \quad \hat{a}^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(\omega^\top x^{(i)} + b)$$

$$\frac{\partial}{\partial \omega_i} J(\omega_0) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial \omega_i} f(d^{(i)}, y^{(i)})}_{\text{d}\omega_i - (x^{(i)}, y^{(i)})}$$

$$\text{algo} \Rightarrow J=0; dw_1=0; dw_2=0; db=0$$

For i=1 to m

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

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

$$J_4 = -[g'' \log(g') + (1-g') \log(1-g')]$$

$$dz^{ij} = d\tilde{z}^{ij} - \tilde{g}^{ij}$$

$$d\varphi_1 = \pi_1^* df_1$$

$$\deg_2 + = n_2^{(ii)} \alpha_2^{(i)}$$

$$d\varphi_2 = \pi_2^{(ii)} d\varphi_1^{(i)} \quad |_{n=2} \\ d\theta = df^{(i)} \quad |_{\text{if } n \text{ more than loop.}}$$

$$db + = df^{-1}$$

卷之三

$\vdash \text{dis}_2 (= m)$

J / = m

$$d\omega_1 \wedge \dots \wedge d\omega_k \neq 0$$

$$\omega_i := \omega_i - \alpha \text{diag}$$

$$\omega_2 := \omega_2 - \alpha \cdot \text{d}\omega_2$$

$$b := b - \alpha db$$

* (2) for loops

→ (we have to reduce
for loops)

(The Vegetation is well
soil developed.)

Vectorization:

$$z = w^T n + b$$

(np - many)

Non-vectorized

$$z = 0$$

for i in range(0-n):

$$z += w[i] * n[i]$$

$$z += b$$

vectorized (very fast):

$$z = np.\text{dot}(w, n) + b$$

w, n - array.

$$v = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_m \end{bmatrix} \rightarrow u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_m} \end{bmatrix}$$

$$u = np.\exp(v)$$

np. log(v)

np. abs(v)

np. softmax():

row normalize with hadamard.
C^n labeled states.

np.sqrt

np. maximum()

VK & Z (small)

1/V (inverse
of softmax)

Thus previous one, ~~the~~ the second for log can also be removed.

$$X = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ 1 & 1 & \dots & 1 \end{bmatrix} \in R^{n \times m} \quad Z = \begin{bmatrix} z^{(1)} \\ z^{(2)} \\ \vdots \\ z^{(m)} \end{bmatrix} = w^T X + [b \ b \ \dots \ b]^T$$

so $Z = np.\text{dot}(w.T, X) + b$
contains all $z^{(1)} \dots z^{(m)}$ terms. $\in (n \times m)$

$$A = \begin{bmatrix} a^{(1)} \\ a^{(2)} \\ \vdots \\ a^{(m)} \end{bmatrix} = \sigma(z)$$

$$\text{Now } dZ = [dZ^{(1)} \ dZ^{(2)} \ \dots] \quad \text{so } \frac{dZ}{dZ} = A - Y$$

$$\begin{aligned} dw &= 0 \\ dw &+ x^{(1)} dZ^{(1)} \\ &\vdots \\ dw &+ x^{(m)} dZ^{(m)} \end{aligned}$$

$$dZ^{(1)} = y^{(1)} - a^{(1)}$$

$$db = 0$$

$$db + dZ^{(1)}$$

$$db + dZ^{(2)}$$

$$\vdots$$

$$db + dZ^{(m)}$$

$$d\theta = \frac{1}{m} np.\text{sum}(dZ)$$

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

New logistic regression:

$$Z = w^T X + b$$

$$= np.\text{dot}(w.T, X) + b$$

$$A = \sigma(Z)$$

$$dZ = A - Y$$

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

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

$$w_{\theta} = w - \alpha dw$$

$$b_{\theta} = b - \alpha db$$

Now put a for loop
for (suppose 100 iterations)
to reach the min. point
that will be needed.

Broadcasting in python :-

$$\text{cal} = A \cdot \text{sum}(\text{axis}=0) \rightarrow [6 \quad 15 \quad 24] \quad \begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{bmatrix}$$

$$\text{percentage} = 100 * \frac{A}{\text{cal.reshape}(1,3)} \quad \text{so, percentage} = \frac{100 * A}{\text{cal}} \quad (1,3)$$

// each row
new).calculate. Broadcasting this is not needed.
as it is originally
(1,7) only.

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

$$; \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

$$; \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

Broadcasting
so applicable for + - * / .

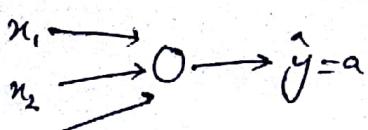
a = np.random.rand(5) \rightarrow rank 1 array. \rightarrow don't use.

use a = np.random.rand(5, 1) \rightarrow column vector use reshape to convert always

Tip: x.reshape(x.shape[0], -1).T \Leftrightarrow (a, b, c, d) \rightarrow (b*c*d, a)

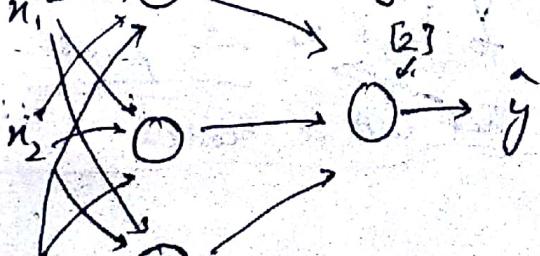
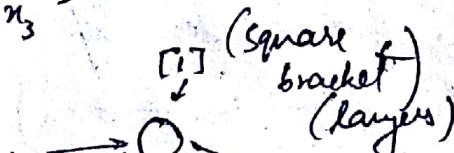
b is last c is last d.
size no juega b*c*d.

Neural Network.



$$z = w^T x + b \rightarrow a = \sigma(z) \rightarrow L(a, y)$$

w z da



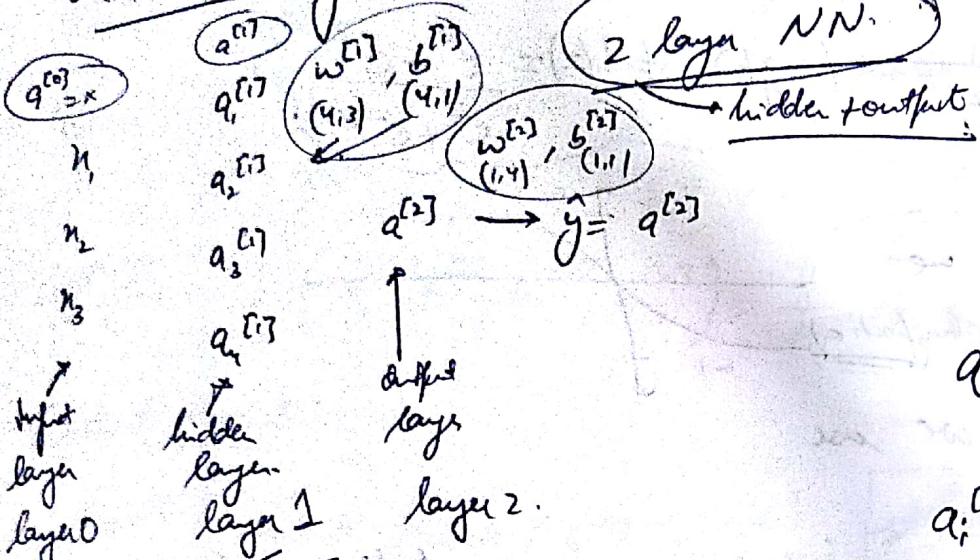
$$L(a^{[1]}, y) \leftarrow a^{[2]} = \sigma(z^{[2]})$$

$$w^{[1]} \rightarrow z^{[1]} = w^{[1]} x + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \rightarrow z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$\frac{\partial L}{\partial a^{[1]}}$ $w^{[2]}$ $b^{[2]}$ $\frac{\partial L}{\partial a^{[2]}}$

$$dL^{[2]}$$

One hidden layer.



$$Z = w^T x + b$$

$$\hat{y} = a = \sigma(z)$$

linear activation is considered as a single layer. σ is ReLU.

$a^{[l]}$ layer.

$a_i^{[l]}$ $i \leftarrow$ node in layer.

$$a_i^{[l]} = \sigma(z_i^{[l]}) = \sigma(w_i^{[l]T} x + b_i^{[l]})$$

layer 2

$$x \rightarrow a^{[2]} = \hat{y}$$

$a^{[2]}(i)$

a example i

$$x^{(m)} \rightarrow a^{[2](m)} = \hat{y}^{(m)}$$

⇒ How not to use

loop

Now loop mean

Put in data.

for 1 to m

$[2]^{[1]}(i)$

multiple examples

$$\text{Basically, } Z^{[1]} = \begin{bmatrix} | & | & | & | \\ z^{1} & z^{[1](2)} & z^{[1](3)} & z^{[1](m)} \\ | & | & | & | \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} | & | & | \\ a^{1} & a^{[1](2)} & a^{[1](m)} \\ | & | & | \end{bmatrix}$$

$$Z^{[1]} = w^{[1]} X + b^{[1]}$$

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

$$Z^{[2]} = w^{[2]} A^{[1]} + b^{[2]}$$

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

corresponds to different training examples.

different nodes
in neural
network.

$$w^{[1](i)} x^{(i)} + b^{[1]} = Z^{[1](i)}$$

(No w, b for different set of training examples do.)

$$\text{Basically, } Z^{1} = w^{1} x^{(1)} + b^{[1]} ; Z^{[1](2)} = w^{[1](2)} x^{(2)} + b^{[1]} ; Z^{[1](3)} = w^{[1](3)} x^{(3)} + b^{[1]}$$

Activation functions

Example: sigmoid ($\sigma(z)$)

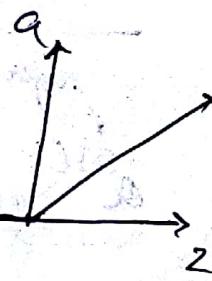
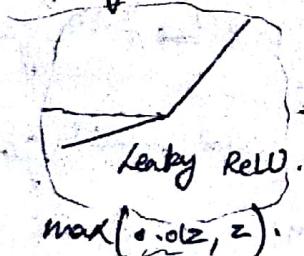
$$\frac{1}{1+e^{-z}}$$

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

when $y \in [0,1]$ so $\hat{y} \in [0,1]$ we can use sigmoid. (binary classification)

* instead of sigmoid if we use $g(z)$ it works better.

Other functions:



$a = \max(0, z)$
ReLU (rectified)
linear unit

Max people use it.
(as $g(z)$ near +1, the slope is very less)

* $y \in \mathbb{R}$, then we can use linear function, otherwise we will use non-linear only.

Non-linear Activation function is a critical part.

$$a = g(z) = \frac{1}{1+e^{-z}} \Rightarrow g(z) = g(z)(1-g(z)) = a(1-a) \quad [\text{sigmoid}]$$

$$a = g(z) = \tanh z = \frac{e^z - e^{-z}}{e^z + e^{-z}} \Rightarrow g'(z) = 1 - (\tanh(z))^2 = 1 - a^2 \quad [\tanh]$$

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

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

Gradient descent for neural networks.

Param: $\omega^{(1)}, b^{(1)}, \omega^{(2)}, b^{(2)}$
 $(n^{(1)}, n^{(2)})$ $(n^{(1)}, 1)$ $(n^{(2)}, n^{(1)})$ $(n^{(2)}, 1)$

$$\text{Cost } f : J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i) \quad \begin{cases} z^{[1]} = w^{[1]} A^{[1]} + b^{[1]} \\ A^{[2]} = g^{[2]}(z^{[1]}) \sigma(z^{[1]}) \end{cases}$$

Gradient descent :

Repeat { Compute predicts ($\hat{y}^{(i)}$, $i=1\dots n$)
 $d\omega^{(i)} = \frac{\partial J}{\partial \omega^{(i)}}$, $db^{(i)} = \frac{\partial J}{\partial b^{(i)}}$ }

$$\omega^{[i]} := \omega^{[i]} - \alpha d\omega^{[i]}$$

$$\begin{aligned} w^{[i,j]} &:= w^{[i,j]} - \alpha d w^{[i,j]} \\ b^{[i,j]} &:= b^{[i,j]} - \alpha d b^{[i,j]} \end{aligned} \quad \boxed{\text{deriving}}$$

$$\omega^{[i]} = \omega^{[i]} - \alpha d\omega^{[i]}$$

$$dZ^{[1]} = \underbrace{w^{[1]T} dZ^{[2]}}_{(n^{[2]}, m)} * \underbrace{g^{[1]'}(z^{[1]})}_{\text{element wise.}} \quad (n^{[1]}, 1)$$

$$dw^{[i]} = \frac{1}{m} dz^{[i]} x^T$$

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

$$\hookrightarrow (n^{(1)}, 1)$$

~~#~~ backward propagation:

$$dZ^{[2]} = A^{[2]} - Y \quad \text{and} \quad dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$d\omega^{E,T} = \frac{1}{m} d\sum^{E,T} A^{E,T}$$

$$db^{[2]} = \frac{1}{m} \text{ np.sum } (dz^{[1]}, \underbrace{\text{axis}=1, \text{ keepdims=True}}_{\text{if } dz^{[1]} \text{ is a column vector}}$$

$$\rightarrow (n^{ij}, 1)$$

Rank 1 array
not one is large.
 $(n^{1/2}, e)$

6 imp for
back propagation

Backpropagation

The diagram illustrates a neural network layer with three neurons. The forward pass starts with input x entering neuron 1. The output of neuron 1 is $a^{[1]} = \sigma(z^{[1]})$, where $z^{[1]} = w^{[1]}x + b^{[1]}$. This value is passed to neuron 2 as input. The output of neuron 2 is $a^{[2]} = \sigma(z^{[2]})$, where $z^{[2]} = w^{[2]}a^{[1]} + b^{[2]}$. The output of neuron 2 is then passed to neuron 3 as input. The output of neuron 3 is $a^{[3]} = \sigma(z^{[3]})$, where $z^{[3]} = w^{[3]}a^{[2]} + b^{[3]}$.

Backpropagation is shown with arrows indicating gradients:

- $\delta z^{[1]} = da^{[1]} / dw^{[1]}$
- $\delta a^{[1]} = da^{[1]} / db^{[1]}$
- $\delta z^{[2]} = da^{[2]} / dw^{[2]}$
- $\delta a^{[2]} = da^{[2]} / db^{[2]}$
- $\delta z^{[3]} = da^{[3]} / dw^{[3]}$
- $\delta a^{[3]} = da^{[3]} / dy$

A handwritten note at the bottom left says "denecke wir produziert".

Initialization

w to random, b to 0 value.
but due to all 0 data to have
some use usage - problem

$$w^{[1]} = \text{np.random.rand}(2, 2) * \frac{1}{\sqrt{2}}$$

$$b^{[1]} = \text{np.zeros}(2, 1)$$

$$w^{[2]} = \dots$$

$$b^{[2]} = 0$$

Deep Neural Network

4 hidden layers

4 layer NN \rightarrow 3 hidden layers.

$L = 4$ as $n^{[L]} = \# \text{ units in layer } L$.
 $(\# \text{ layers})$ so $n^{[1]} = 5$; $n^{[2]} = 5$; $n^{[3]} = 3$; $n^{[4]} = 1$
also $n^{[0]} = n_x = 3$

$a^{[l]}$ = activations in layer l . \hat{y} = $g^{[L]}(z^{[L]})$ $w^{[l]}$ = weights for $z^{[l]}$

basically; $z^{[k]} = w^{[k]} a^{[k-1]} + b^{[k]}$; $a^{[k]} = g^{[k]}(z^{[k]})$; $a^{[4]} = \hat{y}$.

ad ~~case~~ case $a \rightarrow A$, etc for vertically.

ad along along layers to i.e. for loop use inner product.

Dimensions: $w^{[k]} \Rightarrow (n^{[k]}, n^{[k-1]})$

~~dimensions~~ $a^{[k]} \Rightarrow (n^{[k]}, m)$

Same as z .

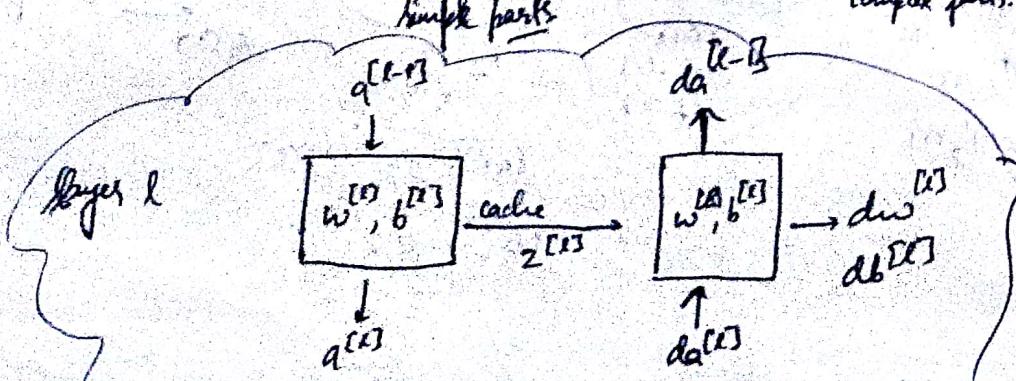
$b^{[k]} \Rightarrow (n^{[k]}, 1)$

$db^{[k]} \Rightarrow (n^{[k]}, m)$ (raise broadcasting to same width along)

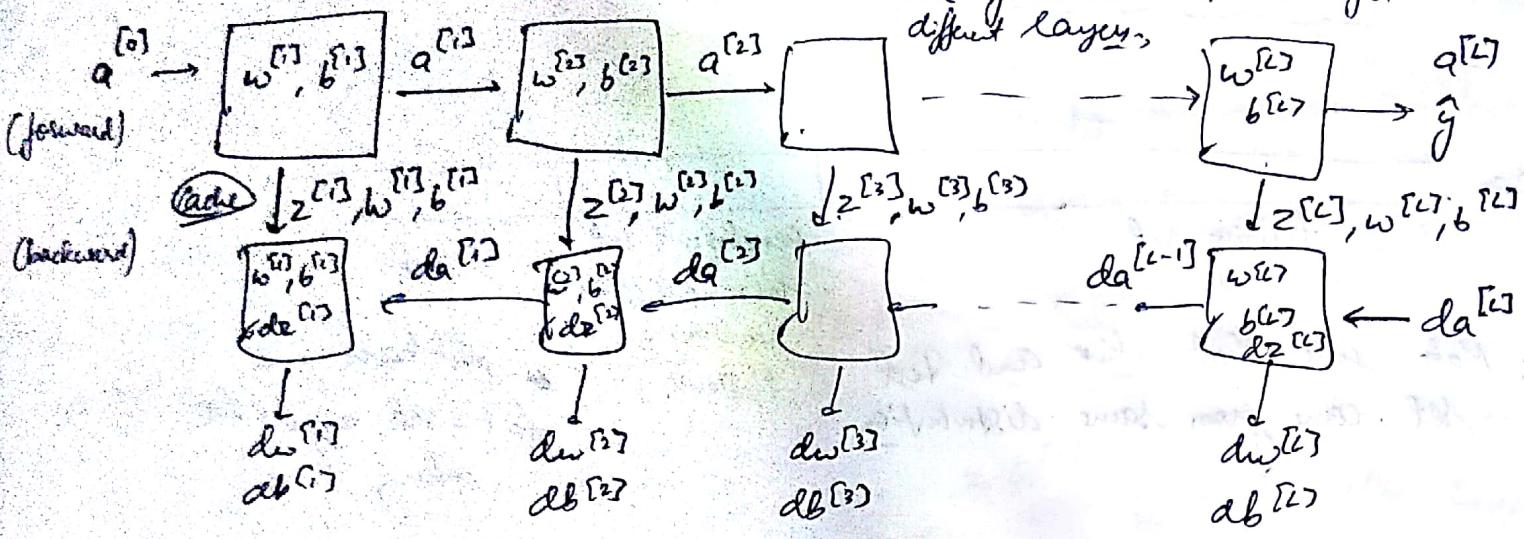
Different levels \rightarrow first low level, then high level recognition

of simple parts

complex parts.



In each layer, there is forward propagation step & backward propagation step.
different activation functions also possible for different layers.



add $z^{[i]}$ like below manner

be w^{T^3} , b^{T^3} the area for

$$\omega = \omega - \alpha d\omega$$

$$b^{[k]} = b^{[k]} - \alpha d b^{[k]}$$

FORWARD

\rightarrow Output $a^{[L-1]}$ $w^{[L]}, b^{[L]}$

→ Output $\alpha^{[k]}$, Cache $z^{[k]}$

$$Z^{[x]} = w^{[x]} \cdot A^{[x-1]} + b^{[x]}$$

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

Backward

→ Input da [P3]

→ Output $da^{[L-1]}, dw^{[L]}, db^{[L]}$

$$dZ^{[e]} = dA^{[e]} * \overleftarrow{g^{[e]}}'(Z^{[e]})$$

$$dW^{[k]} = \frac{1}{\gamma} dZ^{[k]} \cdot A^{[k-1]T}$$

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

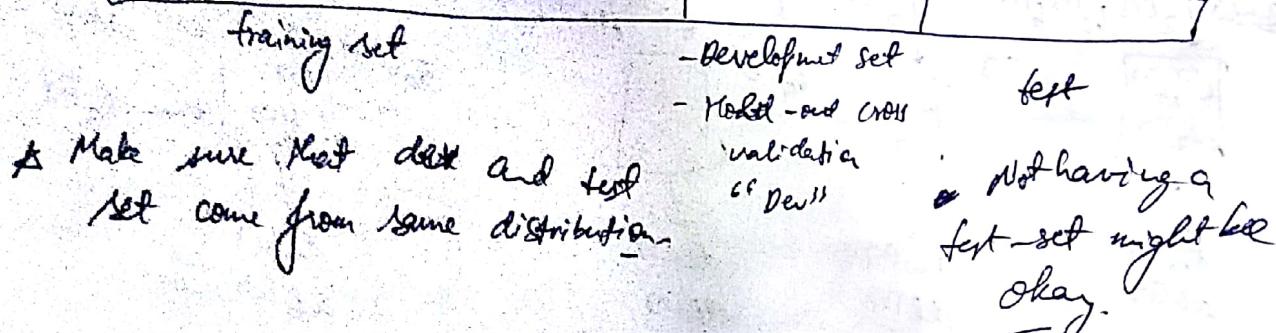
$$dA^{[L-1]} = W^{[L]T} \cdot dz^{[L]}$$

$$dA^{[2]} = \left(-\frac{y^{(1)}}{a^{(1)}} + \frac{(-y^{(1)})}{1-a^{(1)}} \dots \frac{y^{(m)}}{a} \right)$$

Hyperparameters: L , # iterations, ηL , $\#_{\text{RIB}}, \#_{\text{CIT}}$. -> Choice of activation f^k , etc.
 (we have to adjust them ourselves and check the performance).

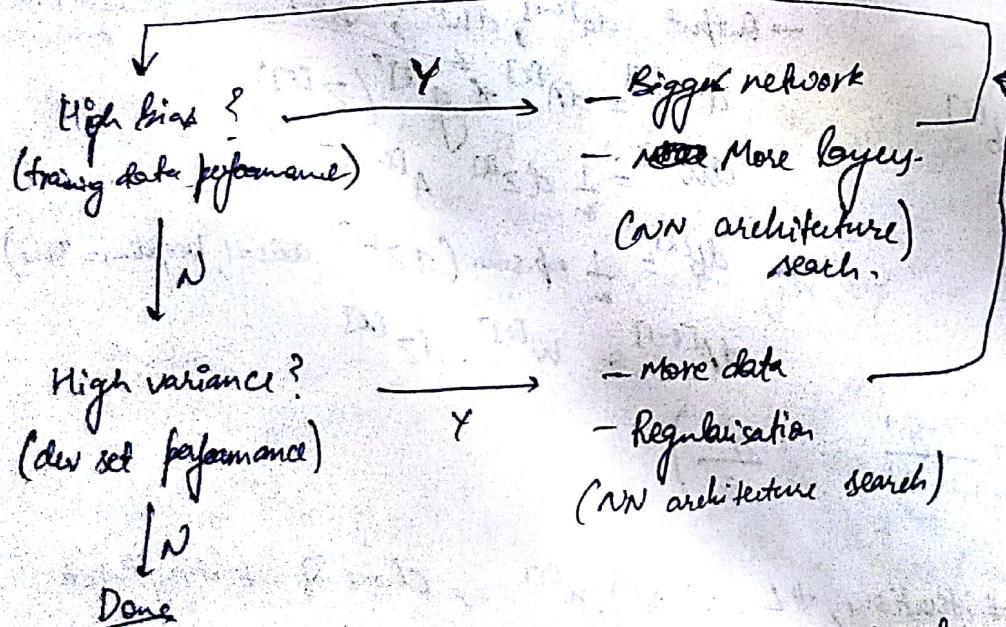
IMPROVING DEEP NEURAL NETWORKS.

Data:



underfitting → high bias] → (train set error vs dev set error = relatively low.)
 overfitting → high variance] → higher dev set error. (high.)

high bias + high variance → high train set error + much higher dev set error.



Dont use if not overfitting.

Regularization? (prevents overfitting)

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

regularization parameter. λ also a hyperparameter.

regression: $\lambda \|w\|_1 = \frac{1}{2m} \sum_{i=1}^m |w_i|$

classification: $\lambda \|w\|_2^2 = \frac{1}{2m} \sum_{i=1}^m w_i^2$

also a hyperparameter.

$$J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

$$\|w^{[l]}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^n (w_{ij}^{[l]})^2 \quad w \in (\mathbb{R}^m, \mathbb{R}^{n \times 1})$$

(Frobenius norm)

$$\text{Now, } dw^{[l]} = (\text{from backprop}) + \frac{\lambda}{m} w^{[l]} \Rightarrow w^{[l]} := \overset{\longrightarrow}{w^{[l]} - \alpha \left[(\text{from backprop}) + \frac{\lambda}{m} w^{[l]}\right]}$$

$$= w^{[l]} - \frac{\lambda}{m} w^{[l]} - \alpha (\text{from backprop})$$

$$w^{[l]} := \left(1 - \frac{\lambda}{m}\right) w^{[l]} - \alpha (\text{from backprop}).$$

① Dropout regularization: note, makes no backprop loss.

Keep-prob ≤ 0.8

② Inverted dropout:

$d_3 = \text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1]) < \text{keep-prob. prob.}$

→ also works for vectorisation → so 80% of being 1, 20% of being 0

dropout vector for layers and $a_3 = \text{np.multiply}(a_3, d_3)$ # jacks rank into order.

$a_3/ = \text{keep-prob}$ # scaling it up.

→ to keep expected value to remain same.

To zero-out some hidden units.

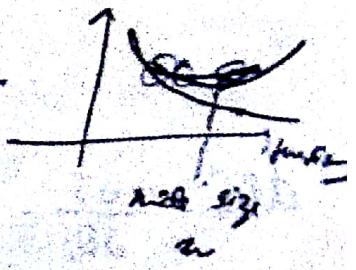
Keep-prob = 1.0 → keeping every unit (no dropout).

always for input keep-prob can vary for every layer.

③ Data augmentation: same image to invert, distort for more input

④ Early stopping: plot training error, dev set error also.

not good. → Prefer to use L_2 than this.
↓ better.



Normalising training set: Basically variance in my equal. 1 or turn see lie.

$$\text{Subtract mean} \Rightarrow \mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$x := x - \mu$$

$$\text{normalise variance} \Rightarrow \sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)T} x^{(i)} \quad \begin{matrix} \leftarrow \\ \text{client wise} \end{matrix}$$

$$x' (= \sigma^2)$$

Use same my σ^2 to normalise test set.

Used so that features come on almost similar scales.

Cost function changes easier and in better way.

Vanishing/exploding gradients: In case of very deep neural networks $L \gg 100$.

The gradients can be very large or very small, thus causing a problem.

Also consider prone to lie, there are various research papers.

like if w be random initialisation then multiply by

$\sqrt{\frac{2}{n^{(L-1)}}}$ when using ReLU function, or multiply by

$\sqrt{\frac{2}{n^{(L-1)} + n^{(L)}}}$ for $w^{(L)}$ when using tanh, etc.

for derivative computation: use $f'(0) = \lim_{\epsilon \rightarrow 0} \frac{f(0+\epsilon) - f(0-\epsilon)}{2\epsilon}$ much more accurate.

$$\text{rather than } \frac{f(0+\epsilon) - f(0)}{\epsilon}$$

$$d\theta[i] = \frac{\partial J}{\partial \theta_i} = \underline{J(\theta_0, \theta_1, \dots, \theta_i + \epsilon, \dots)} - \underline{J(\theta_0, \theta_1, \dots, \theta_i - \epsilon, \dots)}$$

Gradient Checking implementation notes:

- ① Don't use in training - only for debug $d_{approx}[i] \leftarrow d[i]$
- ② If algo fails grad check, look at the components to identify bug
- ③ Remember regularisation $\lambda_{reg} \sum_i \|w^{[i]}x^i\|_F^2$
- ④ Doesn't work with dropout. \rightarrow keep prob = 1.0,
- ⑤ Run at random initialisation; perhaps again after some training.

~~Notes~~

Mini-batch gradient descent

Mini-Batches

$$x^{(1)}, x^{(2)}, \dots, x^{(k)}$$

of training set: $x^{\{1\}}, x^{\{2\}}, \dots, y^{\{1\}}, y^{\{2\}}$
 (when very large training set, make mini batches). $x^{\{t\}}, y^{\{t\}}$

before $R=1000$, then use batch size 1000 and $t \in 1 \rightarrow 5000$

for $t = 1, \dots, 5000$ {

Forward prop. on $x^{\{t\}}$.

$$\begin{aligned} z^{[1]} &= w^{[1]} x^{\{t\}} + b^{[1]} \\ A^{[1]} &= g^{[1]}(z^{[1]}) \\ A^{[L]} &= g^{[L]}(z^{[L]}) \end{aligned} \quad \left. \begin{array}{l} \text{vectorial.} \\ \text{implementation.} \end{array} \right\}$$

This is doing
 1-epoch \rightarrow single pass
 through the training set.

Compute cost $J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^L l(g^{[i]}, y^{[i]}) + \frac{\lambda}{2 \cdot 1000} \sum_i \|w^{[i]}\|_F^2$

Back propagation to compute gradients $w^{[L]} := w^{[L]} - \alpha d w^{[L]}$, $b^{[L]} := b^{[L]} - \alpha d b^{[L]}$ for 2nd check

Choosing your mini-batch size \Rightarrow

Stochastic gradient descent \rightarrow size = 1.
 might never converge, etc.

Batch gradient \rightarrow size = n
 descent (too long per iteration)

Small train set ($n \leq 2000$): use batch gradient descent.

Typical mini-batch sizes \Rightarrow

$$64, 128, 256, 512, 1024 \\ 2^6, 2^7, 2^8, 2^9, 2^{10}$$

Make sure mini-batches fit in CPU/GPU memory.

Exponentially moving weighted averages.

Initialize $v_0 = 0$

$$v_t = \beta v_{t-1} + (1-\beta)d_t$$

If we know d_t for individual days.

$\beta \approx 0.9$ - average over 10 day

$\Rightarrow \beta \approx 0.98$ - average over 50 day

$\beta \approx 0.5$ - average over 2 day

Initially, because $v_0 = 0$, subsequent causes problems.

Initial phase of learning with bias correction from page.

* Gradient descent with momentum

Momentum:

$$v_{dw} = v_{db} = 0$$

On iteration t :

Compute dw, db on current mini-batch

$$v_{dw} = \beta v_{dw} + (1-\beta)dw$$

$$w = w - \alpha v_{dw}$$

$$v_{db} = \beta v_{db} + (1-\beta)db$$

$$b = b - \alpha v_{db}$$

Consider it like inertia velocity acceleration. Hyperparameters $\Rightarrow \beta, \alpha$.

(*) If have different learning rates changing both w, b)
for horizontal & vertical directions

$$\left. \begin{array}{l} \text{corrected} \\ v_{dw} = \frac{v_{dw}}{1-\beta^t} \end{array} \right\}$$

$$\left. \begin{array}{l} \text{corrected} \\ v_{db} = \frac{v_{db}}{1-\beta^t} \end{array} \right\}$$

RMS prop (Root mean square prop)

On iteration t :

Compute dw, db on current mini-batch.

$$S_{dw} = \beta S_{dw} + (1-\beta)dw^2 \leftarrow \text{element wise.}$$

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

$$w: w - \frac{dw}{\sqrt{S_{dw} + \epsilon}} \quad b: b - \frac{db}{\sqrt{S_{db} + \epsilon}}$$

ϵ can also be 10^{-8} . Just to make sure 0 not to jst.

$$\text{corrected } S_{dw} = \frac{S_{dw}}{1-\beta^t}$$

$$\text{corrected } S_{db} = \frac{S_{db}}{1-\beta^t}$$

Adam optimization algorithm

$$v_{dw} = 0; s_{dw} = 0; v_{db} = 0; s_{db} = 0$$

On iteration t :

Compute dw, db using current mini-batch

$$v_{dw} = \beta_1 v_{dw} + (1 - \beta_1) dw; v_{db} = \beta_1 v_{db} + (1 - \beta_1) db \quad \leftarrow \text{momentum } \beta_1$$

$$s_{dw} = \beta_2 s_{dw} + (1 - \beta_2) dw^2; s_{db} = \beta_2 s_{db} + (1 - \beta_2) db \quad \leftarrow \text{RMS prop } \beta_2$$

$$v_{dw}^{\text{corrected}} = v_{dw} / (1 - \beta_1^t); v_{db}^{\text{corrected}} = v_{db} / (1 - \beta_1^t)$$

$$s_{dw}^{\text{corrected}} = s_{dw} / (1 - \beta_2^t); s_{db}^{\text{corrected}} = s_{db} / (1 - \beta_2^t)$$

$$w := w - \alpha \frac{v_{dw}^{\text{corrected}}}{\sqrt{s_{dw}^{\text{corrected}}} + \epsilon}; b := b - \alpha \frac{v_{db}^{\text{corrected}}}{\sqrt{s_{db}^{\text{corrected}}} + \epsilon}$$

Adam: Adaptive moment estimation

Learning rate decay: Slowly decaying α to get better result.

So can put $\alpha = \frac{1}{1 + (\text{decay-rate}) * \text{epoch-num}}$ do formulae can vary

$$\alpha = 0.95^{\text{epoch-num}}$$

$$\alpha = \frac{k}{\sqrt{\text{epoch-num}}} \leftarrow \begin{matrix} \uparrow \text{initial} \\ \text{decay rate} \end{matrix}$$

OR discrete also possible

OR manual decay also possible.

plates make learning slow.

Training approaches:

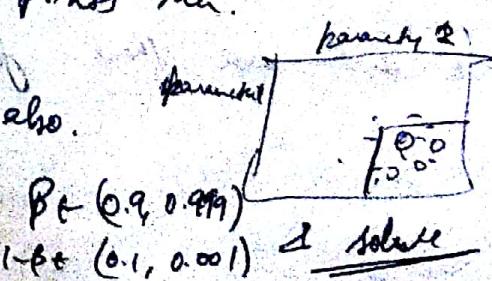
- ① Baby sitting one model (rare)
- ② Training many models in parallel (common)

Hyperparameters: $\alpha, \beta, (\beta_1, \beta_2, \epsilon), \# \text{layers}, \# \text{hidden units}, \text{learning rate}, \text{mini-batch size}, \text{decay}$

Try few random values for better, then choose a region where better hyperparameters, reduce the size of that area, the more points there.

Scaling for hyperparameters: can use

log scale



$$Y = \log_{10}(\alpha) \quad (\text{eg. for } \alpha)$$

$$\beta_1 \in (0.9, 0.999)$$

$$1 - \beta_1 \in (0.1, 0.001) \quad \downarrow \text{scale}$$

Batch Normalisation : It normalises $z^{[1]}, z^{[2]}$ also for first or $a^{[1]}, a^{[2]}$.

$$\mu = \frac{1}{m} \sum_i z^{(i)}$$

BN
considered to
be part of
a ~~next~~ layer.

$$\sigma^2 = \frac{1}{m} \sum_i (z_i - \mu)^2$$

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

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

learnable parameters
of model

Now we use $\tilde{z}^{(i)}$ instead of $z^{(i)}$

$$a^{(i)} = g^{[i]}(\tilde{z}^{(i)})$$

thus normalising the hidden layers, not only input layer.

Now new param : $w^{[x]}, b^{[x]}, \gamma^{[x]}, \beta^{[x]}$

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

$$\tilde{z}^{[x]} = w^{[x]} a^{[x-1]} + \beta^{[x]}$$

~~(can be eliminated, can be considered in $\beta^{[x]}$)~~

thus also works with various other optimisation algo.

- * Batch norm reduces the problem of input-value changing. It keeps mean and variance same.
- * Batch norm regularizes. (Adds noise to each hidden layer's activation). (slight regularization effect).
 Using bigger mini-batch size reduces regularization effect.

* μ, σ^2 are running average like moving average.

Juts from mini-batches change take, μ, σ^2 will be different for different mini-batches. So while computing, take average of μ, σ^2 and use that.

$$\text{for } z_{\text{norm}} = \frac{z - \mu}{\sqrt{\sigma^2 + \epsilon}} \text{ & } \tilde{z} = \gamma z_{\text{norm}} + \beta$$

Softmax Regression

Recognizing different objects (3 objects + others)

Number $C = \# \text{ classes} = 4$

so $\hat{y} \rightarrow (4, 1)$

element-wise.

Activation function for softmax layer:

$$a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{j=1}^C t_j}$$

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

$$t = e^{(z^{[L]})} \quad a_i^{[L]} = \frac{t_i}{\sum_{j=1}^C t_j}$$

find probability

→ Boundary to separate different classes is linear.

→ hard-max → highest value $\rightarrow 1$, others $\rightarrow 0$.

Softmax regression generalizes logistic regression to C classes.

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

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

to \hat{y}_2 make it larger for less L (loss)
prob.

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

Now $y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]^T$

Back prop: $\frac{\partial z^{[L]}}{\partial w} = \hat{y} - y$ & so on.

Tensor flow

tf. variable (,)

tf. placeholder (,)

tf. train. Gradient Descent Optimizer (), minimize (cost)

* cost ref² brain
back hab field.
with (w, b)

init = tf. global_variables_initializer()

session = tf. Session()

session.run(init).

Example of tensorflow code

$$\text{Cost } f \Rightarrow w - 20w + 25$$

import numpy as np

import tensorflow as tf

coefficients = np.array([[-1], [-20], [25]])

w = tf.Variable([0], dtype=tf.float32)

x = tf.placeholder(tf.float32, [3, 1])

cost = x[0][0]*w**2 + x[1][0]*w + x[2][0] $\# w^2 - 20w + 25$

train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)

init = tf.global_variables_initializer()

session = tf.Session()

session.run(init)

print(session.run(w))

for i in range(1000):

session.run(train, feed_dict={x: coefficients})

print(session.run(w))

} with tf.Session() as session:
session.run(init)
print(session.run(w))

It manages back-propagation
as well.

Convolution Neural Network

for vertical edges,

Vertical edge detection \Rightarrow convolution, filter, buffer $\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$

python: conv-forward

softmax: ff.nn.conv2d

multiply type adds all the elements

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix} \rightarrow \text{for horizontal edges.}$$

$6 \times 6 \times 3 \times 3 \rightarrow 4 \times 4$
(Grayscale)
(image)

Sobel filter:

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

Reid filter:

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

Image padding: 6×6 image is converted to 8×8 image using padding.

It is used because the ~~size~~ features of corner elements are usually lost, so to prevent it.

$p = \text{padding}$ $6 \times 6 \rightarrow 8 \times 8 \rightarrow \text{padding} + 1$.

now 4×8 .

Valid convolution: no padding

Same convolution: pad until the output image is of same size as the input image.

$p = \frac{f-1}{2}$ if f is usually odd.

Strided convolution:

Stride = 2: instead of convolution just jumping by 1 step. it jumps by 2 steps.

No stride = 2

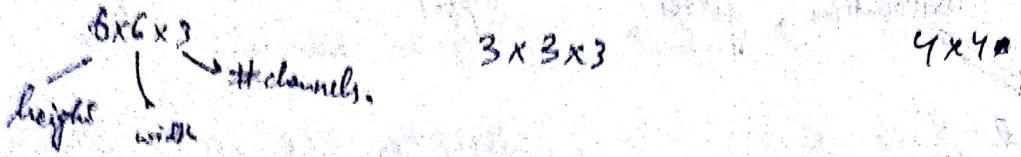
$$7 \times 7 * 3 \times 3 \rightarrow 3 \times 3$$



$n \times n * f \times f \rightarrow \left[\frac{n+2p-f}{s} + 1 \right] \times \left[\frac{n+2p-f}{s} + 1 \right]$

$\frac{n+2p-f}{s}$ is not an integer then

Convolutions on RGB images



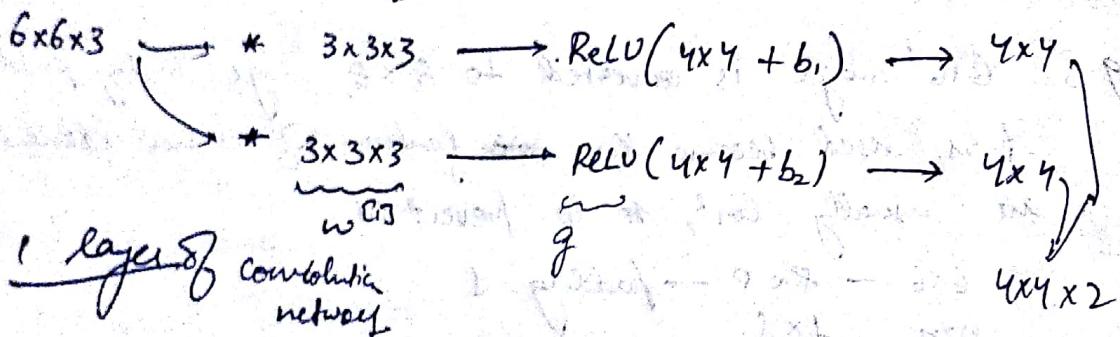
To detect vertical edges only in red

$$\text{channel} \Rightarrow R \begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix} \quad G \begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix} \quad B \begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix}$$

* as we use two different filters, we can stack them.

up after doing both individually, final will be $4 \times 4 \times 2$,

2 filters.



summary

If layer l is a convolution layer:

f^{L3} = filter size

Input: $n_H^{[L-1]} \times n_W^{[L-1]} \times n_C^{[L-1]}$

p^{L3} = padding

Output: $n_H^{[L3]} \times n_W^{[L3]} \times n_C^{[L3]}$

s^{L3} = stride

$$n_W^{[L3]} = \left\lceil \frac{n_W^{[L-1]} + 2p^{[L3]} - f^{[L3]}}{s^{[L3]}} + 1 \right\rceil$$

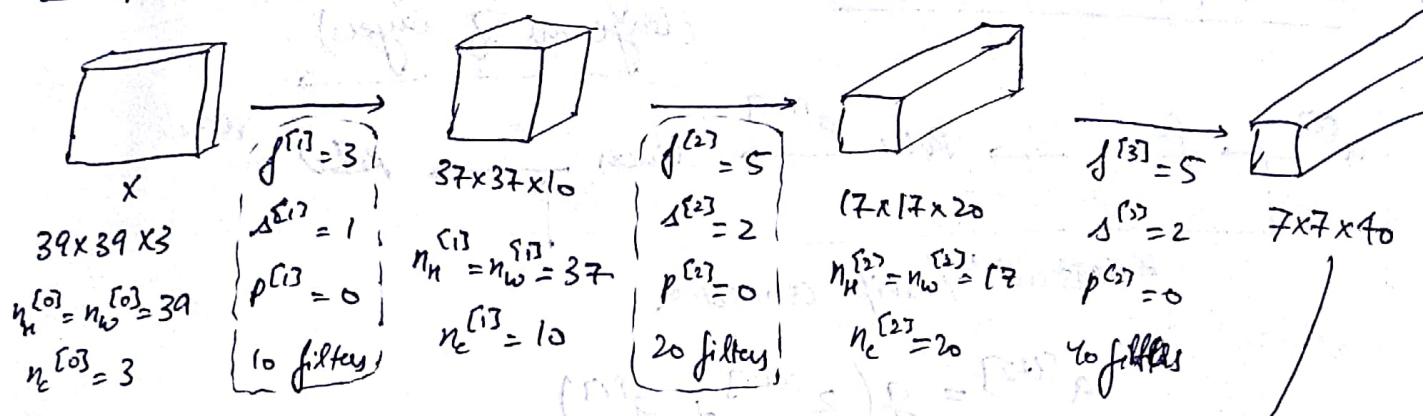
Each filter is: $f^{[L3]} \times f^{[L3]} \times n_C^{[L-1]}$

Activations: $a^{[L3]} \rightarrow n_H^{[L3]} \times n_W^{[L3]} \times n_C^{[L3]}$ $A^{[L3]} \rightarrow m \times n_H^{[L3]} \times n_W^{[L3]} \times n_C^{[L3]}$

Weights: $f^{[L3]} \times f^{[L3]} \times n_C^{[L-1]} \times n_C^{[L3]}$

bias: $n_C^{[L3]} \rightarrow (1, 1, 1, n_C^{[L3]})$ $\# \text{filters in layer } l$

Example of a convolution Neural Network



$$\frac{n+2p-f}{s} + 1$$

1960

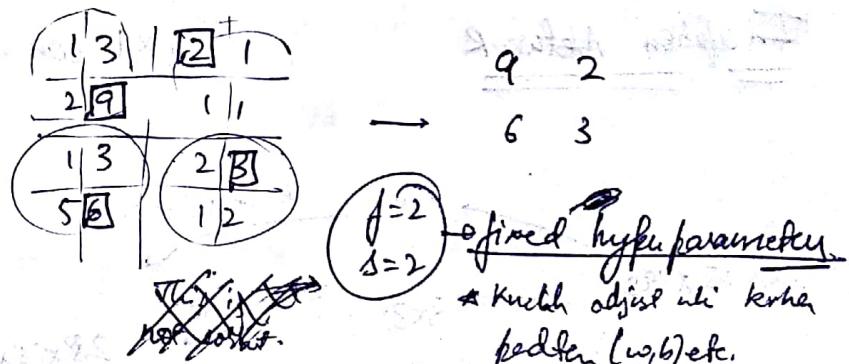
Type of layers in a convolution network:

- Convolution (CONV) → less parameters
- Pooling (POOL) → no parameters
- Fully connected (FC) → lots of parameters.

it is finally
vector longer, is to
be softmaxed.

Pooling layers → Max Pooling

f, s can change also
acc. to ex. This is wrong



Average Pooling: instead of max, we take average.

Pooling? No parameters to learn? Fixed function → parameters once set is set.

individually to weight which is isely.

* Conv + Pool is usually nila ~~defe~~ hai and make them as a single layer.

Suppose finally 400 units aa hii hai

so, then $400 \rightarrow 120 \rightarrow 84$.

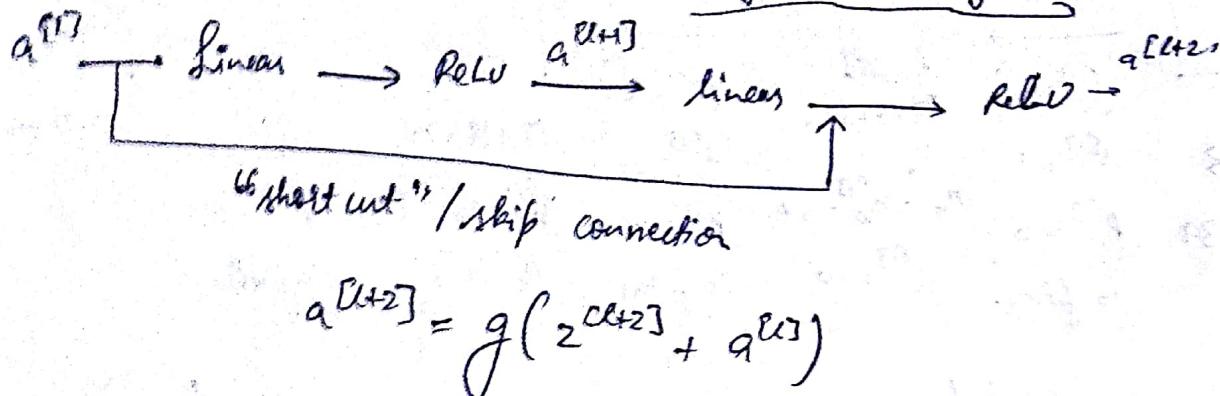
They are connected
(fully connected
layers).

to feed into
softmax.

⇒ Fully connected layer (Saara public
krite the
 w, b wala)

Ques

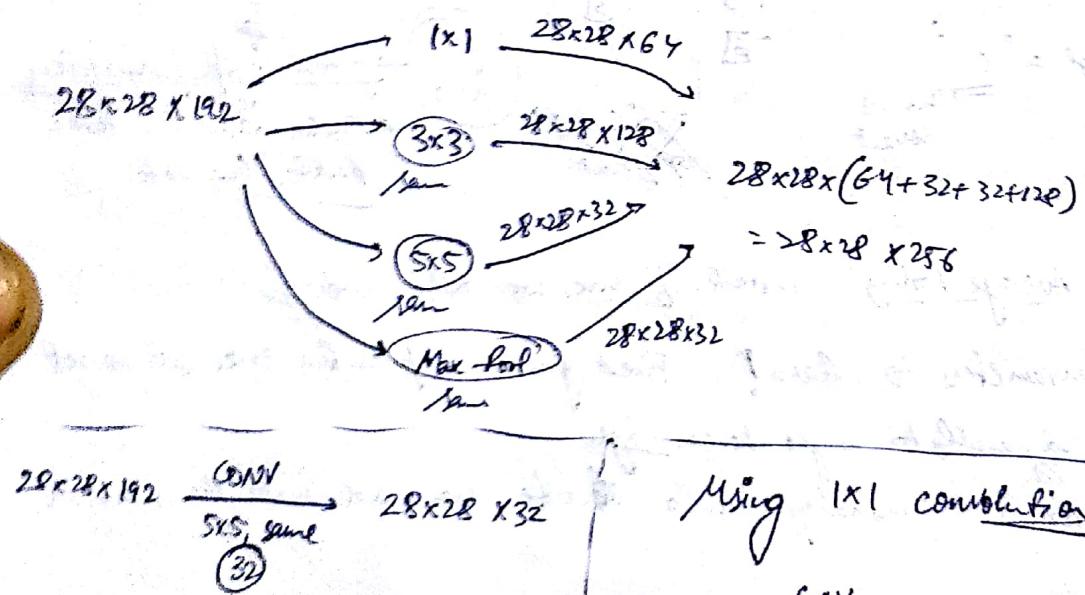
ResNet (Residual Network) : To train very-deep layered networks.
 (large no. of layers)



So we add short-cut connection to plain-network
 we continuously get decreased training error with
 increase in no. of layers.

Inception Network

\square 1×1 convolution.



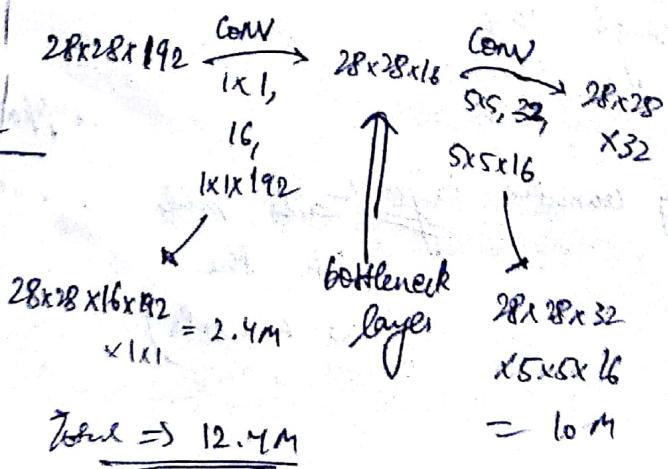
$$28 \times 28 \times 192 \xrightarrow[\text{CONV}]{5 \times 5, \text{ same}} 28 \times 28 \times 32$$

$$\text{Ap. Comput.} \approx 28 \times 28 \times 32 \times 5 \times 5 \times 192$$

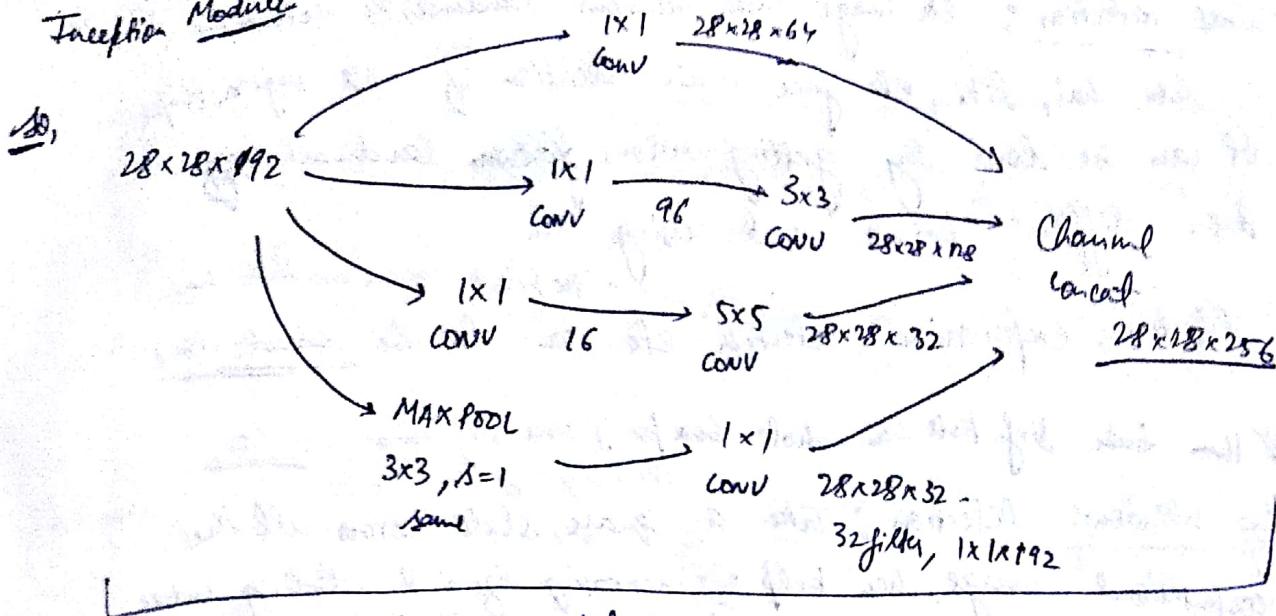
$$= 120M$$

significantly reduced

Using 1×1 convolution



Inception Module



Inception Block

Inception Network:
Inception module repeated multiple times.

* When various Inception blocks

are joined, from the ~~concat~~ channel concat, output is taken for softmax layer and check if it is not overfitting

* A regularising effect

- * When making a project, open source se download kroo and pretrained kroo hui, softmax layer change kroo. (freeze layers)
- * Name the parameters to disk. few or all depending on problem statement.

$$y = \begin{cases} p_c & \text{Is there any object?} \\ b_1 & \text{Image loc} \\ b_2 & \text{Lents co-ordinates.} \\ b_3 & \text{height, width of image.} \\ b_{4w} & \text{Probability} \\ c_1 & 1, 2, 3 \\ c_2 & 1, 2, 3 \end{cases}$$

$1, 2, 3 \rightarrow 1$ $1, 2, 3 \rightarrow \text{objects (car, pedestrian, motorcycle, etc.)}$
 $4 \rightarrow 0$ $4 \rightarrow \text{only background}$

If $p_c = 0$, then other elements are don't care.

$$l(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_8 - y_8)^2 & \text{if } y_1 = 1 \\ (\hat{y}_1 - y_1)^2 & \text{if } y_1 = 0 \end{cases}$$

* Object be around CR box.

Landmark detection: Ek image mein various landmarks detect kr leta hai, like ek face mein location of nose, eyes, shape. It can be done by getting various ~~feature~~ landmarks e.g. i.e. different points and using them.

Face detection, expression detection, etc can also be done.

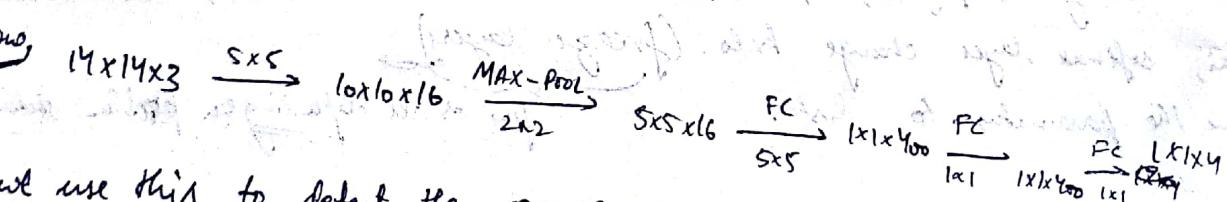
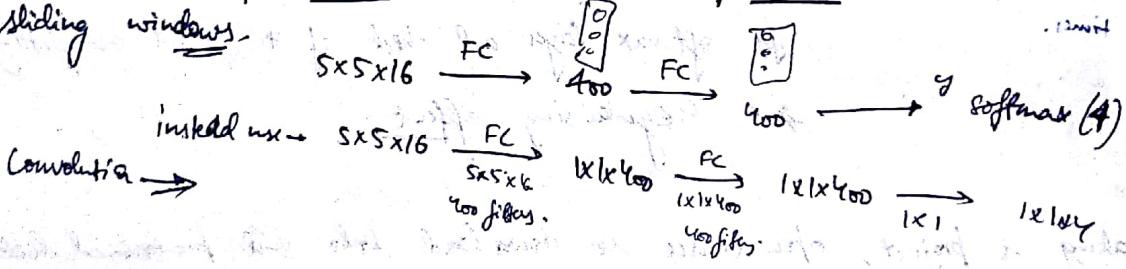
Hum fram sirf koi hai chote box pr jisme sirf ~~one~~ car ho.

Sliding Window Detection: Take a square, slide across all the

High computational usage, then keep an increasing size and sliding until cost the object is detected

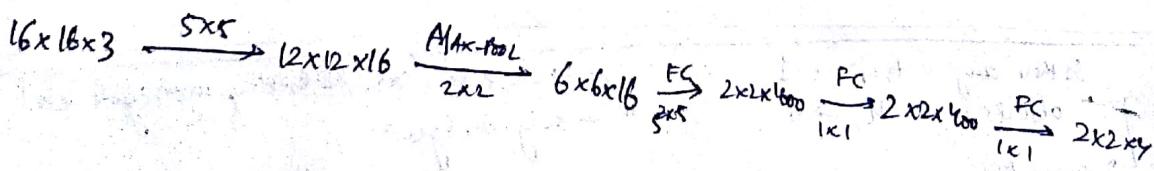
Convolution implementation: To reduce computation

① Sliding windows.



so we use this to detect the ~~output~~ of sliding windows.

Suppose



We used a 14x14x3 box to run over this image and window was slides by 2 pixels each time. So what we would get is ~~one~~ + values for 4 position.

Those four values are
 top-left
 top-right
 bottom-left
 bottom-right

Bounding Box Prediction

YOLO : 3x3 ya 19×19 type min image to divide into then
 the object belongs to that ~~1 box~~ where the center of the image lies.
 very fast to real time the prediction is:
 now for $y = b_x, b_y$ ~~size~~ B should signify the total size of the ~~case~~.

Intersection over Union (IOU) = $\frac{\text{Size of intersection}}{\text{Size of union}}$ \rightarrow if predicted ad what actually is boxes,
 we consider prediction is correct if $\text{IOU} \geq 0.5$.

Non-max suppression : Possibility when image is divided into larger no. of boxes, more boxes may see that center is there.
 Thus multiple detection.

→ See the one with highest prob. of matching and the boxes with high IOU with this box get suppressed. Next we find next higher prob. in the remaining & follow.

$p_c \rightarrow$ prob. Before the above step, we also discard the boxes with $p_c \leq 0.1$.

Anchor boxes: To detect multiple type of objects/ objects located in same frame

→ Predefine two shapes / two anchor boxes (horizontal, vertical)

and put 16 variables. (8 for anchor box 1, 8 for anchor box 2).

What we do here is calculate which anchor box has higher IOU for that part and we use that anchor box for that part.

Algo not good for

- ① 2 anchor box, 3 objects
- ② 2 objects, same type of anchor box.
- ③ Have have midpoint of objects.

YOLO ALGO

- 1 - detection
- 2 - box
- 3 - output

J - radio box

image divided finally in 3x3

$y \in 3 \times 3 \times 2 \times 8$

of 3x3 Algo Sunday S + Monday
Classes

1 + New order boxes can also be used acc to our need.

Region proposal (R-CNN): Rather than running CNN on all windows you just select few windows

to represent objects and no overlaps

task for selecting the # finds

region for running a few areas at the same time

different time

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: pre convolution network to propose regions.

ROI → You Only Look Once

Face recognition
we have only 1 image, not sufficient to train a neural network.

We use similarity function $\Rightarrow d(\text{img}^1, \text{img}^2) = \text{degree of difference between images}$.

If $d(\text{img}^1, \text{img}^2) \leq \tau$ "same"
 $d(\text{img}^1, \text{img}^2) > \tau$ "different" } verification

we we have to find d .

This can be used when a new face joins ad network have to learn acc. to it.

Siamese Network: Feed two images into the system of neural network and calculate the matrix of those two images. we then calculate

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

$f(x^{(1)})$ if $x^{(1)}, x^{(2)}$ are the same person $\|f(x^{(1)}) - f(x^{(2)})\|_2^2$ is small & vice-versa

Triplet loss

Learning

Anchor A

Positive P

Anchor N

Negative N

$$\text{is } \|f(A) - f(P)\|^2 + \underbrace{d}_{\text{margin}} \leq \|f(A) - f(N)\|^2$$

parameters \rightarrow needed because the network forms such that $f(A) \neq f(N)$ always or $f(A) = 0$ always.

To prevent that d is and.

$$L(A, P, N) = \max \left(\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \epsilon, 0 \right)$$

To do problem and loss there, calculate.

Training set: 10K pictures of 1K people is to be used.

Training: if A, P, N are chosen randomly ; the eqn is easily satisfied.

Thus we have to choose A, P, N which are hard to train.

Choose such that $d(A, P) \approx d(A, N)$ and thus hard to predict

Neural style transfer: Image C (content) to be painted in style S to obtain G (generated image)

Deeper layers calculate more complex parts of the image

$$J(G) = \alpha J_{\text{Content}}(C, G) + \beta J_{\text{Style}}(S, G)$$

cost?

Process

① Initiate G randomly

② Use gradient descent to minimize $J(G)$ $G \leftarrow G - \frac{\partial J(G)}{\partial G}$

Content Cost function ($J_{\text{Content}}(C, G)$)

- A hidden layer l is used to compute content cost. It should neither be too shallow nor too deep
- Use pre-trained ConvNet
- Let $a^{[L](C)}$ and $a^{[L](G)}$ be the activation of layer l on images C and G . If $a^{[L](C)}$ and $a^{[L](G)}$ are similar, both images have similar content

$$J_{\text{Content}}(C, G) = \frac{1}{2} \|a^{[L](C)} - a^{[L](G)}\|_1^2$$

Style Cost function ($J_{\text{Style}}(S, G)$)

Style matrix

height, width

layer l

def $a^{[L](s)}_{ijk} = \text{activation at } (i, j, k)$. $G^{[L]}$ is $n_h \times n_w \times n_c$

$$G_{kk'}^{[L](s)} = \sum_{i=1}^{n_h} \sum_{j=1}^{n_w} a_{ijk}^{[L](s)} a_{ijk'}^{[L](s)} \quad \{ \text{image } S \text{ at layer } l \}$$

$k=1, 2, \dots, n_c^{[L]}$

$$G_{kk'}^{[L](G)} = \sum_{i=1}^{n_h} \sum_{j=1}^{n_w} a_{ijk}^{[L](G)} a_{ijk'}^{[L](G)} \quad \{ \text{image } G \text{ at layer } l \}$$

Overall $J_{\text{Style}}(S, G)$
will be sum, i.e.

$$\sum_L \lambda^{[L]} J_{\text{Style}}^{[L]}(S, G)$$

$$J_{\text{Style}}(S, G) = \|G^{[L](s)} - G^{[L](G)}\|_F^2$$

$$= \frac{1}{(2n_h n_w n_c)^2} \sum_L \sum_{k'} \left(G_{kk'}^{[L](s)} - G_{kk'}^{[L](G)} \right)^2$$

Sequence Models

⇒ Recurrent Neural Networks

$$\underline{X^{(i)}}$$

t^{th} training element is like multiple words of a sentence.

i^{th} training example -

T_n → length of input

T_y → length of output.

(Toh Kannaa word)
 $n^{<1>} \quad n^{<2>}$

$T_x^{(i)}$ → input length of i^{th} training example

$T_y^{(i)}$ → output length of i^{th} training example

Use a dictionary {Vocabulary} → take fixed amount of words, index those words.

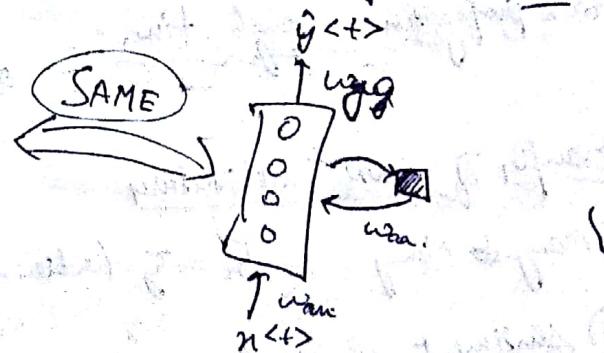
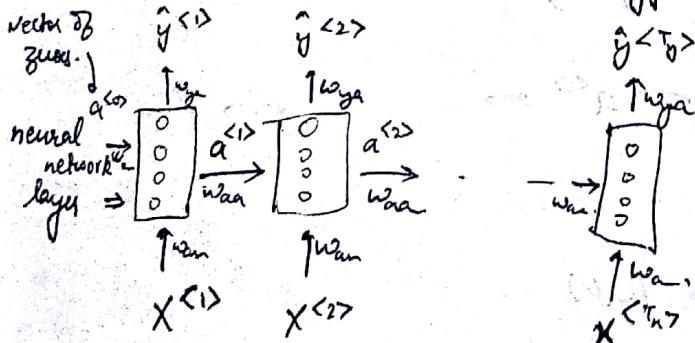
Now for each word, make a vector of $0, 1$ of the size of dictionary.

(where the word matches the dict, rest all 0. (This only 1 one in each vector))

like for 'a' $\Rightarrow \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$ when dict $\Rightarrow \begin{bmatrix} a \\ \text{car} \\ \text{sun} \end{bmatrix}$

* if the word is not in dict., make a fake token like <UNK>

⇒ Not the standard network as different layers does not share features learned



BRNN (Bidirectional RNN): To use all the words in the sentence.

$$a^{<1>} = g_1(w_{aa} a^{<0>} + w_{an} n^{<1>} + b_a) \quad \text{fwd/RNN}$$

$$\hat{y}^{<1>} = g_2(w_{ya} a^{<1>} + b_y) \quad \leftarrow \text{sigmoid.}$$

$$a^{<6>} = g_1(w_{aa} a^{<5>} + w_{an} n^{<6>} + b_a); \quad \hat{y}^{<6>} = g_2(w_{ya} a^{<6>} + b_y)$$

WAN uses n type entity to compute a type ent'l

Simplified notation

$$a^{<t>} = g(w_a[a^{<t-1>}, n^{<t>}] + b_a); \hat{y}^{<t>} = g(w_y a^{<t>} + b_y)$$

$$w_a = [w_{aa} \mid w_{an}]$$

Now if a is (l, o) dimensional.

$$\text{so } w_{aa} \rightarrow (l, l, o)$$

$$\text{Thus } w_a \rightarrow (l, l, o)$$

and if n is (k, o) dimensional

$$\text{so } w_{an} \rightarrow (l, k, o)$$

$$\text{also } E^{<t+1>}[n^{<t>}] = \begin{bmatrix} a^{<t+1>} \\ n^{<t>} \end{bmatrix} \xrightarrow{\text{size}} \begin{bmatrix} l \\ l \\ l \\ l \end{bmatrix}$$

Back propagation

$$\text{Here we will be using loss function: } \mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>}) = -y^{<t>} \log \hat{y}^{<t>} - (1-y^{<t>}) \log (1-\hat{y}^{<t>})$$

$$\text{Now total loss, } L(\hat{y}, y) = \sum_{t=1}^T \mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

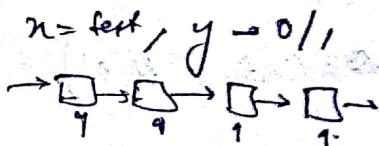
Back propagation through time:

$$a^{<1>} \rightarrow a^{<2>} \rightarrow \dots \rightarrow a^{<T>}$$

Example of RNN architecture

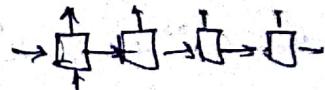
① Many-to-Many $T_n = T_y$ (when same length)

② Sentiment classification - (many-to-one) $n = \text{feat}, y = 0/1$

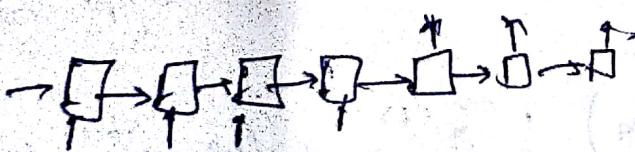


③ One-to-one \rightarrow basic word

④ One-to-many (e.g. music generation)



⑤ ~~Many~~ Many-to-many - different lengths

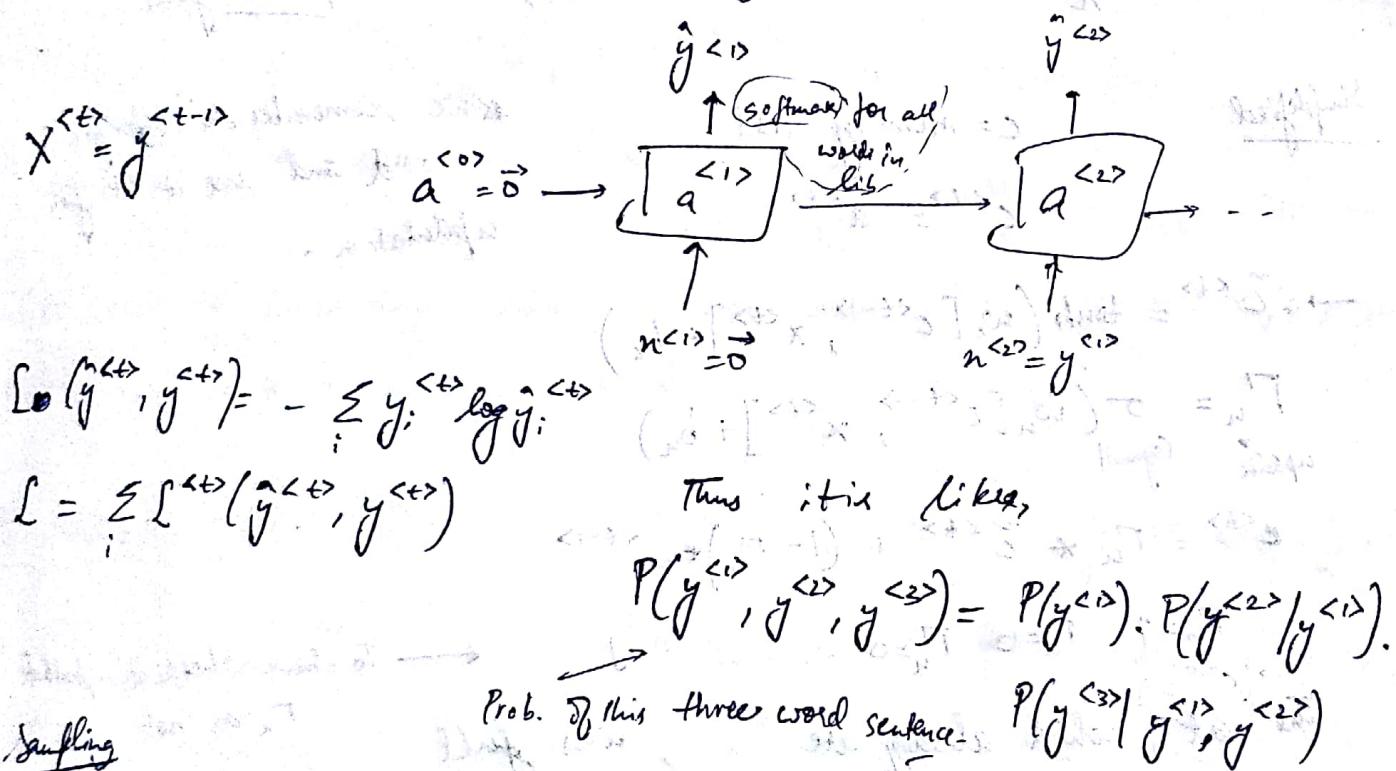


Language Modelling : determines the probability of being that sentence
of all the similar sounding sentences.

Training Set : large corpus of English text.

Approach

- * Tokenise the words, end the sentence with `<EOS>` token.
- * `<UNK>` token for words not in library.



Character-level language model : Every character is in vocab, and the sentence is divided on the basis of every character and not word.

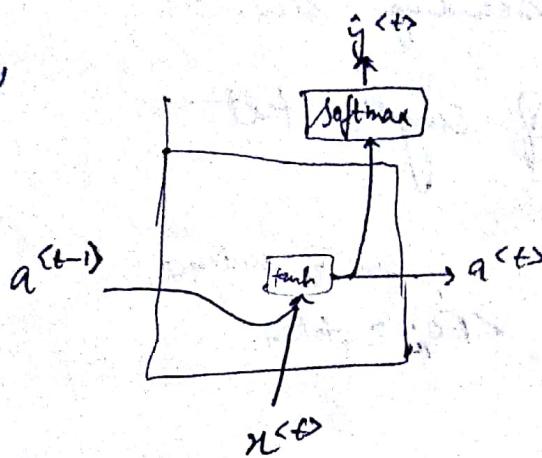
* No prob. of UNK.
* very computationally expensive.

Vanishing gradients : Kabhi kabhi, because of very initial words, the later word in the sentence is affected.

Exploding gradients : Apply gradient clipping.

Gated Recurrent Unit:

RNN



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

The cat, ---, was full
The cats, ---, were full.

Simplified

c = memory cell

$$c^{t+1} = a^{t+1}$$

we remember the specific word and use it like for update.

$$\rightarrow \tilde{c}^{t+1} = \tanh(w_c[c^{t-1}, x^t] + b_c)$$

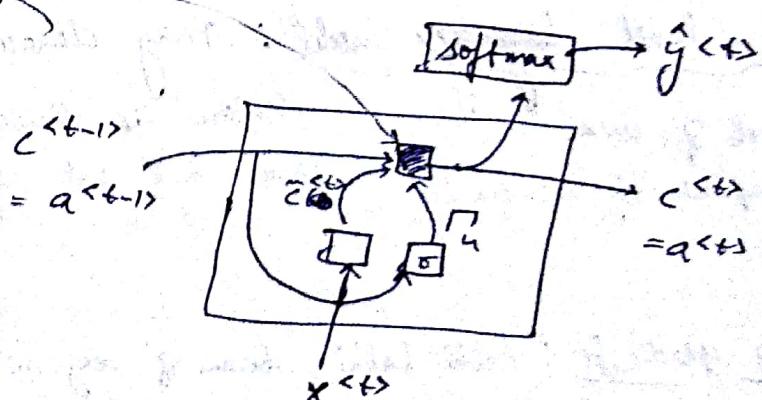
$$\Gamma_u = \sigma(w_u[c^{t-1}, x^t] + b_u)$$

update (sigmoid)

$$c^{t+1} = \Gamma_u * \tilde{c}^{t+1} + (1 - \Gamma_u) * c^{t-1}$$

$\Gamma_u = 1 \quad \Gamma_u = 0 \quad \Gamma_u = 0.5 \dots J.$ ← To check whether to update c^{t+1} or not.
The cat, which already ate ---, was full

This simplified version of GRU unit is



Now full GRU

$$\tilde{c}^{t+1} = \tanh(w_c[\Gamma_r * c^{t-1}, x^t] + b_c)$$

$$\Gamma_u = \sigma(w_u[c^{t-1}, x^t] + b_u)$$

$$\Gamma_r = \sigma(w_r[c^{t-1}, x^t] + b_r)$$

Another model
LSTM

LSTM \rightarrow Long Short Term Memory.

$$\tilde{c}^{<t>} = \tanh(w_c[a^{<t-1>}, n^{<t>}] + b_c)$$

$$a^{<t>} + c^{<t>}$$

$$\Gamma_u = \sigma(w_u[a^{<t-1>}, n^{<t>}] + b_u) \quad (\text{update})$$

$$\Gamma_f = \sigma(w_f[a^{<t-1>}, n^{<t>}] + b_f) \quad (\text{forget})$$

$$\Gamma_o = \sigma(w_o[a^{<t-1>}, n^{<t>}] + b_o) \quad (\text{output})$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>} \quad \left\{ \begin{array}{l} \\ \end{array} \right.$$

$$a^{<t>} = \Gamma_o * \tanh c^{<t>} \quad \left\{ \begin{array}{l} \\ \end{array} \right.$$

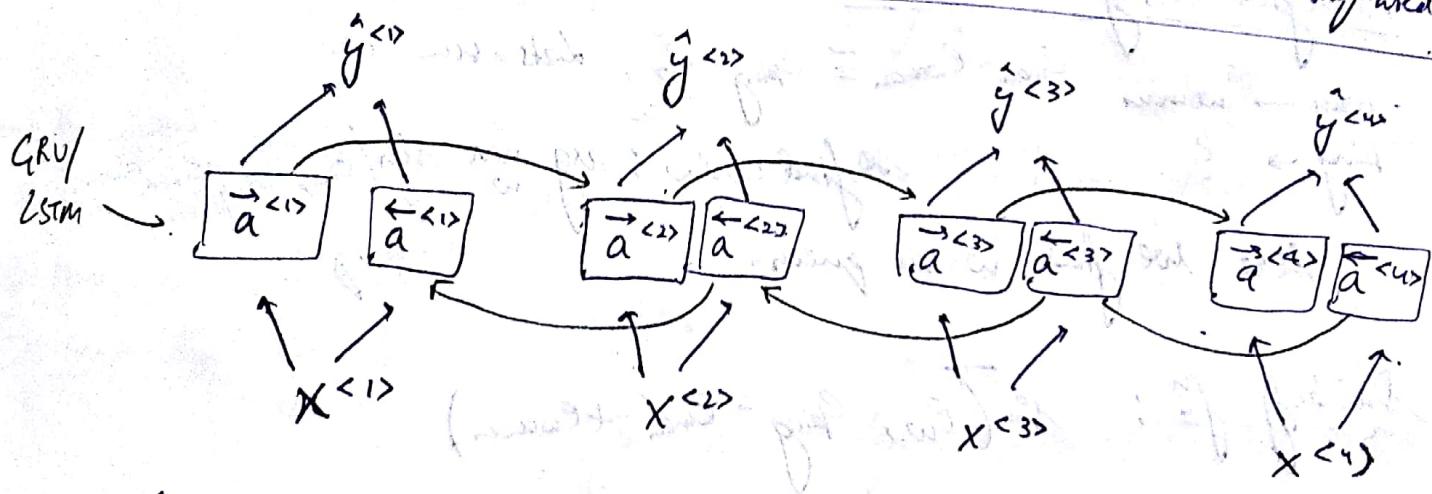
per hole connection

Sometime people put $c^{<t-1>}$
also $a^{<t-1>}$ and $n^{<t>}$,
it also effects.

GRU easier to build bigger models whereas LSTM is more flexible. People prefer to use LSTM but GRU is simpler.

Bi-directional RNN : (BRNN)

\Rightarrow BRNN with LSTM is also commonly used.



layer
 $a^{<t>} < t >$ activation
time

time = no. of words
we.

layer = no. of words
predicted by
the layer

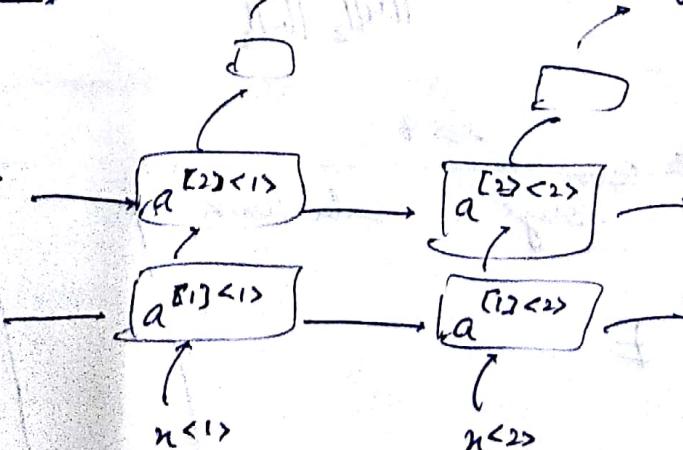
$$a^{(27)<0>}$$

$$a^{(17)<0>}$$

$$n^{<1>}$$

$$n^{<2>}$$

$y^{<1>}$ $y^{<2>}$



1-hot representation: But problem here is that words that might be linked to each other are also placed very far apart.

E. 06257

(06257)
(300, 10K)

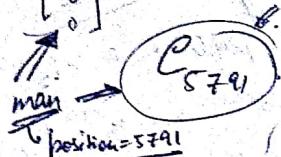
(10K, 1)

$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} = \mathbf{o}_{5791}$

Featureized representation: Word Embedding!

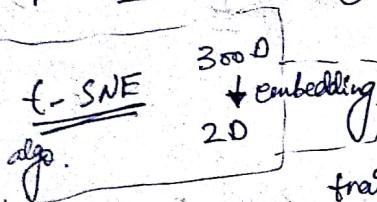
Words are given different types of features. Like

Gender, Royal, Age, Food, ... Suppose there are 300 features
so table consists of those features



* Transfer learning & word embeddings

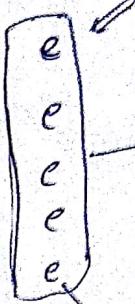
① Learn word embeddings from large text corpus (1-100B words)
(or download pre-trained embedding online)



② Transfer embedding to new task with smaller training set (say, 10K words)

③ Optional: Continue to finetune the word embeddings with new data.

this is
embed



w
Context
 $O_c \rightarrow$

Analogies using word vectors

man - woman \approx king - ?
man \rightarrow woman

$C_{\text{man}} - C_{\text{woman}} \approx C_{\text{king}} - C_w$, let's assume w

king \rightarrow ?
so, find word w: arg max_w sim($C_w, C_{\text{king}} - C_{\text{man}} + C_{\text{woman}}$)
similarity

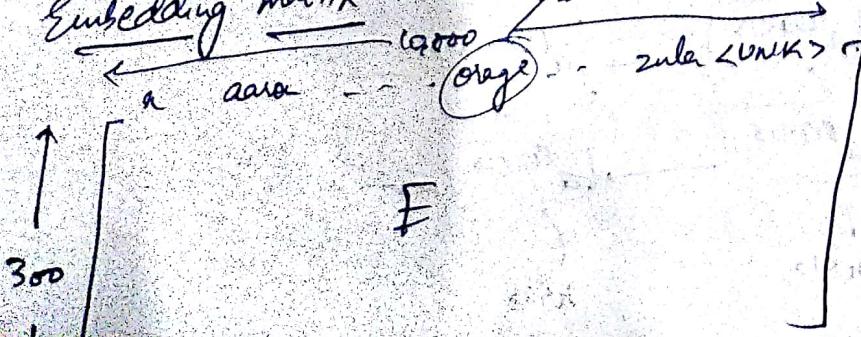
here we find w as queen.

Similarity fn: $\text{sim}(C_w, C_{\text{king}} - C_{\text{man}} + C_{\text{woman}})$

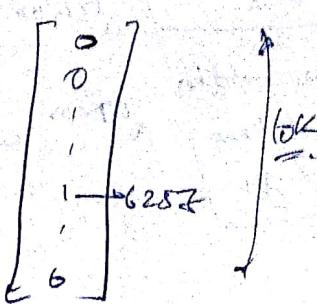
Cosine similarity: $\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2} \Rightarrow \text{cosine b/w } u, v$.

[Maximize of!]

Embedding matrix



\mathbf{o}_{G257}



Now this

P near

② this

① like

② n

Context

be

E. $O_{257} \rightarrow e_{(300,1)}$

$\begin{matrix} / & \backslash \\ (\text{orange}) & (10k_1) \\ \downarrow & \uparrow \\ (300, 10k) \end{matrix}$

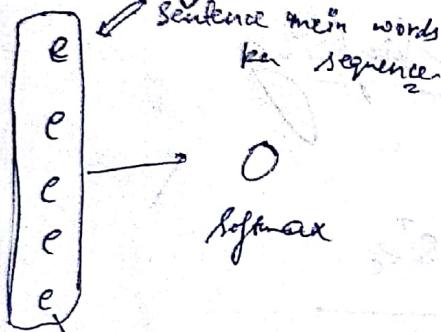
e_j
orange.

(E). (one hot vector)

E. $O_j = e_j$

= embedding for word j

This is kinda inefficient, so use specialised fn to look up an embedding.



There can be various contexts. like.

- last 4 words
- 4 words on left & right
- last 1 word
- nearby 1 word.

Q: Now we put them according to various contexts.

Q_t = parameters associated with output t .

"orange" "juice"

Context c Target t . softmax: $p(t|c) = \frac{e^{Q_t^T e_c}}{\sum_{j=1}^{10,000} e^{Q_j^T e_c}}$

$O_c \rightarrow E \rightarrow e_c \rightarrow O \rightarrow y$
 $e_c = E O_c$ softmax

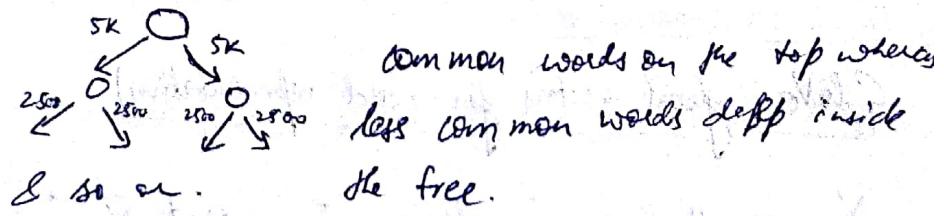
$$L(y, \hat{y}) = - \sum_{i=1}^{10,000} y_i \log \hat{y}_i$$

Now this softmax is slow bcz computing

P mean calculating all. For that purpose,

① There can be used:

① Hierarchical softmax -



② Negative Sampling -

Context ' c ' must not be random, like (jf , ker , a , $-$) when used, e_c is to be updated again and again, calculated ~~by~~ by it takes longer time.

New learning problem

$R = 5-20$ small dataset

$k = 2-5$ large dataset

context	word	target?
orange	juice	1

$\hat{P}: 1 \text{ negative to positive}$

ratio
examples.

orange	king	0	$R=4$
orange	book	0	words taken randomly
orange	of	0	
orange	the	0	
q	q	q	
c	t	y	

Negative Sampling

$$P(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

$$P(y=1|c, t) = \sigma(\theta_t^T e_c)$$

Instead of calculating for 10,000 words we randomly take words and calculate for just $(1+k)$ words. [($k+1$) binary classification problems.]

This is 10,000 binary classification problem.

Selecting negative example: $f(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10,000} f(w_j)^{3/4}}$ $f \rightarrow$ frequency of that word.

Glove (global vectors for word representation)

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

$X_{ij} = X_{ji} \rightarrow$ if in range of 10 words

Q_i, e_j are symmetric

$X_{ij} \neq X_{ji}$ if context like "before j", "after j"

minimize $\sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij}) (\theta_i^T e_j + b_i - b_j - \log X_{ij})^2$

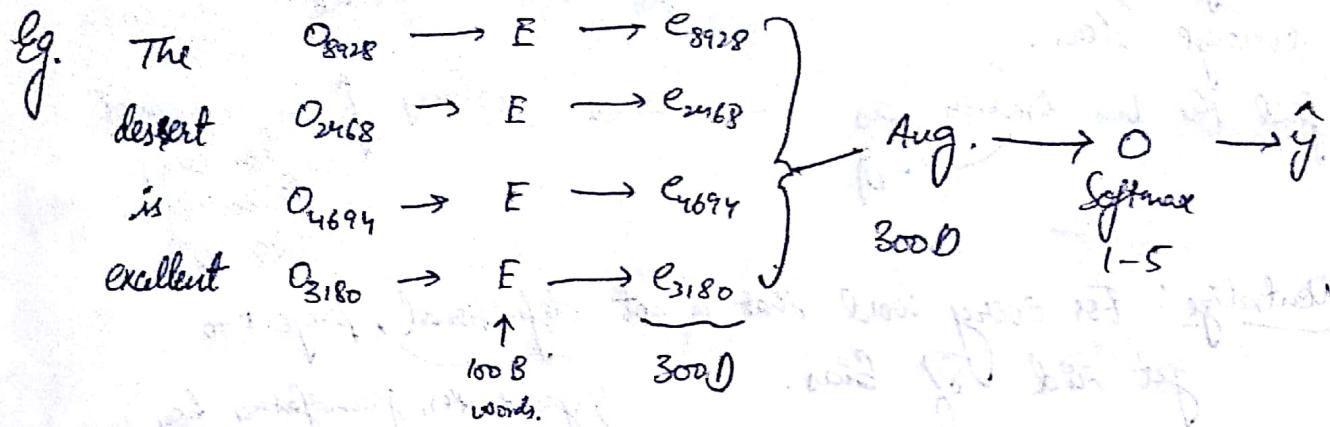
$f(X_{ij}) = 0 \text{ if } X_{ij} = 0 \quad \theta \log 0 = 0$

$e_w^{(\text{final})} = \frac{e_w + Q_w}{2}$

$$\text{minimize} \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(x_{ij}) (\mathbf{e}_j^T \mathbf{e}_i + b_i - b_j - \log x_{ij})^2$$

Sentiment Classification

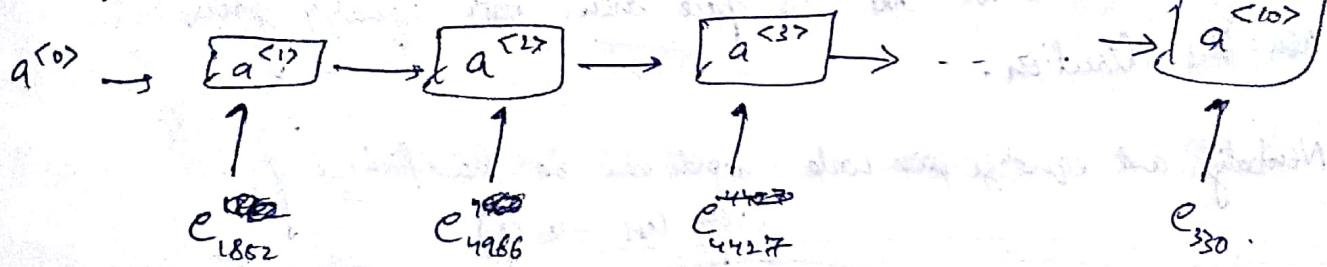
(review
in words. \rightarrow stay
convert to nra).



Now, "completely lacking in [good] food, [good] service, and [good] ambience."

\rightarrow Here just taking an average will lead in wrong results, bcz too many of "good".

RNN for sentiment classification



many-to-one RNN.

Debiasing Word Embeddings

Word embeddings can reflect gender, ethnicity, age, sexual orientation, and other biases of the text used to train the model. which is wrong like

Father: Doctor as Mother: Nurse X.

So we have to debias.

Distilling Box

babysitter

Doctor

- ① Identify tree structure

spouse {
wife
husband
son of your dad
daughter

we find for your brother as
(-D)

grandmother

girl

bro

• Grandfather

boy

he

lara
(-D)

•

boy

he

children { is unrelated
(non-bias
direction)
(299-D)

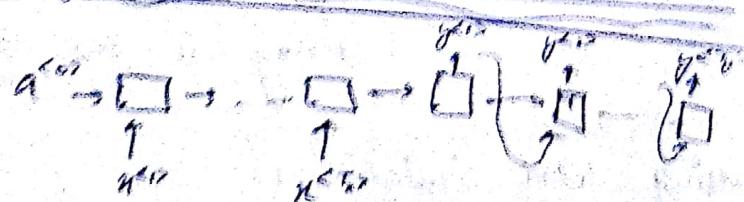
- ② Ranking: For every word dict is not definitional, project do
not add by bias.

{ grandmother, grandfather, boy,
girl, etc. }

- ③ Equity Bias: Eg. grandmother and grandfather with babysitter.
although grandmother has a higher chance of being a babysitter,
dictator grandfather will be given a priority for babysitter which is
wrong. So what we do is give them both equally from
non-bias direction.

[Ranking and equity bias words can be handpicked.]
(very few words)

Sequence to Sequence Model

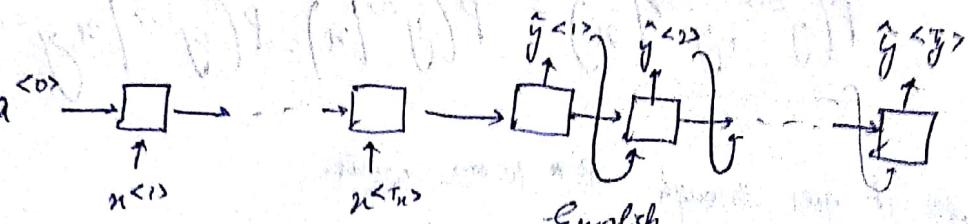


- ④ Translation (French to English)

- ⑤ Image Captioning

Machine translation as building a conditional language model

Machine translation



"conditional language model"

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

Finding the most likely translation

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

+ Randomly choose first words.

Greedy Search

① best $\hat{y}^{<1>}$ word choose $y^{<1>}$, then best $\hat{y}^{<2>}$ word choose $y^{<2>}$, and so on. But not a good approach.

$\rightarrow P(\hat{y}^{<1>} | x)$ is similar to obtain $P(\hat{y}^{<1>} | x), P(\hat{y}^{<2>} | x), P(\hat{y}^{<3>} | x), \dots, P(\hat{y}^{<Ty>} | x)$

but no sake has that we choose initial words that are better fitting but sentence that forms is not good.

Approximate Search

$$\arg \max_y P(\hat{y}^{<1>} | x), \dots, \hat{y}^{<Ty>} | x)$$

② We can't choose the words randomly, so we use various types of search algorithms for that purpose.

Beam Search

B = beam width \Rightarrow no. of possibilities considered by beam search at a time.

Step 1 $P(y^{<1>} | x)$ is used to search for ③ most likely words that could fit.

for first word and keep track of them $a^{<0>} \rightarrow \hat{y}^{<1>} \rightarrow \hat{y}^{<2>} \rightarrow \dots$

Step 2 $a^{<0>} \rightarrow \hat{y}^{<1>} \rightarrow \hat{y}^{<2>} \rightarrow \hat{y}^{<3>} \rightarrow \dots$ calculate $P(y^{<1>} | x), P(y^{<2>} | x), P(y^{<3>} | x), \dots$. Then calculate $P(y^{<1>} | x), P(y^{<2>} | x), P(y^{<3>} | x), \dots$

Calculate all B possibilities.

$\xrightarrow{1 \rightarrow 3 \rightarrow 12}$
huge gain

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

so it goes through $B * 10,000$ words
After this, we select B such combinations of 1st and 2nd words

In the BFS faster

Also, only B copies of network are to be used.

Step 3: 3rd word $P(y^{<3>} | x, y^{<1>}, y^{<2>})$

and it keeps going.

To be terminated by <EOS>

finally take the highest $P(y^{<1>} y^{<2>} \dots y^{<T>} | x)$

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

usually this product is very tiny.

so instead this is used $\sum_{y=1}^T \log(P(y^{<1>} | x, y^{<2>}, y^{<3>} \dots, y^{<T-1>}))$

thus $\arg \max_y \sum_{y=1}^T \log(P(y^{<1>} | x, y^{<2>}, y^{<3>} \dots, y^{<T-1>}))$

Then, we normalize to some extent,

we divide by T^α $\alpha \in (0, 1)$
 not normalizing $\alpha \rightarrow$ hyperparameters
 normalizing fully

finally, $\arg \max_y \frac{1}{T^\alpha} \sum_{y=1}^T \log(P(y^{<1>} | x, y^{<2>}, \dots, y^{<T>}))$

Beam width selection:
 Large B: better result, slower
 Small B: worse result, faster.

$100 \rightarrow 5000$
 very less gain

Like BFS (Breadth First search) or DFS (Depth first search), Beam search runs fast but is not guaranteed to find exact max. for $\arg \max_y P(y|x)$.

Error analysis on beam search.

(increasing) B never hurts.

y^* translation by human.
 \Rightarrow very good. $P(y^*|x) > P(\hat{y}|x)$

\hat{y} by algo \Rightarrow not that good
 \rightarrow to be compared by RNN.

Case 1: RNN predicts $P(y^*|x) > P(\hat{y}|x)$

Beam search chooses \hat{y} . But y^* attains higher $P(y|x)$.

\rightarrow Thus, Beam search is at fault.

Case 2: $P(y^*|x) < P(\hat{y}|x)$

y^* is a better translation than \hat{y} . But RNN predicted $P(y^*|x) < P(\hat{y}|x)$.

\rightarrow Thus, RNN model is at fault.

For error analysis, for example see go through his and check which is majorly at fault, Beam or RNN and then optimize that.

Blue Score (bilingual evaluation under way)

Other give references as a part of test set and Machine Text Output is seen.

Modified precision : $\frac{\text{max count of word in } \hat{x} \text{ references}}{\text{Count of word in output}}$

(for pair of words: bigram)

(Count clip) (Count for each pair, calculate Count and Count clip)

and sum those, then calculate modified precision

$$P_i = \frac{\sum_{\text{unigram } \in j} \text{Count}_{\text{clip}}(\text{unigram})}{\sum_{\text{unigram } \in j} \text{Count}(\text{unigram})}$$

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

$$P_1 = P_2 = \dots = P_n = 1.0$$

So equal to one if the reference is machine output

~~Combined bleu score %~~: $BP = \exp\left(\frac{1}{n} \sum_{k=1}^n p_k\right)$

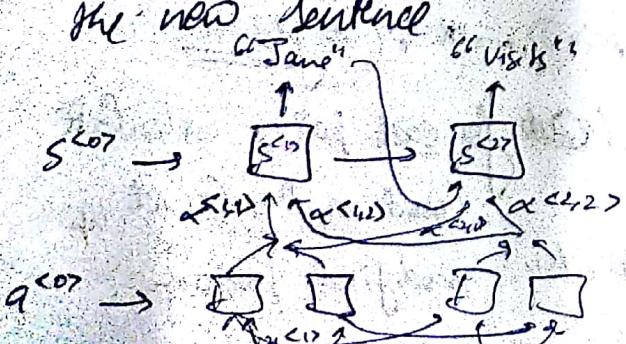
BP = Brewey penalty

$$BP = \begin{cases} 1 & \text{if MT_output_length} \leq \text{reference_output_length} \\ \exp(1 - \text{MT_output_length}/\text{reference_output_length}) & \text{otherwise} \end{cases}$$

Attention Model intuition

Problem of long sequences: Model have to memoise a long sentence, so its bleu score drops as it is not that efficient.

Every word is given a certain amount of attention for forming the new sentence



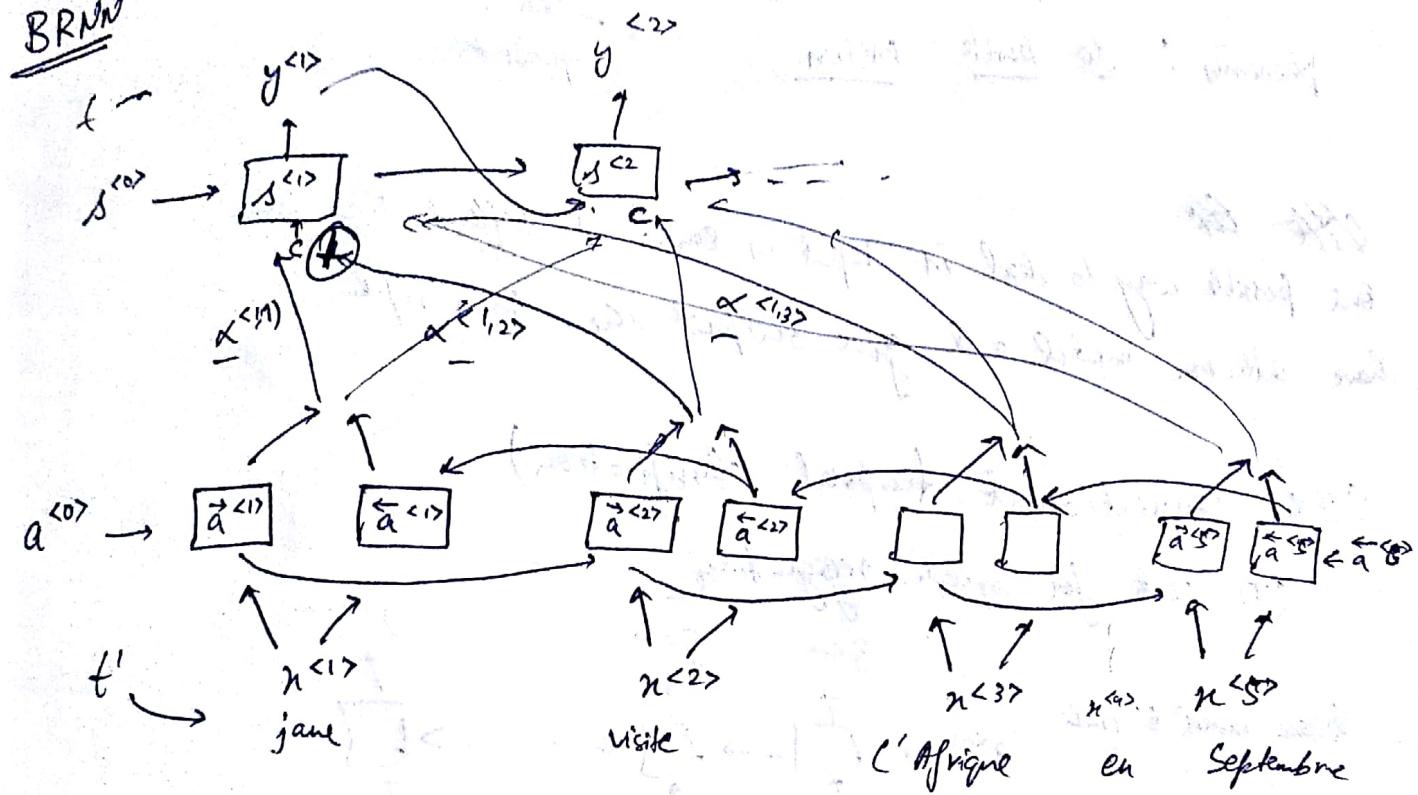
$$\alpha^{1,1}, \alpha^{1,2}, \alpha^{2,1}, \alpha^{2,2}$$

Are the ~~weights~~ different given to different words

$\alpha^{(t,t')}$ are the attention weights

f word bname be lie t' ki taf se kites attention.

BRNN



$$a^{(t,t')} = (\vec{a}^{(t)}, \underline{a}^{(t)})$$

weighted sum of attention weight $\Rightarrow \sum_{t'} \alpha^{(t,t')} = 1$

$$c^{(1)} = \sum_{t'} \alpha^{(1,t')} a^{(t')}$$

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

$$c^{(2)} = \sum_{t'} \alpha^{(2,t')} a^{(t')}$$

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

$\Rightarrow T_n$ is the amount of attention, so high cost.

$$s^{(t-1)} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow e^{(t,t')} \quad \& \alpha^{(t,t')}$$

small NN to obtain $e^{(t,t')}$

phonemes: smallest units of cell.

like "the quick brown fox"

phonemes: de quick brown

→ hard-to-write
representation

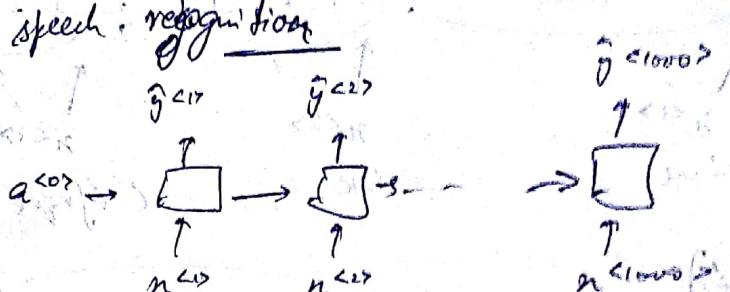
~~CTC~~ ~~CTC~~

One possible way to deal is input is divided at different time frames.
have attention model and give output the transcript.

CTC (Connectionist temporal classification)

CTC cost for speech recognition

Basic model is trah



possibility that there is a cont along t.o.
(Output is ~~way~~ too less)

so output is ttt - h - eee - - l - - gggg -

Basic Rule: collapse repeated characters not separated by "blank".

⇒ the → t

the quick brown fox → 19 characters

so 19 to expand like 1000.

Trigger Word Detection

Turn on something by saying a word like "Alexa play despacito" ↗

RNN output to be set to 1 just after the trigger word was said, else

0.

