

- Also in backpropgⁿ, we need each value in chain rule, so we need some memory to store our temporary results.

→ Challenges in DL

- Deep learning is data hungry
↳ can be solved by augmentⁿ
- Overfitting / Lack of generalizⁿ
↳ solved by using bottleneck layer
↳ dropout
↳ generalizⁿ loss (regularizⁿ)
- Vanishing / exploding gradient problem
↳ severe for deep network
↳ depends on activⁿ funcⁿ (vanishing problem) & the weight initializⁿ (exploding problem)
• So, weights are chosen randomly.

$\frac{w_i}{\sqrt{n_{nodes}}} \rightarrow 0$... we try to initialize weights with mean=0 & variance = $\frac{1}{N}$

$$\therefore \mu=0, \sigma^2=\frac{1}{N}$$

- This is a good initializⁿ technique.
- For solving vanishing problem, we use ReLU, as its derivative is 1 (0 or 0)

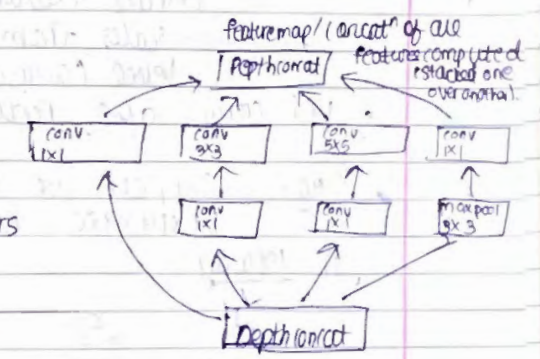
→ GoogleNet

- It has 22 layers with parameters (Conv. + FC)
- With maxpool, it has 27 layers
- There are 9 ~~to~~ identical modules / units
↳ called Inception modules

- Inception modules calculate features of the network.

Inception module:-

- There are 9 such modules
- They use 1x1 filters



- The depthconcat then stacks up of all the conv layers to give feature map / o/p

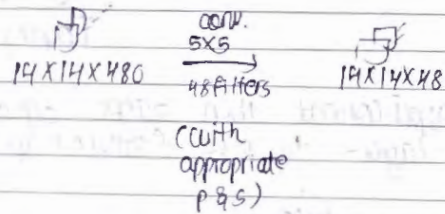
→ Inception module

- Computing 1x1, 3x3, and 5x5 convolutⁿs within the same module of the network.
- Network inside network.
- Covers bigger area, at the same time preserves fine resolution for small information on the images

- uses conv. kernels of different sizes in parallel from the most accurate detailing (1x1) to a bigger one (5x5).
- 5x5 covers a big receptive field; so higher level feature detection, where 1x1 detects fine grain features.
- From 1x1 to 3x3 to 5x5 we extract features on different scales, from minute to higher level (coarse features)
- 1x1 conv. also reduces computation

eg- suppose we want to go from $14 \times 14 \times 480$ to $14 \times 14 \times 48$.

1) 1st way

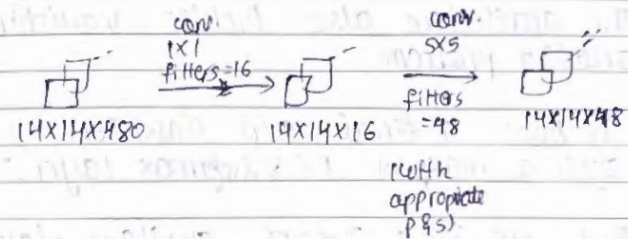


$$\therefore \text{total ops} = (14 \times 14 \times 48) \times (5 \times 5 \times 480)$$

\downarrow for each o/p cell \downarrow the 5x5 filter convolved over 480 channels

$$= 112.9 \text{ million}$$

2) 2nd way

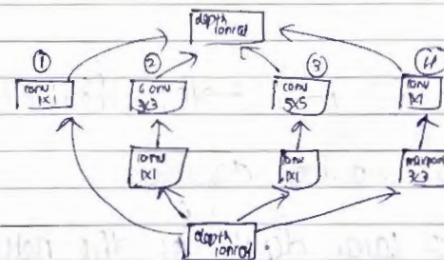


$$\therefore \text{total ops} = \text{for } 1 \times 1 + \text{for } 5 \times 5$$

$$= (14 \times 14 \times 16) \times (1 \times 1 \times 16) + (14 \times 14 \times 48) \times (5 \times 5 \times 16)$$

$$= 1.5 + 3.8 = 5.3 \text{ millions}$$

Thus, no. of ops significantly reduced, and info. was still preserved as 1x1 filters detected fine features



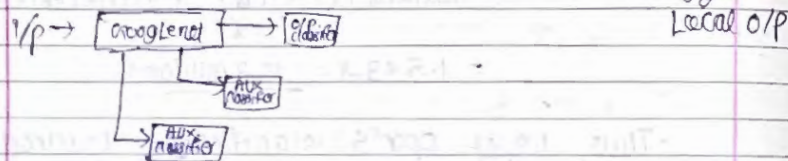
- o/p of ①, ②, ③ & ④ are stacked along the channel dimension to get final o/p of module
- So, we have a multilevel feature extractor
- In the full architecture, there are 9 such inception modules

- The top 5 error rate is less than 7%.

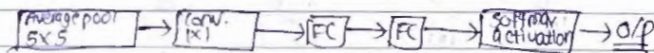
→ The architecture also tackles vanishing gradient problem

- We have a final o/p classifier after a maxpool, FC & softmax layer.

- But we have 2 more auxiliary classifiers in between the architecture



- These aux. classifiers also have various layers:-



- There are 2 such classifiers

→ Due to large depth of the network, ability to propagate gradients back through all the layers was a concern.

- Auxiliary classifiers are smaller CNNs put on top of middle inception modules

Aux. classifiers in the middle exploits

the discriminative power of the features produced by the layers in the middle

- They also give classification error.

So, while backpropagating error, we add the 2 auxiliary losses to the total loss.

$$\therefore \text{total loss} = \text{final classifier loss} + \text{aux. 1 loss} \times 0.3 + \text{aux. 2 loss} \times 0.3$$

- These 2 losses are given a weightage of 0.3

- Then the model is trained.

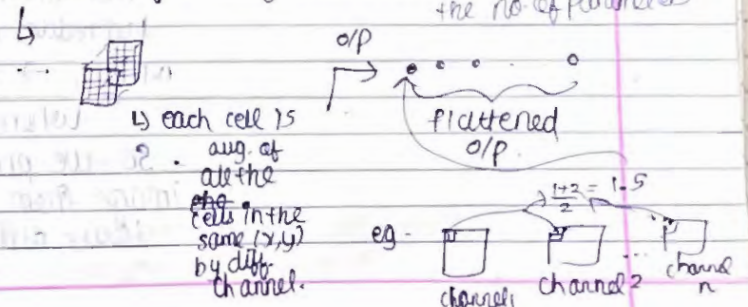
- Aux. classifiers are used only during training time, to solve vanishing gradient problem

during testing / inference time, o/p of aux classifiers are discarded & only final classifier's o/p is considered.

→ Extra points

- The model also used dropout rate of 0.7

- They also used global average pooling instead of FC layers. → This reduced the no. of parameters



• so we had 4 types of layers:-

- 2) stem layers

2) stem layers, like the general conv.-maxpool

2 conv. + maxpool \rightarrow reduce image size.

(four normal layers are called stem layers in GoogLeNet)

- 30) Inception modules

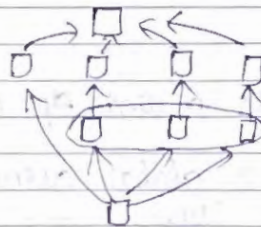
(9 modules)

- 4) Auxiliary classifiers

- Finally, avg. pooling, dropout, softmax etc.

- The multiple local networks work parallelly

* \Rightarrow In inception module:-



→ These are 1x1 filters,
our bottleneck layers,
who ~~maintain~~ maintain
maintain image size,
but reduce channels
 $n \times n \times c \rightarrow n \times n \times c_2$

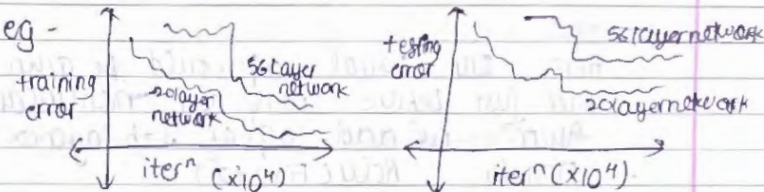
where $n_{c2} < n_{c1}$

So, we project image from higher to lower dimension.

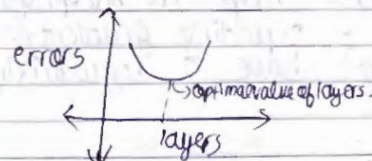
→ Resnet (2015 winning ILSVRC)

- Inductively, we may think that deeper networks are better, as they are more complex & have more parameters to learn.
- But this is not the case
- Also, we expect that deeper network may perform better during training (as even if overfit occurs, then training accuracy ↑), but this is not the case

• eg -



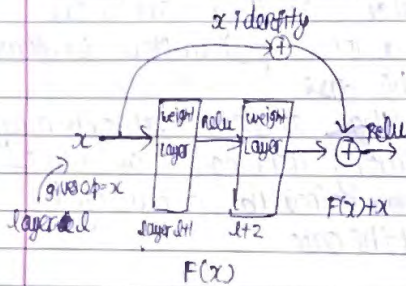
- Here, we observe that for 56 layer network, the training error is more than 20 layer network, ~~also~~ even for testing error is more.
- So, our hypothesis that deeper networks are better fails.
- In reality,



- The failure of this hypothesis is due to vanishing & exploding gradient problem.

So blindly \uparrow layers is not good.

→ The ~~Res~~ Residual Block

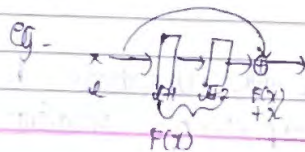


Here, our usual o/p would be given by $F(x)$.
But just before using our nonlinear function, we add o/p of l th layer (x).
Finally, $\text{ReLU}(F(x) + x)$

So effectively

This is a residual connection, or skip connection

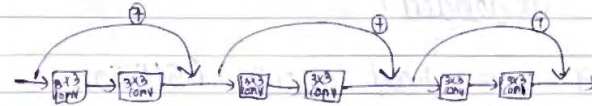
- The original Resnet architecture used 152 layers
- The skip connects help in backprop (help vanishing & exploding gradient problems)
- They also have a regularizing effect



In $F(x) + x$, if we make $F(x) = 0$, then o/p is x .

So, the weights of $l+1$ & $l+2$ layers did not contribute to the effective o/p.
So, overfitting due to these weights was prevented.

- So, skip connections help the gradients to flow & thus makes training possible.
- Thus, with skip connections, even if the network is deep, it can behave as shallow network.
- So, the 152 layer original layer also works



$$H(x) = F(x) + x$$

$$\therefore F(x) = H(x) - x$$

we try to learn $F(x)$

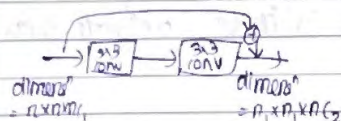
- Also, we are not multiplying, only adding.
- So, gradients will flow properly, & no vanishing or exploding gradient problems occurs.

- Less than 1% error in ImageNet Database.
- They also used bottleneck Layer (1×1 conv) to keep the no. of computations sustainable.

→

* SKIP CONCEPT: solve overfitting
 $H(x) = F(x) + \lambda$, $\therefore F(x) = H(x)$
 Here: $H(x)$ is a complex function
 which when take I/P, gives o/p similar
 to without $H(x)$
 so weights must be such that they produce same I/P.

One point needs to be taken care of -



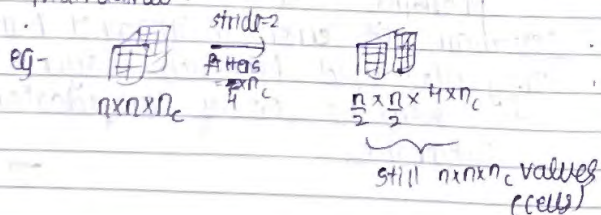
During backprop, we can skip & propagate so, no vanishing gradient.

- So, if dimensions do not match, then additⁿ not possible
- So, using various heuristics, this concern is solved.
- Also, how many layers to skip is important (according to the problem this is decided)

* Resnet 2015 had about 65 Million parameters

same imp points

- We have 3 residual connect^{ns}, each with 2 filters of 3×3 . Their o/p's are stacked.
- Using strides=2, the total values are maintained.



error = $\text{RMSE} + W_2$
 (weights)²
 now, for derivⁿ, we have

a weight factor subtract from respective wgt.
 so, it is a penalty

* weight decay = kind of regularizⁿ.
 When overfitting, weight doesn't vary.
 so, we bound weight but adding weight to o/p's, so they get regularized.

- No use of FC layers as it has more parameters, so Global Avg. Pooling was used.
- They used bottleneck layer (1x1 conv.) if network deeper & more parameters (more than 60 layers)

- Stochastic Gradient Descent (SGD) + momentum = 0.9
 - momentum 0.9 is multiplied in training gradient (to control displacement in directⁿ of Velocity)
- minibatch size = 256
- use of batch normalizⁿ after each conv.
- Xavier/2 weight initializⁿ
- learning rate initially = 0.1, but if validⁿ error increases, the learning rate divided by 10
- Weight decay = 10^{-5}
 (similar to ML regularizⁿ)
 (like λ regularizⁿ const.)

Inception v4 = GoogleNet + Resnet
 (highest accuracy)

- VGG16 & 19 have more memory consumpⁿ
- GoogleNet consumes less memory. (less parameters)
- Resnet has less memory consumpⁿ & high accuracy.

class
NOTE
outputs.

- We preserve Identity of x .
So, $H(x)$ is a complex functⁿ.

$$F(x) \\ x_1 \rightarrow x_2$$

$$F(x) \\ x_1 \rightarrow \textcircled{+x} \rightarrow x_2$$

So, we have a very complex functⁿ,
which nearly gives same o/p as
without previous layer, without
distorting much data.

This kind of functⁿ are approx^m
function.

- No. of params do not change by
adding skip connections.

In paper's diagram,

----- dotted = dimensⁿ change.
→ solid line = same dimensⁿ.

- They do not use too many maxpooling.
- To reduce dimensⁿ, use of conv. with more strides.
- filters used - 64, 128, 256, 512.
- skip connectⁿ - every 2 conv. layer.

- while changing dimensⁿ, for downsampling,
we use ~~stride~~ bigger stride.
- use of bottleneck layer ~~1x1~~ (1x1) conv.
to downsample (eg - 56×56 to 28×28
using 1x1 conv.)
- No. of operⁿ remain same
(~~28~~ $56 \times 64 = 28 \times 128$)
half double

Recurrent Neural Network

- It works very well with sequential data
- we need sequence modelling.

eg - speech recognitⁿ

- long modelling
- machine transⁿ

- named entity recognitⁿ
- sentiment classifier
- video activity analysis

→ why RNN?

eg - NLP (natural lang. processing)

↳ text data i/p (eg - spam classifier)

↳ we can use

- Bow (bag of words), TF-IDF, word2vec

They convert text to vectors.

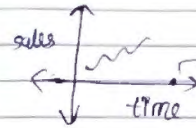
word1 word2
sentence 1 0.3 0.9 ...
vector

now we can use ML algo. (like Naive Bayes).

- Here, the model understands the words, ~~but~~ the sequence informⁿ has been discarded

- so, accuracy may be low.
- we have seen many Chatbots / voice assistant answer us in proper sequence.

eg - Time series data



- for predictⁿ, we want the sequence of previous sales
- so, here RNN can be used (this will give better accuracy than usual ML models)

eg - google image search



→ explain this image in text.
- needs RNN

eg - English to Spanish

- RNN needed to maintain sequencing

→ If we use feed forward NN (ANN) then we need to pad the words to make a fixed i/p size.

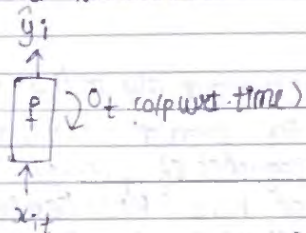
eg - Pankaj

↓ lives in Munich

Pankaj lives in Munich DE

- while feeding these sentences to the ANN, we need to pad.
- When data is very large, ANN may not work properly, as context may change ahead in the sentence, and previous weights may not be useful.
- For context knowledge, we need previous ops, but ANN or CNN doesn't store it.

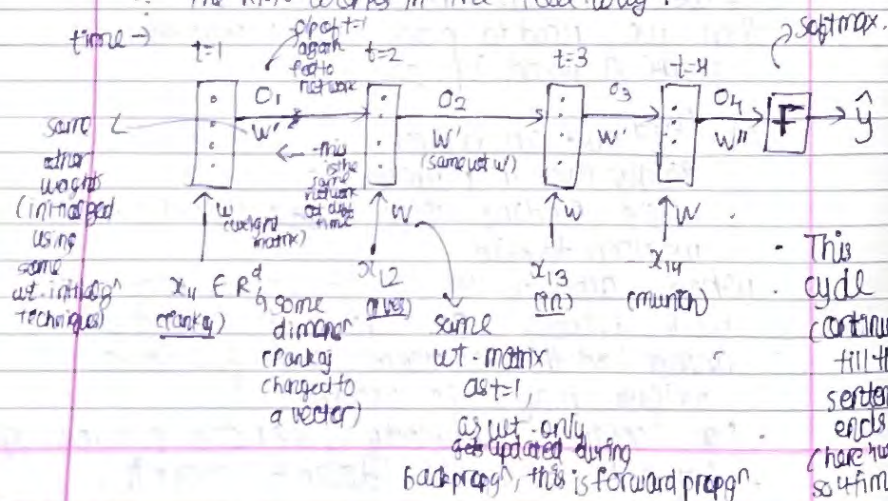
Recurrent Network:-



- The main idea is that we basically feed the o/p of the previous time stamp again to the neuron/layer.
- Suppose we are doing sentiment analysis.

sentence = Panhaji lives in Munich
 $x_{1,1}$ $x_{2,1}$ $x_{3,1}$ $x_{4,1}$ → let it be sentiment = "happy"
 4 words

∴ The RNN works in the following way:-



activation function

$$O_1 = f(x_{1,1} \times w)$$

$$O_2 = f(x_{1,2} \times w + O_1 \times w')$$

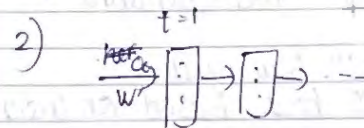
→ we can see that O_2, O_3, O_4 preserves the seq. info.

$$O_3 = f(O_2 \times w + O_1 \times w')$$

$$O_4 = f(O_3 \times w + O_2 \times w')$$

This o/p O_4 is then finally given to a softmax function F , (O_4 multiplied with some w'')
 (permuted sentiment classes like happy, sad etc.)

- Then we get find predicted \hat{y} .
- We calculate loss ($\hat{y} - y$)
- Then we backpropagate this loss, to minimize it. (also, we maintained sequence info.)
- NOTE - 1) we represented each word using in the vocabulary using a vector, this needs special techniques



here, we need O_0 as dummy i/p, as all the words have some i/p.
 O_0 can be 0 padded values, or some randomly initialized values.

$$\Rightarrow O_1 = f(x_{1,1} \times w + O_0 \times w')$$

- But here the numbers are not normalized.
- $a = 0$, zebra = 1000

↳ not in same range

So, we may use one hot encoding.

a	0
ant	1
happy	25
zebra	1000

→ For happy, [000 1 00]
 (word)
 25th position

$$a = [100 \dots 0]$$

Adv: of this method

easy to use

disadv.

- very huge dimension for big corpus
- doesn't have any meaning

eg - I am happy, I am glad
 similar, but we can't detect these semantics

2nd way)

- We use featureset
- Each cell of the matrix has normalized value (in range 0 to 1 or -1 to 1 depending on applicⁿ).

eg- words = apple, mango, elephant.

we define some features for these words

word \ feature →	Is fruit	Is eatable	Is animal
apple	0.9	1	0.01
mango	0.92	1	0.01
elephant	0.01	0.01	0.9

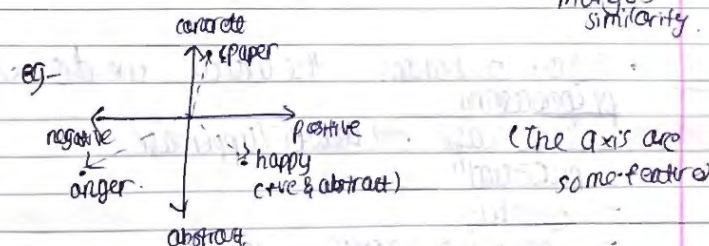
now, apple's → [0.9 1 0.01]

↳ this vector

- How many features to choose is a hyperparameter (usually 300 features for practical applicⁿ).
- Now, we can see that Apple & mango are similar (cosine similarity)

so now if sentence → I like apple juice
 I like mango juice

we can predict this based on apple's & mango's similarity.



eg- so, every word can be shown as a n dimensional vector (n = no. of features)

Adv. of this method

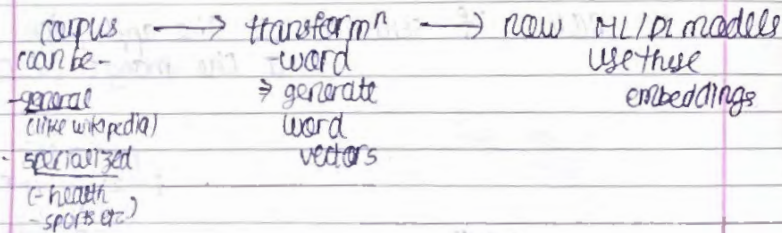
- Low dimension
- meaning given to encoded words

- Some ML models for word embeddings:-
 - word2vec (google, 2013)
 - Continuous Bag of words (CBOW)

- DL models - GPT2, BERT, ELMo

→ The task of creating our own embedding is self as well as unsupervised.

∴ Flow of the method:-



NOTE - Before embedding the words, we do some preprocessing:-

- Letter case → Lower/uppercase
- punctuation
- numbers
- special characters (@, #)
- special words

eg- 😊 : happy : ☹️
symbols

New apple → $\begin{bmatrix} 0.9 \\ 1 \\ 0.01 \end{bmatrix}$

fruit edible animal

But how do we get these values?

- How do we generate these rel values / feature values.

- 2 solutions

- 1) - either take readymade weight
- 2) - or learn them from data

- For learning, there are 2 methods-

- 1) CBOW (Continuous Bag of words) } classical models
- 2) skip gram model

② Skip gram model → Predict context word, given centre word

eg- sentence = The cat is eating the mouse it caught.

 context centre context

- We have a window ~~that~~ (with a window size let's say 2).
- So, we have a centre word.
- To the left side, we have a window with specified size, i.e. 2, & similarly to the right we have the window.
- The window have our context words, & the centre word is in the middle.
- So, the centre has the full relⁿ:-
 (cat, is), (cat, mouse), (cat, eat), (cat, is)
 - combinⁿ of centre with each window word forms arelⁿ.

- Choosing ~~off~~ window size is a hyperparameter
- If window size is too large, then too many rel's which may not be proper.
- If too small, then we may miss out imp rel's
- Usually window size is 2-5.

→ The entire corpus, ~~is~~ is our raw data.
 From that, we generate a training dataset.
 For ML model to learn embeddings.

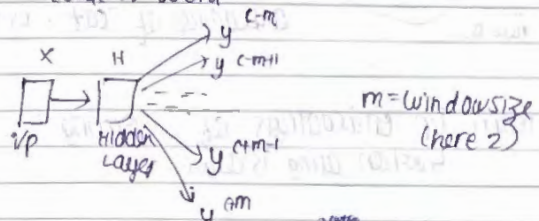
eg- corpus = The cat is eating the mouse it caught.
 (1st time) (context words)

- next time, cat becomes centre, & "the", "is eating" becomes corpus
- So, we keep on sliding this window. & we get rel's.
 eg- 1st time, rel's:- (the, cat), (the, is) & similarly, we do for entire corpus, from 1st to last words
- These rel's are our training examples.
- We can see that we can learn rel's through this, eg- eating & mouse are related, mouse & caught are related, etc

[Corpus is usually very large (billions of words)]

Training process → we feed centre as i/p word

- Let the 1st example / relⁿ generated be (the, cat)
- Here "the" is the centre word, & "cat" is the context word.

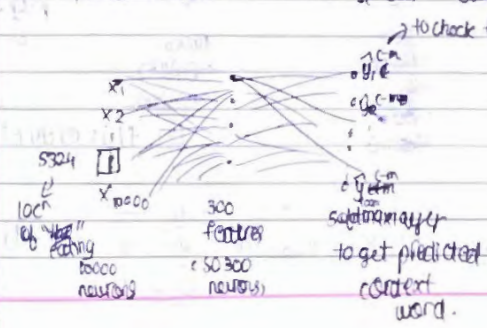


→ The cat is eating the mouse it caught.
 (centre) (context words)

eg- when eating is centre, if window (m=2), then we get 2m context words

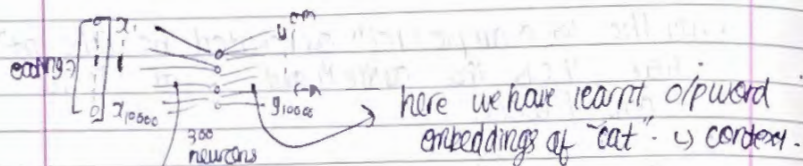
So, when we feed y^c as "the", we want o/p as y^{cm} is cat (in that case)

- Now, again, we can't directly feed ~~cat~~ ^{eating} as centre, so we use 1-hot encoding.
- So, for 10000 words in vocab, we create 1-hot vector for the ~~cat~~ ^{eating} & cat.



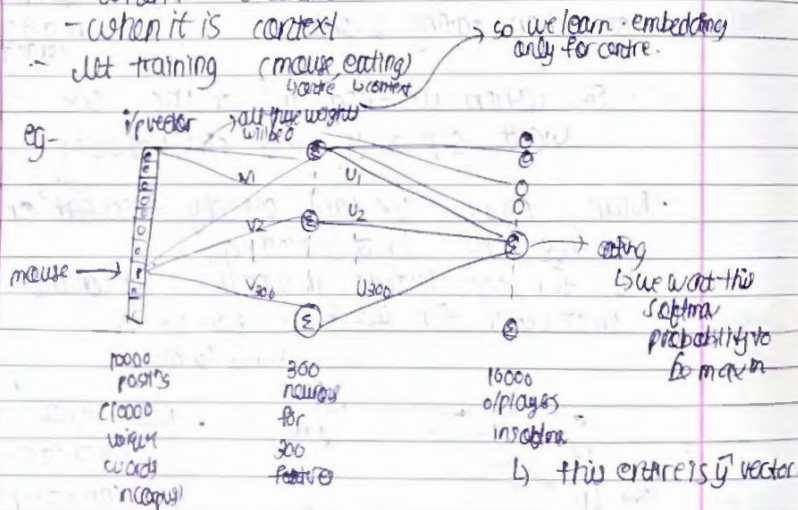
- Now we have actual 1-hot encoding of cat.
- We can compare it with this y^{cm} vector (of same dimensⁿ) to calculate loss.

xx. for each example, we learn 2 word embedding



here we learn 1/p embeddings of "coating"
↳ when coating is center.

- So, we learn centre & context embeddings
- So, for every word, we will learn 2 embeddings:
 - when it is centre
 - when it is context
- Let training (mouse eating cheese) → so we learn embedding only for centre.



Ex 4a

$$\text{prob. of eating} = [V_1, V_2 - V_{300}]^T \begin{bmatrix} U_1 \\ U_2 \\ U_{300} \end{bmatrix} = V^T U$$

$$J = i/p$$
$$\sigma = \alpha/p$$

2 word is here missing

$$\therefore \text{prob. of eating by softmax} = e^{\frac{v_{T, w} u_w}{\sum_{w=1}^{10000} e^{v_{T, w} u_w}}} \xrightarrow{\text{words}} \text{eating}$$

$= P(W_0/W_I)$
↳ This is for pair

- For next training example, we repeat this.
- Also, if prob of eating is not max, then we calculate loss using y (original ans).
- We then backpropagate this loss.

① CBOW (continuous bag of words) → predict entire word given context words

- Alternative to skipgram
- Predict centre word given context word window.
- Now we are given context words, we find prob. of centre words.

eg - The cat is eating the mouse it caught.
context centre context

- We had 2m context words for window m_s
- But here the 4 words "eating the ~~meat~~ it caught" is treated as continuous - conti. bag of words.

∴ training example \rightarrow ((casing, the, it, caught), mouse)

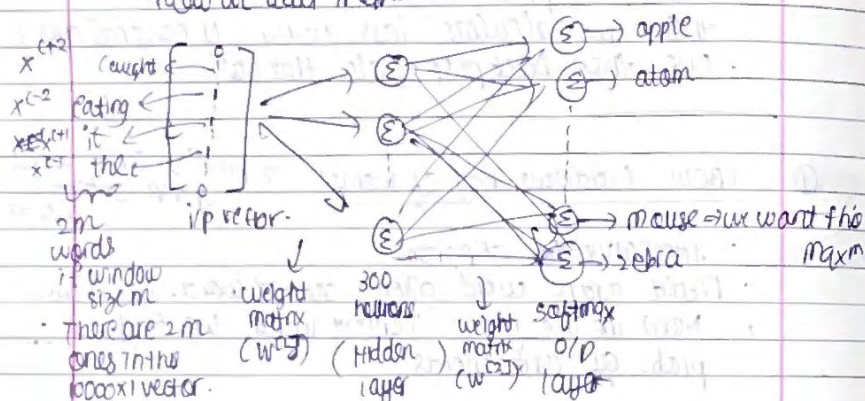
order of →
(cont. words doesn't matter,
i.e. (the, it, eating, caught), mouse)
is also valid (y all other combin^{ns})

\therefore i/p vector \rightarrow

$$\text{caught} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}, \text{eating} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}, \text{it} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \text{the} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

one hot ~~vec~~ encoding of words

Now we add them.



\rightarrow In this softmax o/p, we want probability of "mouse" to be the highest.

$$\therefore V^{(2)} = W^{(2)} X^{(2)}$$

vector for word $X^{(2)}$ \rightarrow 300x1
Gauged

as $W^{(1)}$ is 300x1000

$X^{(1)}$ is 1000x1

$$\therefore W^{(1)} X^{(1)} = 300 \times 1$$

similarly, $V^{(1)} = W^{(1)} X^{(1)}$

$$V^{(1)} = W^{(1)} X^{(1)}$$

$$V^{(2)} = W^{(2)} X^{(2)}$$

$$\text{net } \hat{V} = \frac{V^{(2m)} + \theta V^{(2m-1)} + \dots + V^{(2m-1)} + V^{(2m)}}{2m}$$

$$= \hat{V} \rightarrow 300 \times 1$$

\rightarrow i/p context word vector.

\Rightarrow This is not i/p

Now for final layer, we get o/p of 1000x1, where we want mouse to get max value.

$$\therefore \text{mouse} = W^{(2)} \hat{V} = V_c'$$

final o/p \rightarrow 1000x300 \rightarrow 300x1

$$= 1000 \times 1 \rightarrow \text{o/p}$$

we want softmax of z to have max value of mouse
 \rightarrow in those 1000 values

\Rightarrow so, we want to maximize -1

$$\text{Prob} \propto \frac{P_c W_c}{W_{c-m}, W_{c-m+1}, \dots, W_{c-1}, W_{c+1}, W_{c+m}}$$

\downarrow
center word

\rightarrow given these context words (2m words)

\therefore maximize this
 \Rightarrow or maximize \log of the
 \Rightarrow or minimize $-\log(P(w_c / w_{c-m}, \dots, w_{c+m}))$

\hookrightarrow This is our objective function.

$$J = -\log P(v_c / \hat{v})$$

\downarrow (i/p)
o/p

- minimize J
- This was for 1 training example
- we do this for all training examples.

[NOTE] \rightarrow • In both the models we saw, the embeddings from i/p layer to hidden layer give us our word embeddings.