

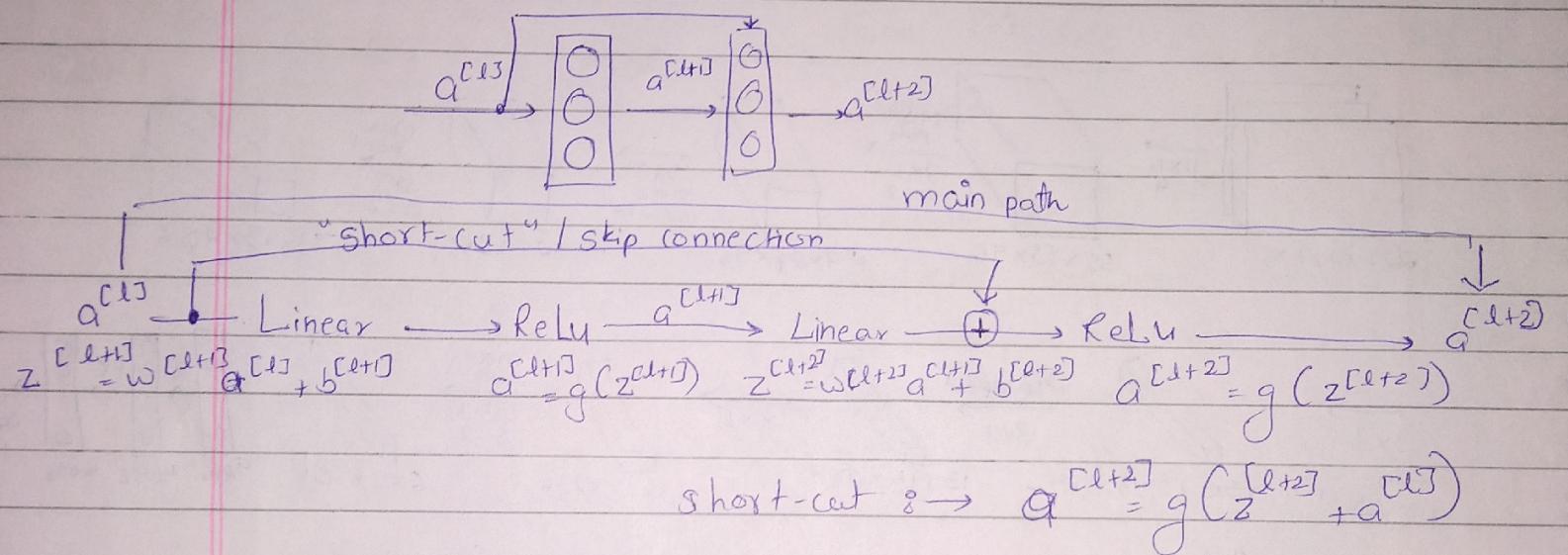
Residual Networks (ResNets)

Very big NNs are difficult to train, due to vanishing and exploding gradient descent problem.

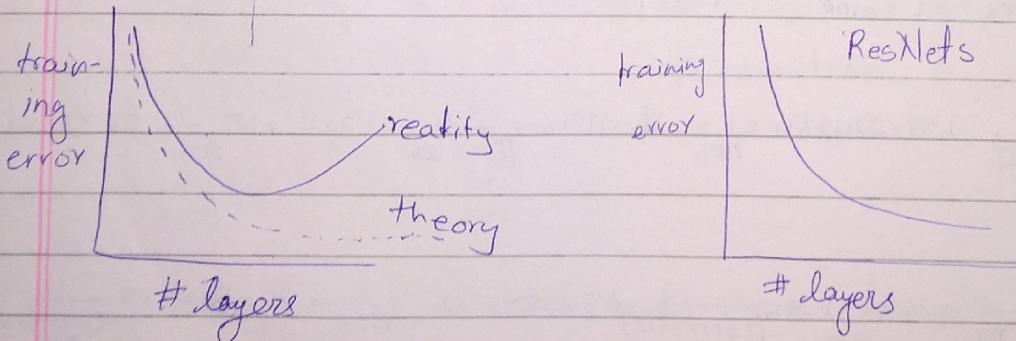
ResNets

ResNets are built on residual blocks

A



In theory, training error should decrease as we make our NN bigger.



This skipping of connections helps solve gradient descent problems. But ResNets are not ~~very great~~ at very big NN

$$X \rightarrow [\text{Big NN}] \rightarrow a^{[L]}$$

$$X \rightarrow [\text{Big NN}] \xrightarrow{\text{ReLU}} \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} \rightarrow a^{[L+2]}$$

Identity f^n is easy for Residual block to learn.

$$a^{[L+2]} = g(z^{[L+2]} + a^{[L]}) = g(\omega^{[L+1]} a^{[L+1]} + b^{[L+2]} + a^{[L]}) = g(a^{[L]})$$

due to L2 regularization, $\omega^{[L+2]}$ and $b^{[L+2]}$ will be close to zero.

Reason why big NN don't work, is because it is very difficult to choose parameters that learn even the identity f^n . And so, many layers make the result worse.

Since resNets can learn identity f^n easily, it is guaranteed that extra layers won't hurt.

Here, we assume that $z^{[L+2]}$ and $a^{[L]}$ have same dimension So, a lot of same convolutions are used in resNets. In case they are not of same size, one could add WS matrix to balance it. (Padding is used.)

1x1 Convolution

$$\begin{matrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{matrix} \quad * \quad [2] = \begin{matrix} 2 & 4 & 6 \\ 8 & 10 & 12 \\ 14 & 16 & 18 \end{matrix}$$

$6 \times 6 \times 1$

$$\begin{matrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{matrix} \quad * \quad [1 \times 1 \times 32] = \begin{matrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{matrix} \quad \text{channel}$$

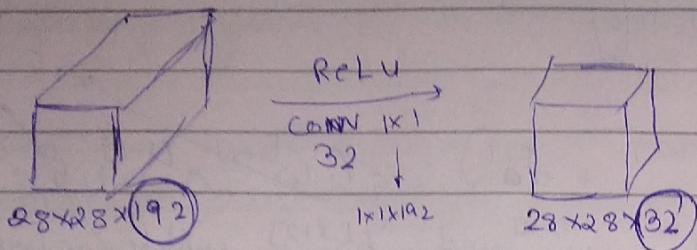
$6 \times 6 \times 32$ $6 \times 6 \times \# \text{filters}$

each element of 32 matrix is multiplied with each channel element of $1 \times 1 \times 32$ and then summed to get a value on output

Called One by One Convolution or Network in Network

Using 1x1 Convolutions

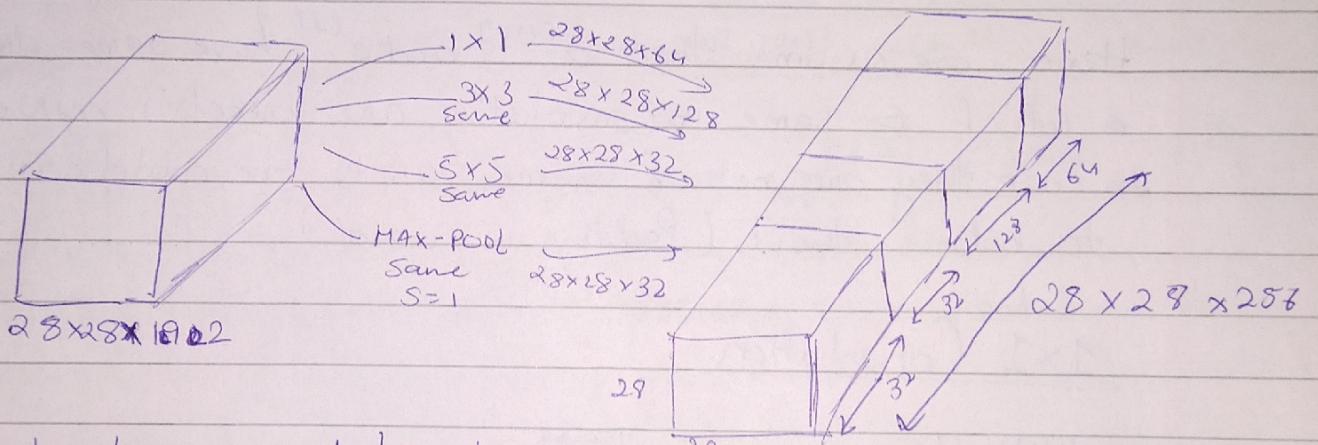
In order to shrink No



Thus, it can help to increase or decrease number of channels.

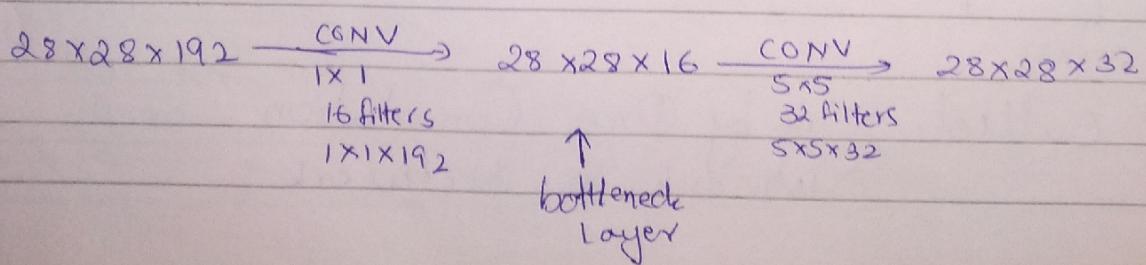
Inceptⁿ Network

When you are confused which layer to use, you use all, and stack them up. This stacked thing is your i/p to model. So that the model can choose whichever parameters it wants to use.



But, the problem here, is computational cost. For e.g., 5×5 , we have 32 filters. So, the total no of multiplications done will be $28 \times 28 \times 32 \times 5 \times 5 \times 192 = 1.20 \text{ m}$. But in 1x1 convolution, it reduces these calculatⁿ by a factor of 10. i.e. $1/10^{\text{th}}$ of 120 million.

So, what we do instead is:

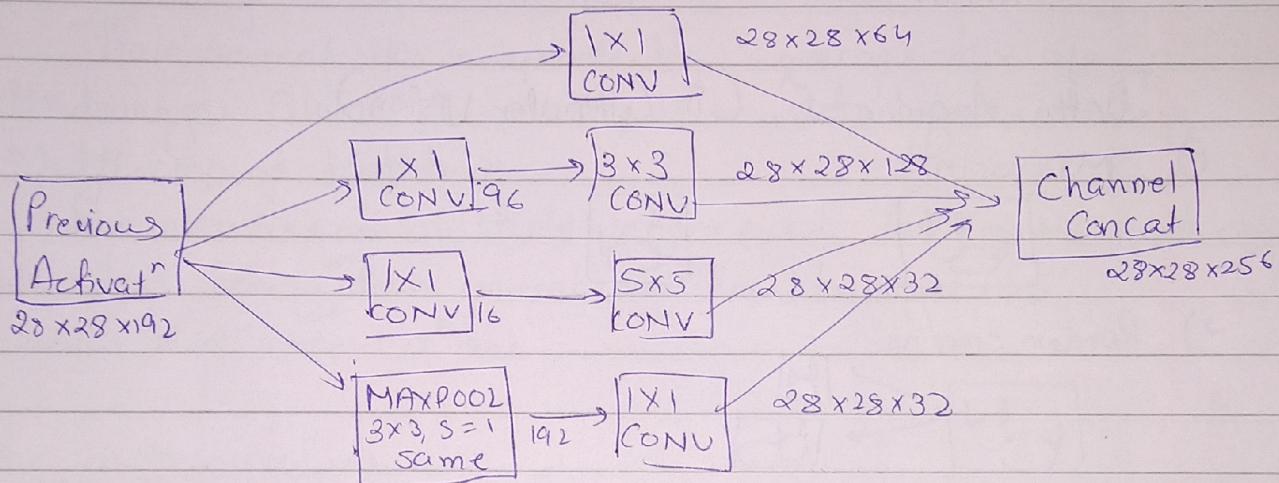


Computation cost:

$$\begin{aligned} 28 \times 28 \times 16 \times 192 &= 2.4 \text{ million} \\ + (28 \times 28 \times 32) \times (5 \times 5 \times 16) &= 10.0 \text{ million} \end{aligned} \quad \left. \begin{array}{l} \text{?} \\ = \end{array} \right\} = 12.4 \text{ million}$$

This shrinking does not hurt the ~~layer~~, ~~in~~ performance but does save a lot of computation.

So, what we actually do is:



These inception blocks can be placed inside the n/w.

There are certain side branches to the inception blocks.

These are actually softmax o/p's they try to predict the o/p layer label. They are present in hidden layers.

Thus, even hidden layers have a chance to contribute to o/p. This appears to have a regularizing effect on incept n/w and helps prevent overfitting.

Transfer learning

More the data, more the layers you freeze.

By freeze, we mean to not change any parameters of that already trained layers.

data size

small

more big

lot of data

what to do

remove last 1-2 layers, add yours & train yours

" " few " add yours & train yours

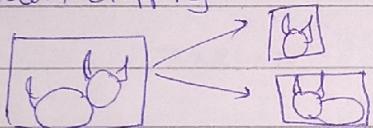
train entire n/w. Previously trained weights work as initialization

Data Augmentation (in Computer Vision)

1) Mirroring



2) Random cropping



3) Rotation

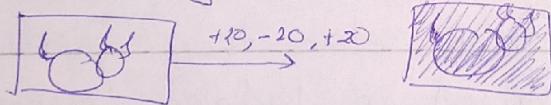
4) Shearing

5) Local warping

6) Color shifting

} Not used in practical world

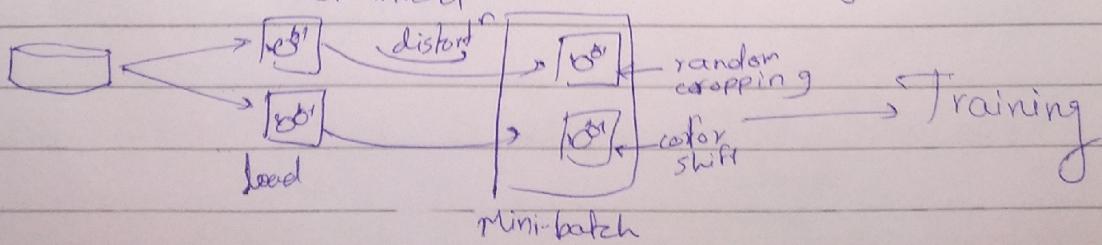
In practice, R,G,B are from a distribution that could be quite small



PCA color augmentation :- Finds which is more out of R,G,B and then adds more of it while subtracts the other.

ref: AlexNet paper

⇒ ④ Implementing distortions during training



~~In~~ do we

We have two sources of knowledge

- 1) Labelled data
- 2) Hand engineered features / new architecture / other components

If you have less of ①, you need to do ②

Tips for doing well on benchmarks / winning competitions

- 1) Ensembling

Train several n/w independently, average their o/p

disadvantage: Slows computⁿ, not used in production

- 2) Multi-crop at test time

Run classifier on multiple versions of test images and avg. results

adv.: less memory ^{used} than ensembling due to single n/w

disadv.: More computⁿ time

Use open source Code

Use architectures of n/w's published in literature

~ Open src. implementatⁿs if possible

~ pretrained models & fine-tune on your dataset