

$$\frac{\partial l}{\partial k_{nm}} = \sum_i^{\text{Size-out}} \sum_j^{\text{Size-out}} \frac{\partial l}{\partial C_{ij}} \cdot \frac{\partial C_{ij}}{\partial k_{nm}} = \sum_i^{\text{Size-out}} \sum_j^{\text{Size-out}} \delta(i,j) \cdot \frac{\partial C_{ij}}{\partial k_{nm}}$$

$$C_{ij} = \sum_n^N \sum_m^M Q^{l-1}(i+n, j+m) k(n, m)$$

$$= \sum_i^{\text{Size-out}} \sum_j^{\text{Size-out}} \delta(i,j) \frac{\partial \sum_n^N \sum_m^M Q^{l-1}(i+n, j+m) k(n, m)}{\partial k_{nm}}$$

$$k_{nm} = k(n, m)$$

$$= \sum_i^{\text{Size-out}} \sum_j^{\text{Size-out}} \underbrace{\delta(i,j)}_{(3,3) \atop (I-N+1, J-M+1)} \cdot \underbrace{Q^{l-1}(i+n, j+m)}_{(I, J) \atop (5, 5)} \quad (\text{only when } n' = n \atop m' = m)$$

interpretation: sliding $N \times M$ times to generate update

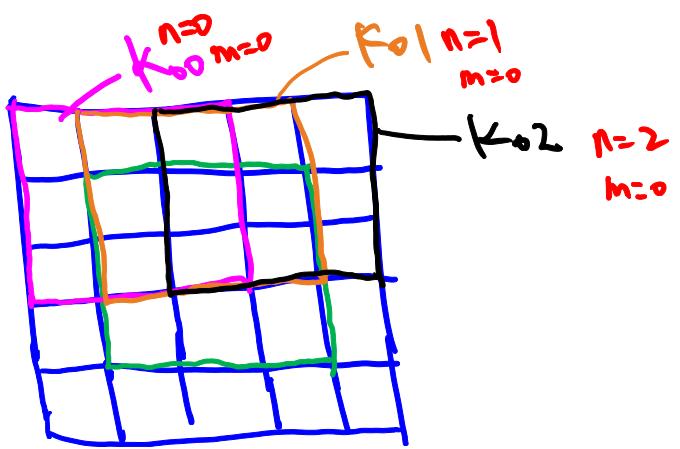
Coefficients

$I \times J$ is the conv output size

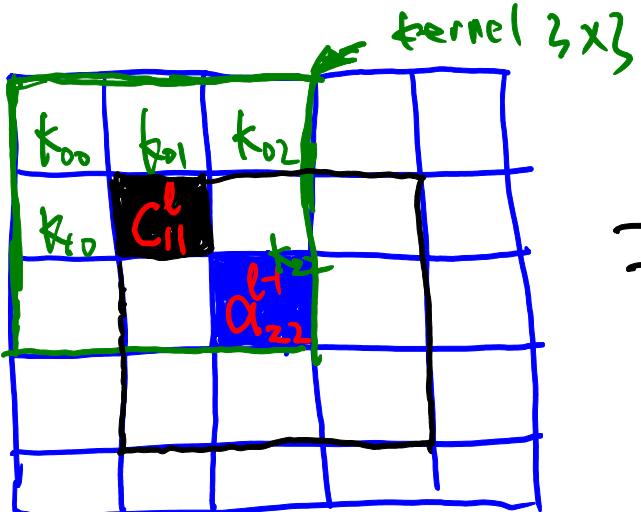
the size of $\delta(i,j)$ is less than that of $Q^{l-1}(i,j)$ which is the input.

$$\text{size-in} = I \times J$$

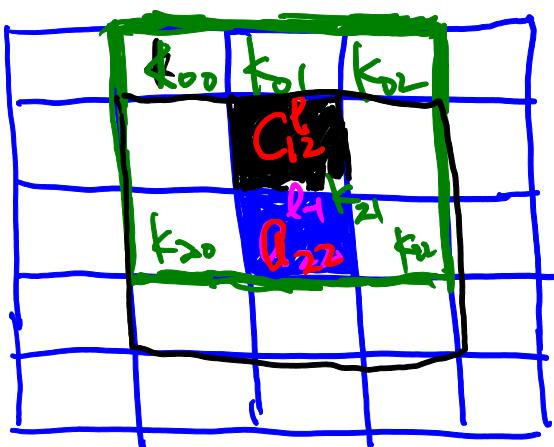
$$\text{size-out} = (I - N + 1) \times (J - M + 1)$$



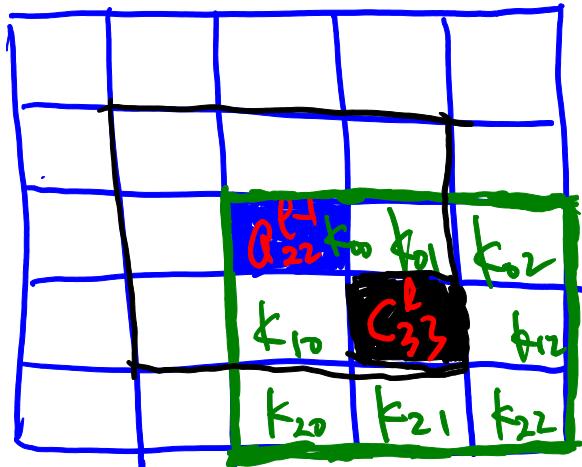
$K(-n, -m)$ is the default kernel layout we store in memory!



$$\frac{\partial l}{\partial a_{22}^{l+1}} = k_{22} * \frac{\partial l}{\partial c_{11}^e}$$



$$\frac{\partial l}{\partial a_{21}^{l+1}} = k_{21} * \frac{\partial l}{\partial c_{12}^e}$$



$$\frac{\partial l}{\partial a_{23}^{l+1}} = k_{23} * \frac{\partial l}{\partial c_{23}^e}$$

$$\begin{aligned} \frac{\partial l}{\partial a_{22}^{l+1}} &= k_{22} * \frac{\partial l}{\partial c_{11}^e} + k_{21} * \frac{\partial l}{\partial c_{12}^e} + \dots + k_{00} * \frac{\partial l}{\partial c_{33}^e} \\ &= k_{22} * \delta^{(1,1)} + k_{21} * \delta^{(1,2)} + \dots + k_{00} * \delta^{(3,3)} \end{aligned}$$

	0	1	2	3	4	5	6
0	0	0	0	0	0	0	0
1	0	k_{11}	k_{10}	k_{1-1}	M		0
2	0	k_{01}	G _{LL}	k_{0-1}			0
3	0	$k_{11} M$	k_{10}	$Q_{33} P_1$	k_{10}	k_{1-1}	0
4	0				C ₄₄		0
5	0					k_{24}	0
6	0	0	0	0	0	0	0

$k(n, m)$

1	1	-1	1	1
0	1	0	0	0
1	-1	1	0	1

padded output img

$$\frac{\partial \ell}{\partial Q_{ij}^l} = \sum_{n=-1}^{N-1} \sum_{m=-1}^{M-1} \delta^l(i+n, j+m) k(1^{-n}, 1^{-m})$$

$$= \sum_{\substack{n=-N \\ n \neq -1}}^{N-1} \sum_{\substack{m=-M \\ m \neq -1}}^{M-1} \delta^l(i+n, j+m) k(-n, -m)$$

move center of kernel to middle yield:

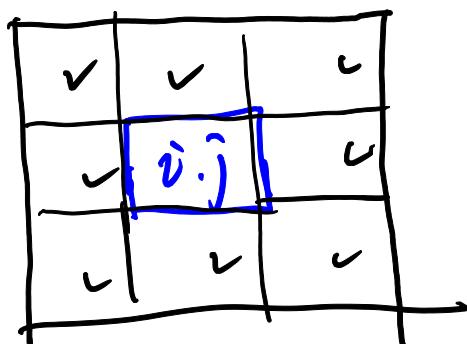
$$= \sum_{n=-N}^{N-1} \sum_{m=-M}^{M-1} \delta^l(i+n, j+m) k(-n, -m)$$

k_{-1}	k_{0}	k_1
k_{-1}	k_0	k_1
k_1	k_0	k_1

$$G_{ij}^l = \sum_{n=-N}^{N-1} \sum_{m=-M}^{M-1} a^l(i+n, j+m) k(-n, -m)$$

$$\frac{\partial l}{\partial a_{ij}^{ft}} = \sum_{i'} \sum_{j'} \frac{\partial l}{\partial C_{i'j'}} \cdot \frac{\partial C_{i'j'}}{\partial a_{ij}^{ft}}$$

a_{ij} is only contributing to output pixels in a $M \times N$ kernel box with its center pixel being (i, j) , which is same as a_{ij}

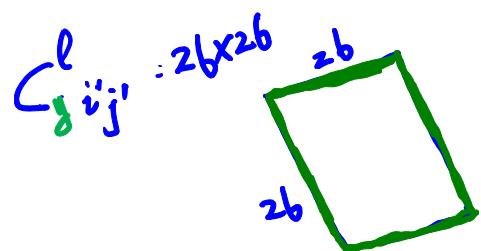
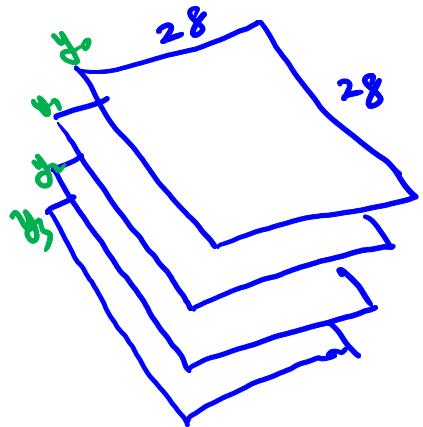


$$= \sum_{\substack{i'=i-\frac{M-1}{2} \\ i'=i+\frac{M-1}{2}}} \sum_{\substack{j'=j-\frac{N-1}{2} \\ j'=j+\frac{N-1}{2}}} \delta(i', j') \frac{\partial C_{i'j'}}{\partial a_{ij}^{ft}}$$

Check gdo DDC notes for further rigorous derivation.

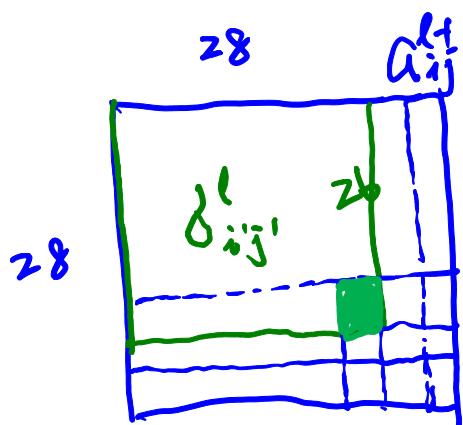
$$a_{ij}^{l-1} = 4 \times 28 \times 28$$

$$k^l(n,m) = 4 \times 3 \times 3$$

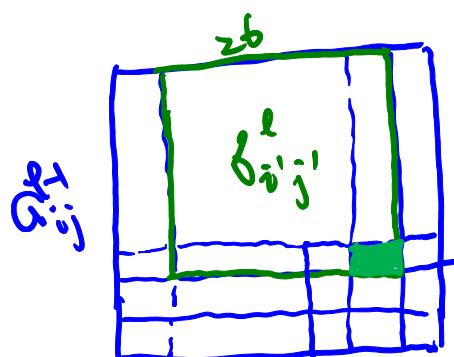


$$\frac{\partial l}{\partial k_{yin}} = \sum_{i'=0}^{25} \sum_{j'=0}^{25} a_y^l(i'+n, j'+m) b^l(i', j')$$

y is the channel index in
input channel or filter channel



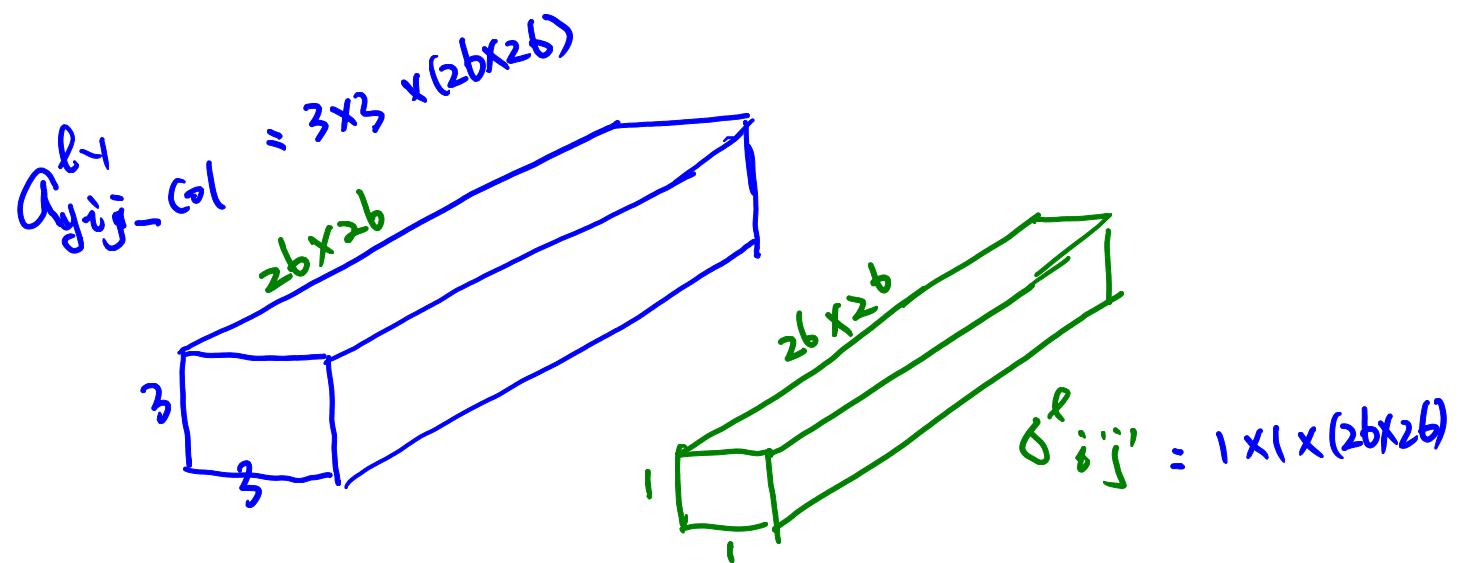
$$= k(0, 0)$$



$$= k(0, 1)$$

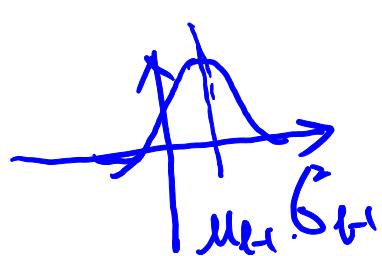
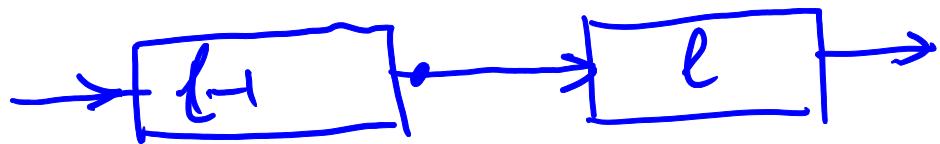
⋮

im2Col for $\frac{8e}{2k_{\text{ynn}}}$

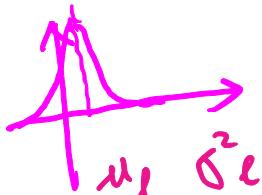


This is for each channel of filter, if filter size is $N \times 3 \times 3$. then we have N such col column.

ON Batch normalization



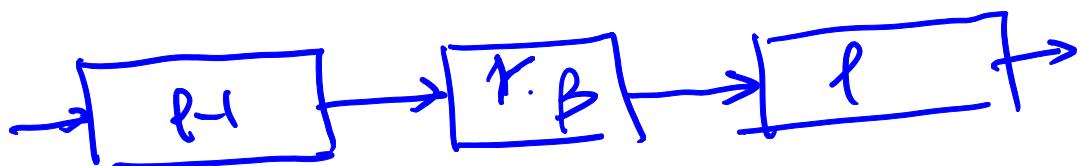
provided from
previous layer



input characteristic

required

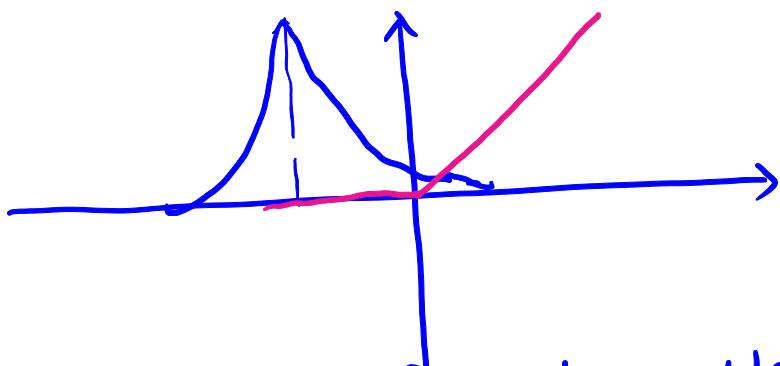
So there is a mismatch between given and requirement which will degrade the system performance
A solution to this is Adding scale & shifting factor to it.



distribution
transformation .

the second benefit from batch normalization
is we can speed up training!

The reason for this is that. We can avoid
output distribution like below



For relu activation function, the output
will be very sparse \nearrow too many zero
 \nwarrow layer with
neurons carries less feature information,
this will tend to cause all zero output in
the layers that follows!

In training without BN we often set
the learning rate to be small to avoid
such sparse situation.

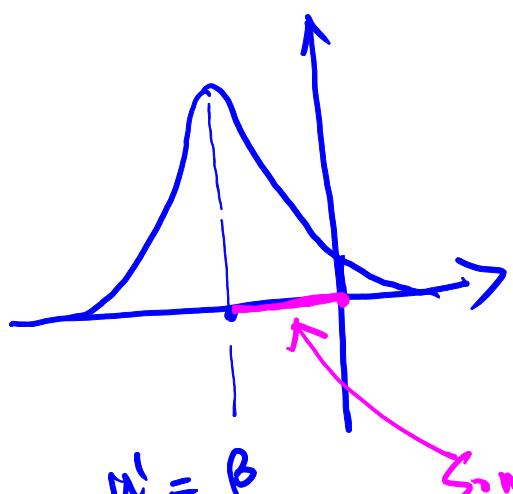
with BN, we tend to move the centroid of the data distribution around zero, which avoids the sparse condition in some degree. so that we can make the learning rate a little bigger.

A suggestion to make the learning rate even bigger

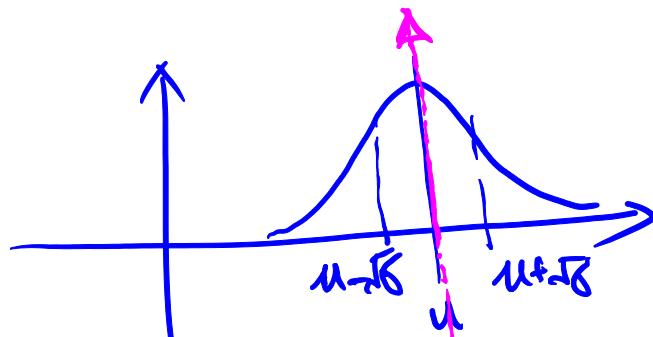
Adding constraint to β . Since we know the negative β tend to move the distribution to the left, we can constrain it so it does not move too much to the left.

Since we know the current output distribution, we can limit it to

$$\beta + \underbrace{\text{some_value} * \sqrt{\epsilon}}_{\text{try to find out this value}} > 0$$



Some-value * $\sqrt{\epsilon}$



β can't be more negative than this.

so in update process:

$$\beta = \beta + df$$

if $\beta < -\text{Some_Value} * \sqrt{\epsilon}$

$$\beta = \text{Some Value} * \sqrt{\epsilon}$$