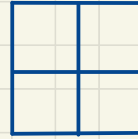Filter 1
Feature Map 1
$A_1$

Filter 2
Feature Map 2
$A_2$

Filter 3
Feature Map 3
$A_3$

# Typical CNN Architecture

$X$ input image

Layers $\rightarrow$ Layers

last conv layer $A^3$

$A^2$

$A^1$

We will focus on this.

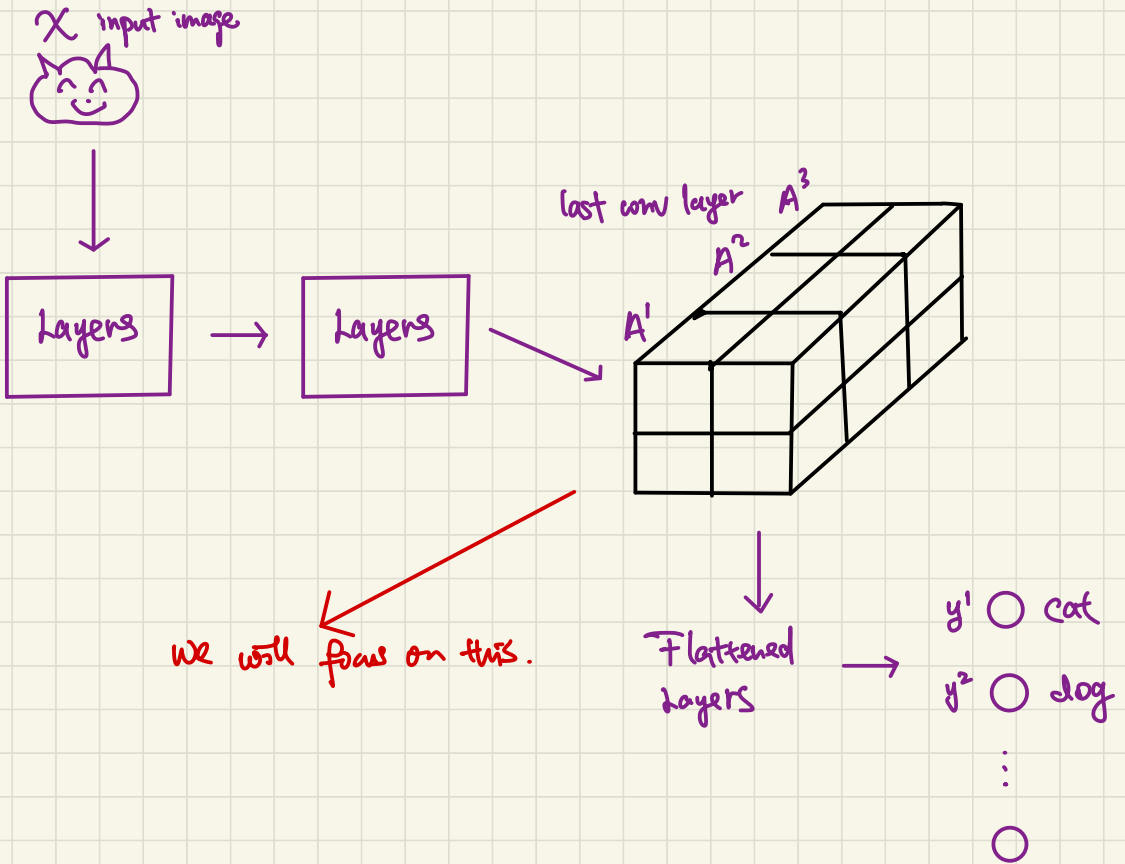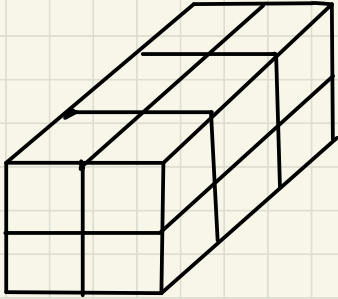Flattened layers $\rightarrow$

$y^1 \bigcirc$ cat

$y^2 \bigcirc$ dog

$\vdots$

$\bigcirc$

Filter 1
Feature Map 1
$A_1$

| 1 | -1 |
|---|---|
| 2 | 0.1 |

Filter 2
Feature Map 2
$A_2$

| 1.5 | 0.1 |
|---|---|
| 2.5 | 0.5 |

Filter 3
Feature Map 3
$A_3$

| -2 | 3 |
|---|---|
| 0 | 2 |

$A^1, A^2, A^3$ are feature maps of a CAT.
By design, values in $A^3$ are "more different"
than those in $A^1$ & $A^2$. This is to highlight
that $A^3$ is the "background kernel" and it
is quite irrelevant to whether a cat is a cat.
We will certainly hope our Grad-CAM does not focus on this.
Our $\alpha_3^c$ should reflect this later.

If input image
is an elephant
on grass!

For eg. this feature map $A^3$
is designed to focus on
the grasses of the
elephant image & thus
is not "important".

Differentiate $y^c$ wrt $A_{ij}^k$

$$\left\{ \frac{dy^c}{dA_1}, \frac{dy^c}{dA_2}, \frac{dy^c}{dA_3} \right\}$$

| 2 | 3 |
|---|---|
| 4 | 2 |

$\frac{dy^c}{dA_1}$

| 3 | 5 |
|---|---|
| 6 | 3 |

$\frac{dy^c}{dA_2}$

| 0.1 | -0.1 |
|-----|------|
| 0.2 | 0.2 |

$\frac{dy^c}{dA_3}$

We perform GAP to get hold of the rate of change of individual feature map $A^k$ wrt $y^c$.

Computing $\frac{dy^c}{dA^k}$ helps us understand how feature map affects the class of interest. In other words, we see that $\frac{dy^c}{dA'} = \begin{bmatrix} 2 & 3 \\ 4 & 2 \end{bmatrix}$ has these values. Then we can intuitively understand the value 2 means a unit change in pixel $A'' = 1$ will cost 2 unit change to $y^c$.

GAP

$\xrightarrow{\text{GAP}}$ $\left\{ \frac{2+3+4+2}{4} = 2.75 \right\}$

$\alpha_1^c$
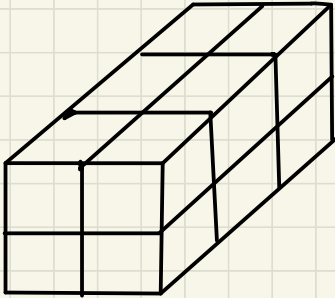
$\alpha_2^c$

$\xrightarrow{\text{GAP}}$ $\left\{ \frac{3+5+6+3}{4} = 4.25 \right\}$

$\alpha_3^c$

$\xrightarrow{\text{GAP}}$ $\left\{ \frac{0.1-0.1+0.2+0.2}{4} = 0.1 \right\}$

$$\left\{ \alpha_1^c \quad \alpha_2^c \quad \alpha_3^c \right\}$$

Feature Maps $\{A_1, A_2, A_3\}$

| 1 | -1 |
|---|----|
| 2 | 0.1 |

Filter 1
Feature Map 1
$A_1$

| 1.5 | 0.1 |
|-----|-----|
| 2.5 | 0.5 |

Filter 2
Feature Map 2
$A_2$

| -2 | 3 |
|----|---|
| 0  | 2 |

Filter 3
Feature Map 3
$A_3$

Differentiate $y^c$ wrt $A_{ij}^k$

$$\left\{ \frac{dy^c}{dA_1}, \quad \frac{dy^c}{dA_2}, \quad \frac{dy^c}{dA_3} \right\}$$

| 2 | 3 |
|---|---|
| 4 | 2 |

$\frac{dy^c}{dA_1}$

| 3 | 5 |
|---|---|
| 6 | 3 |

$\frac{dy^c}{dA_2}$

| 0.1 | -0.1 |
|-----|------|
| 0.2 | 0.2  |

$\frac{dy^c}{dA_3}$

GAP

GAP $\longrightarrow$ $\alpha_1^c$ $\left\{ \frac{2+3+4+2}{4} = 2.75 \right\}$

$\alpha_2^c$

GAP $\longrightarrow$ $\left\{ \frac{3+5+6+3}{4} = 4.25 \right\}$

$\alpha_3^c$

GAP $\longrightarrow$ $\left\{ \frac{0.1 - 0.1 + 0.2 + 0.2}{4} = 0.1 \right\}$

$\left\{ \alpha_1^c \quad \alpha_2^c \quad \alpha_3^c \right\}$

weighted localization MAP.

| 1 | -1 |
|---|---|
| 2 | 0.1 |

$x \quad \alpha_1^c$
$2.75$

| 2.75 | -2.75 |
|---|---|
| 5.5 | 0.275 |

Weighted sum.

| 1.5 | 0.1 |
|---|---|
| 2.5 | 0.5 |

$x \quad \alpha_2^c$
$4.25$

| 6.375 | 0.425 |
|---|---|
| 10.625 | 2.125 |

| 8.925 | -2.025 |
|---|---|
| 16.125 | 2.6 |

| -2 | 3 |
|---|---|
| 0 | 2 |

$x \quad \alpha_3^c \qquad =$
$0.1$

| -0.2 | 0.3 |
|---|---|
| 0 | 0.2 |

Weighted Sum of All Feature Maps

$$L = \alpha_1^c A_1 + \alpha_2^c A_2 + \alpha_3^c A_3 = \begin{bmatrix} 8.925 & -2.025 \\ 16.125 & 2.6 \end{bmatrix}$$

Notice that $\alpha_3^c A_3$ has very small values and hence contribute lesser to how our CNN looks at $y^c$.

Lastly apply ReLU to $L$: $\text{ReLU}(L) = \begin{bmatrix} 8.925 & 0 \\ 16.125 & 2.6 \end{bmatrix}$

L-grad-cam

This has an intuitive meaning, negative values may indicate regions not related to $y^c$.

We overlay L grad cam to the original Image.