

Convolutional Neural Networks

Nicoletta Noceti

Nicoletta.noceti@unige.it

Convolutional Neural Networks

A specialized kind of neural network for processing data with a known grid-like topology

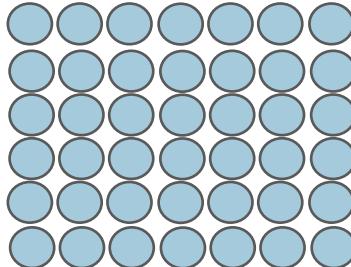
Examples:

- Time-series



1D grid

- Images



2D grid

Motivations

CNNs leverage three important principles:

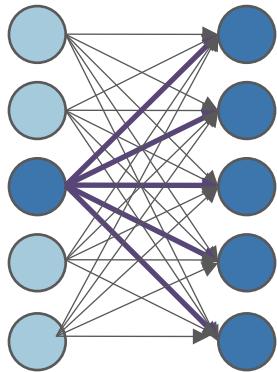
- **Sparse interaction**
- **Parameter sharing**
- **Equivariant representations** ← The layers show equivariance to translation. We will go back to this point later...

Sparse interactions

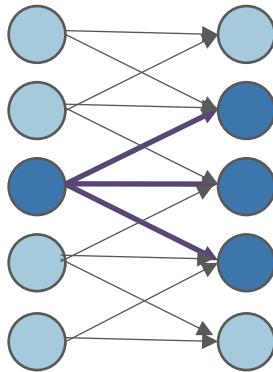
- In traditional DNNs every output unit interacts with every input unit
- In CNNs the “field of view” of each neuron is “limited” and the weights are re-used in different positions of the input
- The pros is that we have to learn fewer parameters, with improvements in
 - Memory requirements
 - Statistical efficiency: statistical strength for more samples per weight, reduced variance when estimating the parameters
 - Computations (less parameters/flops)

Sparse interaction: intuitions

Focus on the input unit



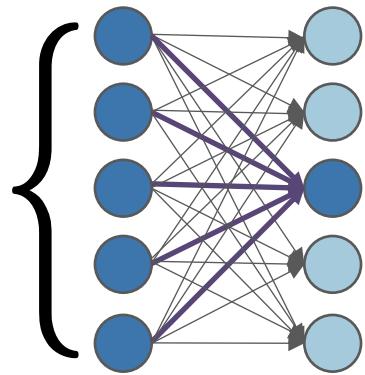
DNN



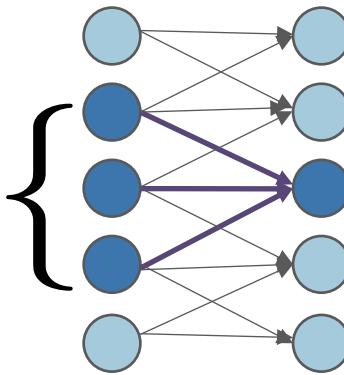
CNN

Sparse interaction: intuitions

Focus on the output unit



DNN

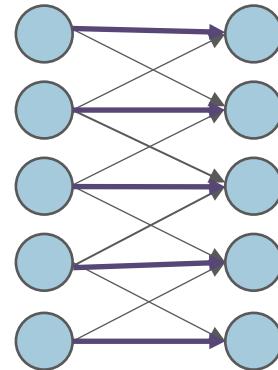
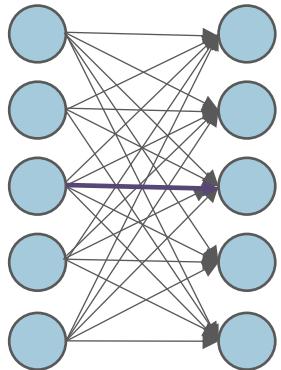


CNN

Receptive fields

Parameter sharing

- The same parameter is used by more than one function in the model
 - In traditional DNNs each weight is used exactly ones when computing the output
 - In CNNs each member of a kernel (a weight) is used at every position of the input



A basic operation

Convolution/Cross Correlation

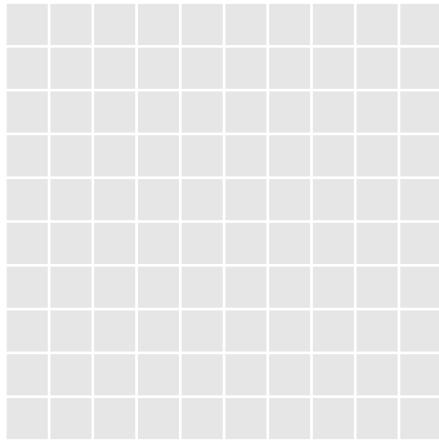
$$\begin{aligned} I_{out}^{conv}(i, j) = (K * I)(i, j) &= \sum_{m=-\infty}^{+\infty} \sum_{n=-\infty}^{+\infty} I(m, n) K(i - m, j - n) = \\ &= \sum_{m=-\infty}^{+\infty} \sum_{n=-\infty}^{+\infty} I(i - m, j - n) K(m, n) \end{aligned}$$

$$I_{out}(i, j) = \sum_{m=-\infty}^{+\infty} \sum_{n=-\infty}^{+\infty} I(i + m, j + n) K(m, n)$$

A basic operation

Convolution/Cross Correlation

$$I_{out}(i, j) = \sum_{m=-W}^{+W} \sum_{n=-W}^{+W} I(i + m, j + n)K(m, n)$$



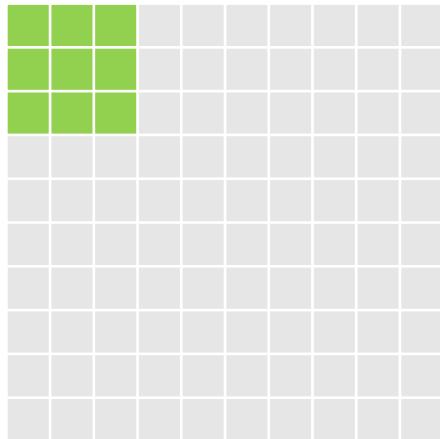
A "feature detector" (kernel) slides over the inputs to generate a feature map

An image of size 10x10

A basic operation

Convolution/Cross Correlation

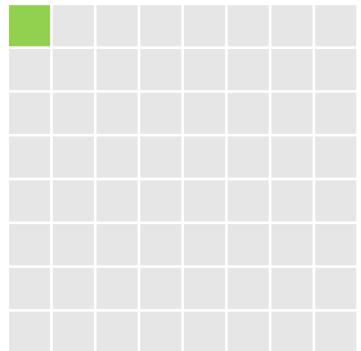
$$I_{out}(i, j) = \sum_{m=-W}^{+W} \sum_{n=-W}^{+W} I(i + m, j + n)K(m, n)$$



INPUT: an image of size 10x10

A kernel (filter) of size 3x3

A "feature detector" (kernel) slides over the inputs to generate a feature map

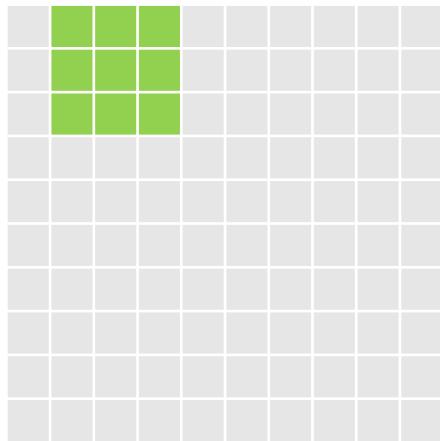


OUTPUT: a structure of size 8x8

A basic operation

Convolution/Cross Correlation

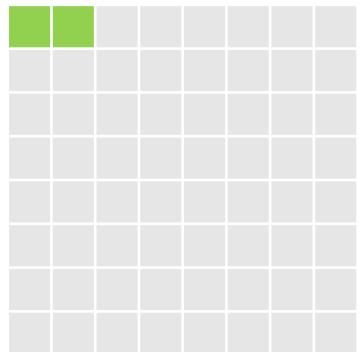
$$I_{out}(i, j) = \sum_{m=-W}^{+W} \sum_{n=-W}^{+W} I(i + m, j + n)K(m, n)$$



INPUT: an image of size 10x10

A kernel (filter) of size 3x3

A "feature detector" (kernel) slides over the inputs to generate a feature map

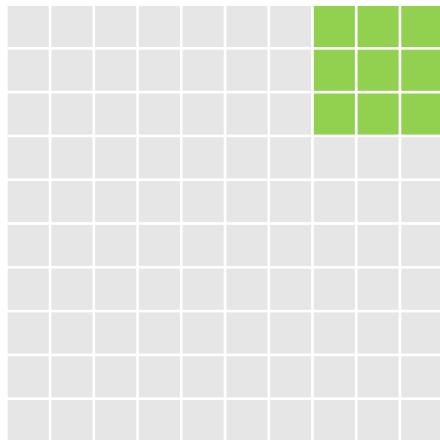


OUTPUT: a structure of size 8x8

A basic operation

Convolution/Cross Correlation

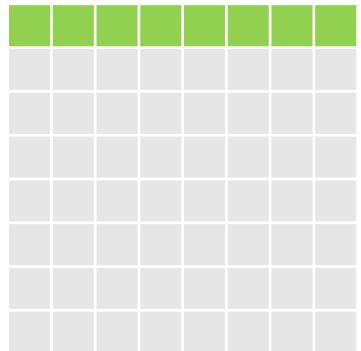
$$I_{out}(i, j) = \sum_{m=-W}^{+W} \sum_{n=-W}^{+W} I(i + m, j + n)K(m, n)$$



INPUT: an image of size 10x10

A kernel (filter) of size 3x3

A "feature detector" (kernel) slides over the inputs to generate a feature map

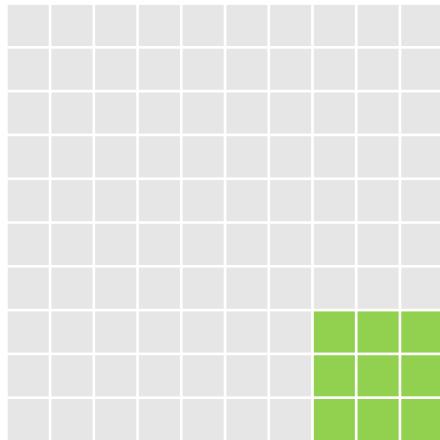


OUTPUT: a structure of size 8x8

A basic operation

Convolution/Cross Correlation

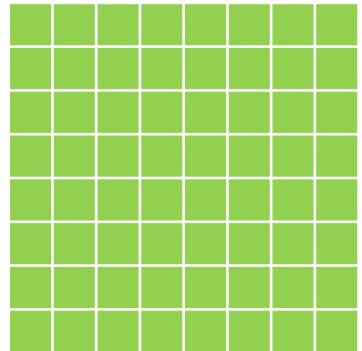
$$I_{out}(i, j) = \sum_{m=-W}^{+W} \sum_{n=-W}^{+W} I(i + m, j + n)K(m, n)$$



INPUT: an image of size 10x10

A kernel (filter) of size 3x3

A "feature detector" (kernel) slides over the inputs to generate a feature map



OUTPUT: a structure of size 8x8

A couple of observations

Dense Neural Networks don't scale to images. Example: Let us consider a fully connected network with a single unit

- With tiny color images of size $32 \times 32 \times 3$ in input:
 - Size of the input layer: $32 \times 32 \times 3 = 3072$
 - Size of the weights: 3072
- With small color images of size $200 \times 200 \times 3$ in input:
 - Size of input layer and weights: $200 \times 200 \times 3 = 120000$

Dense Neural Networks don't fully exploit/capture the specificities of images

CNNs

A change of paradigm

From Feature engineering...



x

→ *Data
Representation
(feature extraction)* →

$\Psi(x)$

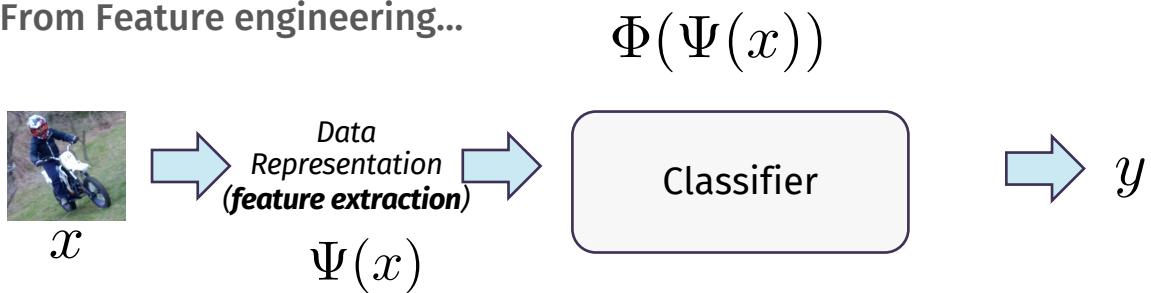
$\Phi(\Psi(x))$



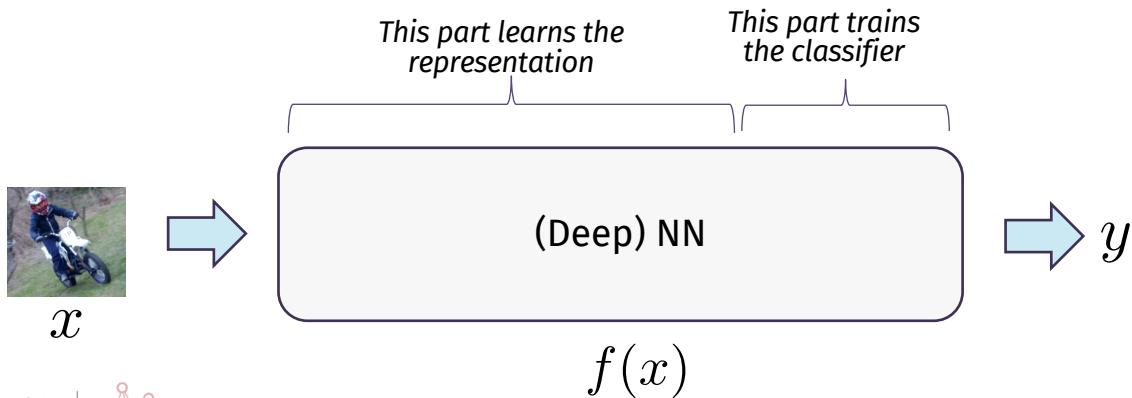
→ y

A change of paradigm

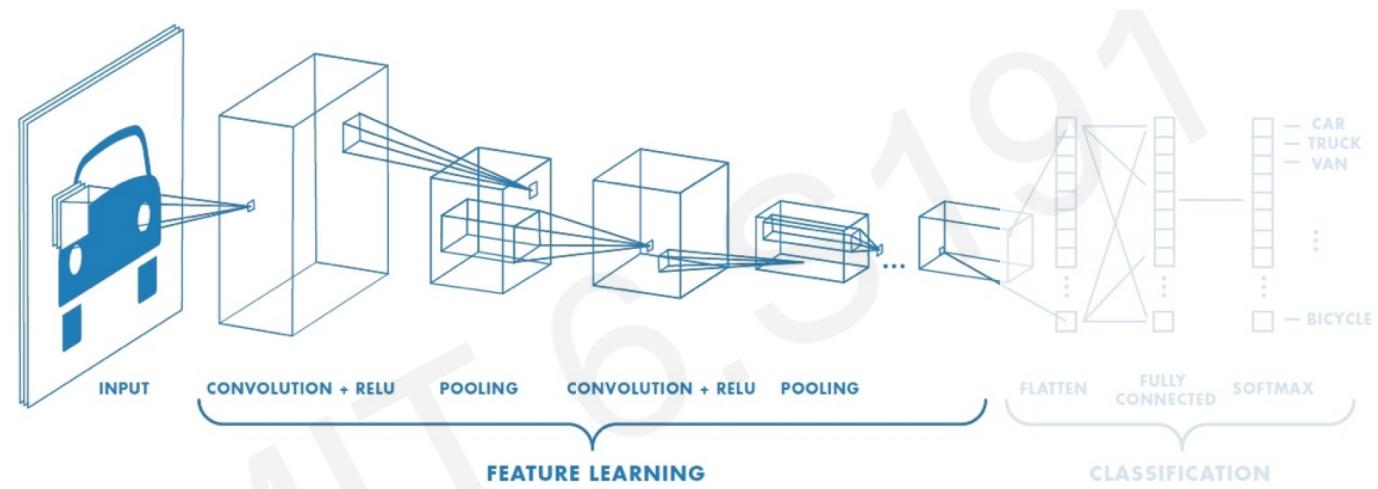
From Feature engineering...



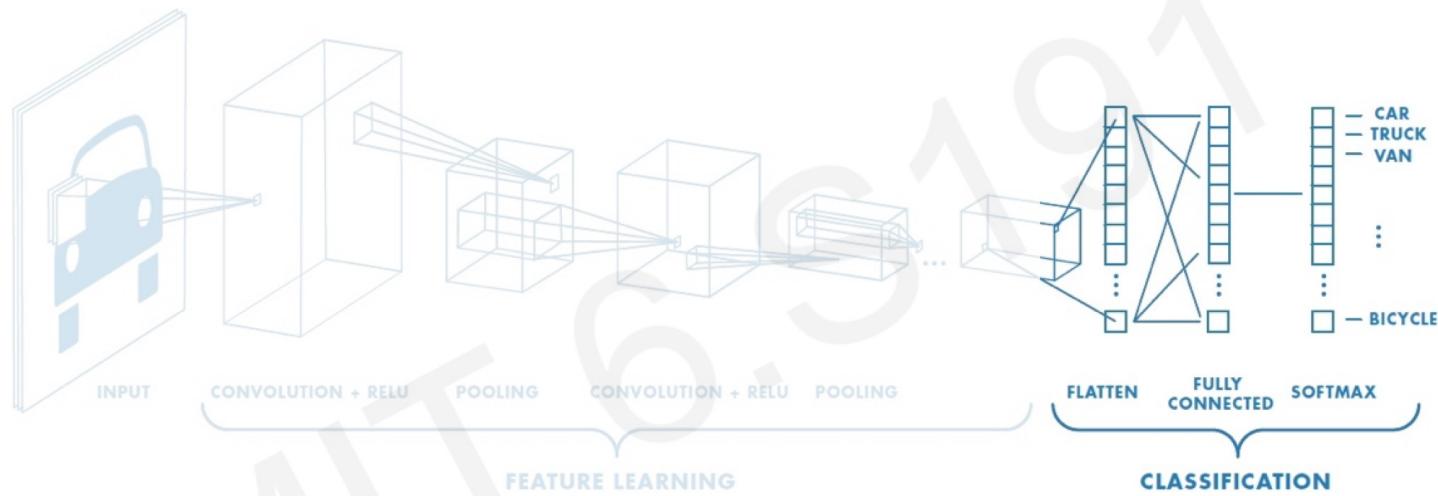
... to Feature learning



A typical CNN



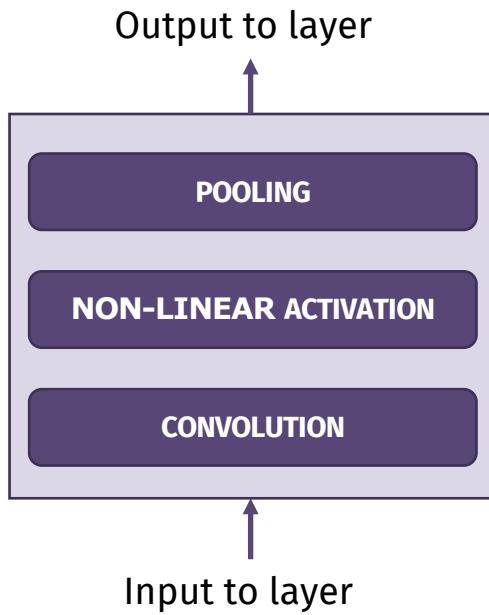
A typical CNN



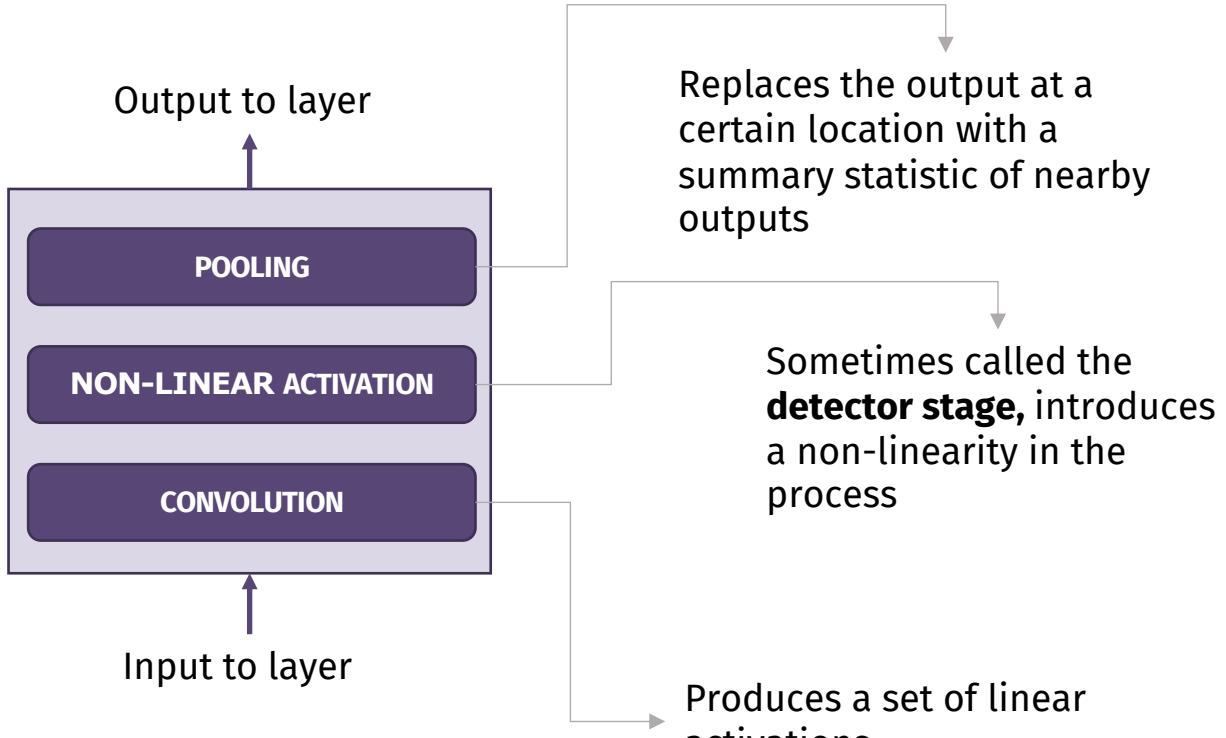
Interlude:
$$\text{softmax}(y_i) = \frac{e^{y_i}}{\sum_{j=1}^{NC} e^{y_j}}$$

[NC is the number of classes]

A typical CNN layer



A typical CNN layer

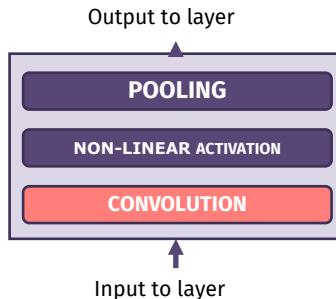


The conv layer

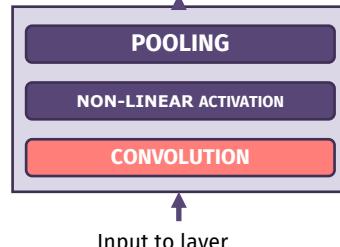
A closer look to convolution and parameter sharing

As the kernel slides on the image, **it is able to capture the same property in different image regions**

Multiple kernels/filters are used → Multiple feature detectors can be used to capture different image properties



Output to layer

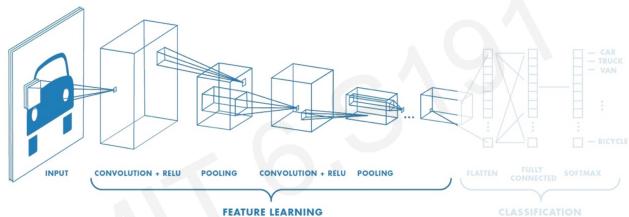


The conv layer

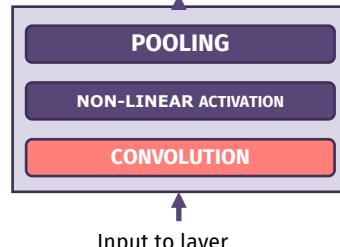
A closer look to convolution and parameter sharing

As the kernel slides on the image, **it is able to capture the same property in different image regions**

Multiple kernels/filters are used → Multiple feature detectors can be used to capture different image properties



Output to layer

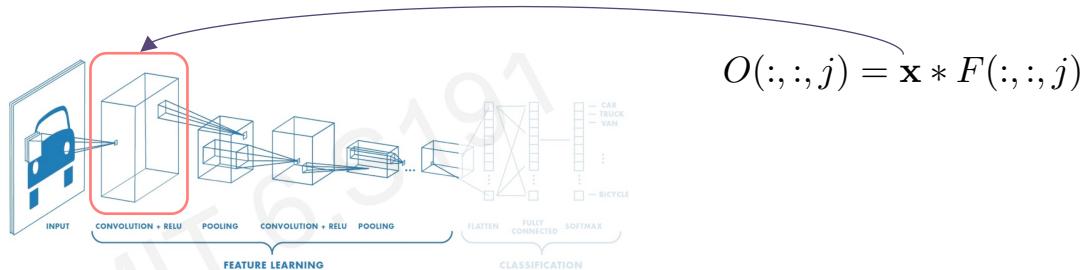


The conv layer

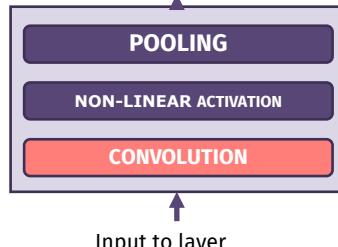
A closer look to convolution and parameter sharing

As the kernel slides on the image, **it is able to capture the same property in different image regions**

Multiple kernels/filters are used → Multiple feature detectors can be used to capture different image properties



Output to layer

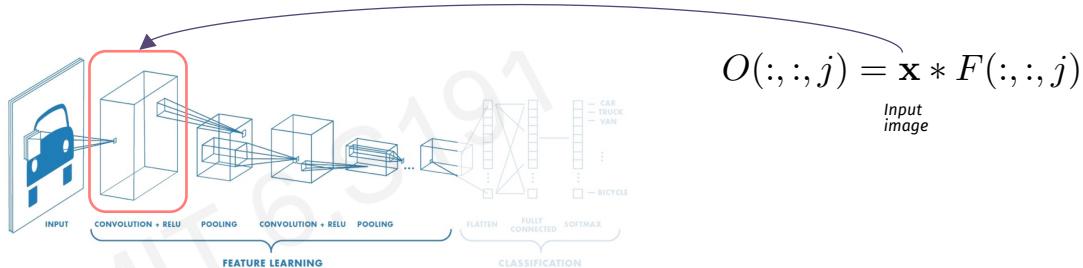


The conv layer

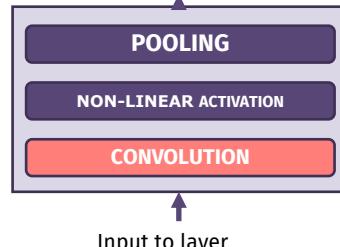
A closer look to convolution and parameter sharing

As the kernel slides on the image, **it is able to capture the same property in different image regions**

Multiple kernels/filters are used → Multiple feature detectors can be used to capture different image properties



Output to layer

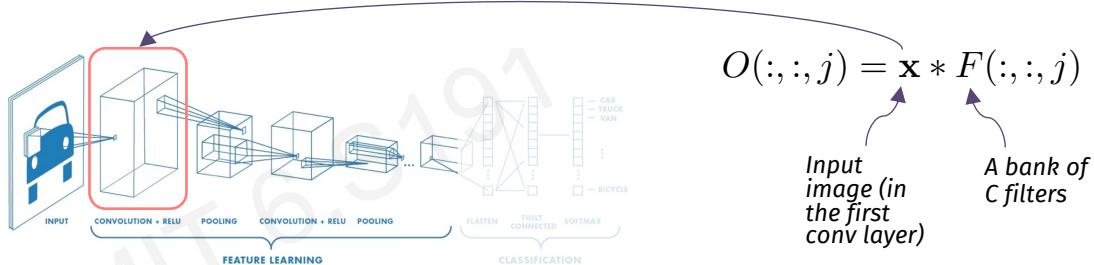


The conv layer

A closer look to convolution and parameter sharing

As the kernel slides on the image, **it is able to capture the same property in different image regions**

Multiple kernels/filters are used → Multiple feature detectors can be used to capture different image properties



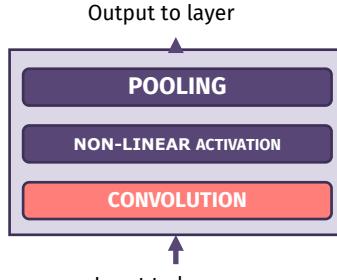
An example from the animation

$x[:, :, 0]$	$w0[:, :, 0]$
0 0 0 0 0 0 0 0	1 1 1 0
0 1 0 2 0 2 0 0	-1 -1 -1
0 1 2 1 2 1 0 0	1 -1 0
0 1 1 1 1 0 0 0	w0[:, :, 1]
0 1 0 0 0 1 0 0	0 -1 0
0 2 1 1 0 0 0 0	1 1 -1
0 0 0 0 0 0 0 0	0 0 0
$x[:, :, 1]$	$w0[:, :, 2]$
0 0 0 0 0 0 0 0	1 1 1
0 1 0 2 2 1 0 0	-1 -1 -1
0 0 2 2 2 2 0 0	0 1 0
0 0 2 2 0 2 0 0	2 0 0
0 0 0 2 2 0 0 0	0 0 1
0 0 2 1 1 1 0 0	1
0 0 0 0 0 0 0 0	
$x[:, :, 2]$	
0 0 0 0 0 0 0 0	
0 1 2 2 0 1 0 0	
0 2 1 1 0 0 0 0	
0 1 0 2 1 1 0 0	
0 2 0 0 2 0 0 0	
0 0 2 0 1 1 0 0	
0 0 0 0 0 0 0 0	

$w1[:, :, 0]$	$o[:, :, 0]$
0 1 1	green 1
1 -1 0	green -2
0 -1 -1	green 2
w1[:, :, 1]	o[:, :, 1]
0 -1 -1	green 1
-1 0 0	green -4
1 -1 1	green 3
3 4 3	green 3
3 -1 4	green -4
$w1[:, :, 2]$	$o[:, :, 2]$
1 1 1	green 1
1 1 1	green -1
0 0 1	green 1
Bias $b0 (1 \times 1 \times 1)$	
0	0
$b1[:, :, 0]$	
0	0

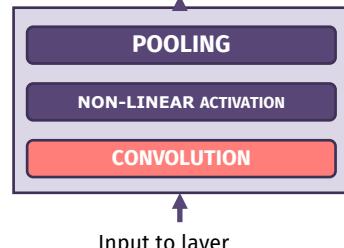
$$\begin{aligned}
 & [2*1 + 1*1 + 1*(-1) + 1*(-1) + 1*(-1)] + \\
 & [2*(-1) + 2*1 + 2*1] + \\
 & [1*(-1) + 1*(-1) + 2*(-1) + 1*(-1)] + \\
 & 1 = \\
 & = 2 + 1 - 1 - 1 - 2 + 2 + 2 - 1 - 1 - 2 - 1 + 1 \\
 & = 2 - 1 - 1 - 1 - 1 = -2
 \end{aligned}$$

toggle movement



See the demo here: <https://cs231n.github.io/convolutional-networks>

Output to layer

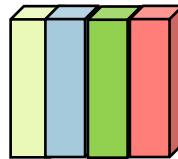


Another animation to clearly understand

$W \times H \times 3$



$K \times K \times 3 \times 4$
($C==4$)



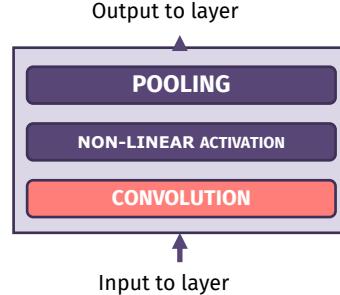
$W^l \times H^l \times 4$
($C==4$)



*C is the number of channels,
a hyper-parameter of the
conv layer*

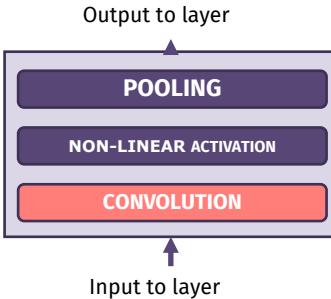
Output feature size of conv layers

Three parameters control the size of the output of a layer



- **Depth**, the number of filters (kernels) of the layer
- **Stride**, the step used to slide the filter on the input
 - When stride > 1 we are down-sampling the input data
 - **Tiling** refers to the special case where stride = kernel span
- **Padding** to enlarge the input and allow for kernel application in each one of the (original) point

Output features size of conv layers



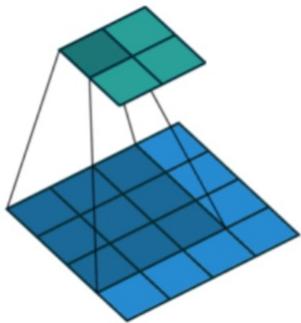
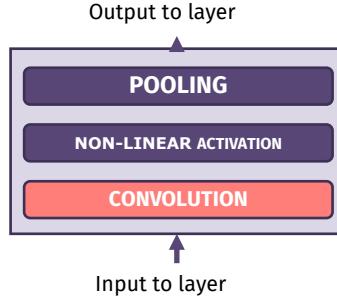
$$O = \frac{W - K + 2P}{S} + 1$$

Annotations for the formula:

- "Size BEFORE convolution" points to the term W .
- "Kernel size" points to the term K .
- "Padding" points to the term $2P$.
- "Stride" points to the term S .
- "Size AFTER convolution" points to the result O .

Output features size of conv layers

Examples



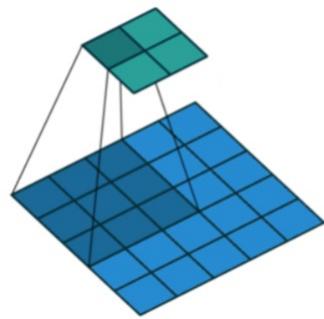
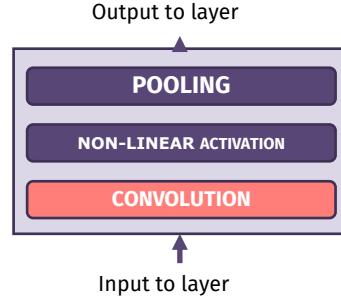
No padding, stride 1

$$O = \frac{W - K + 2P}{S} + 1$$

where O is the output size, W is the input width, K is the kernel size, P is the padding, and S is the stride.

Output features size of conv layers

Examples

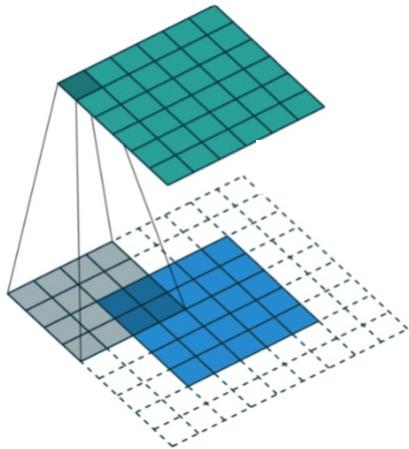
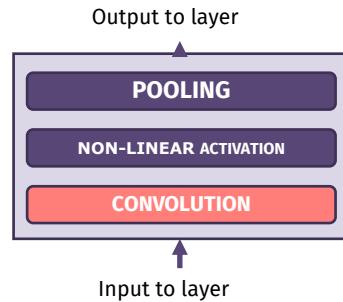


$$O = \frac{\frac{5}{2}W - \frac{3}{2}K + 2P}{S_2} + 1$$

No padding, stride 2

Output features size of conv layers

Examples

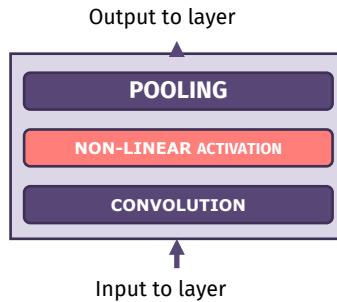


Padding 2, stride 1

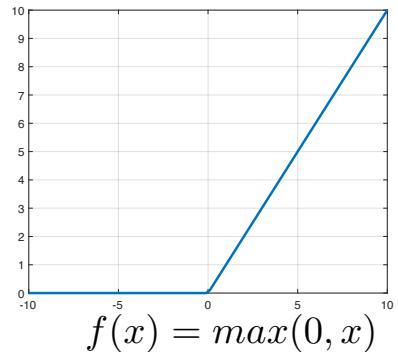
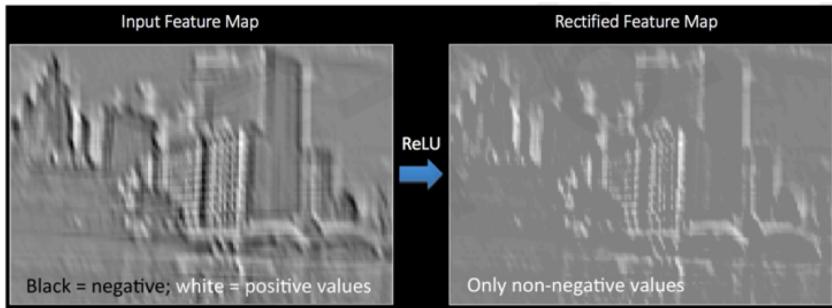
$$O = \frac{W - K + 2P}{S} + 1$$

The equation for calculating the output size of a convolutional layer. The variables are defined as follows:
- O : Output size
- W : Input width
- K : Kernel size
- P : Padding size
- S : Stride size
In the equation, the terms $W - K + 2P$ are highlighted in blue, and the term S is highlighted in red.

Introducing non-linearities

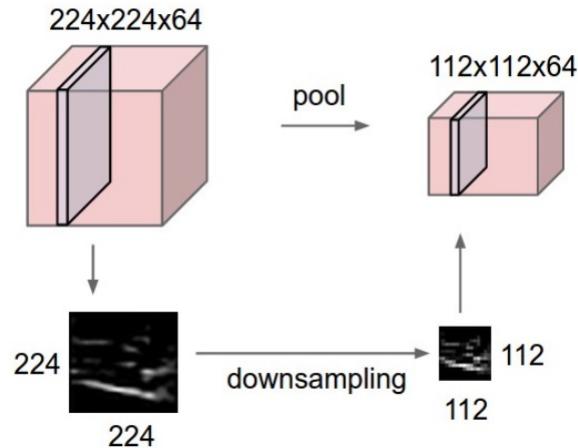
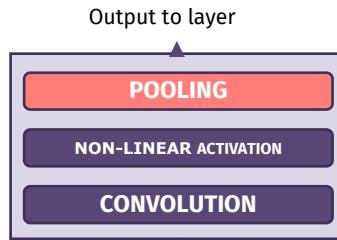


A typical choice is ReLU: it is applied after each convolutional layer and only preserves non-negative values



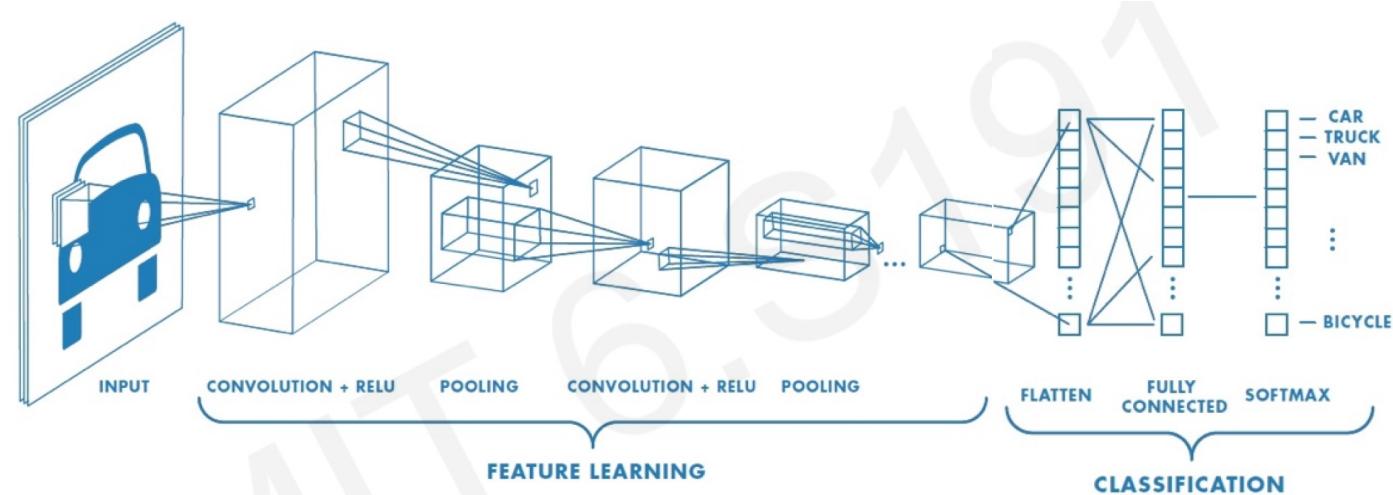
Pooling and invariance

- It provides invariance to small shifts of the inputs
- Pooling functions:
 - **Average pooling:** average activation of the convolutional layer
 - **Max pooling:** max activation of the convolutional layer



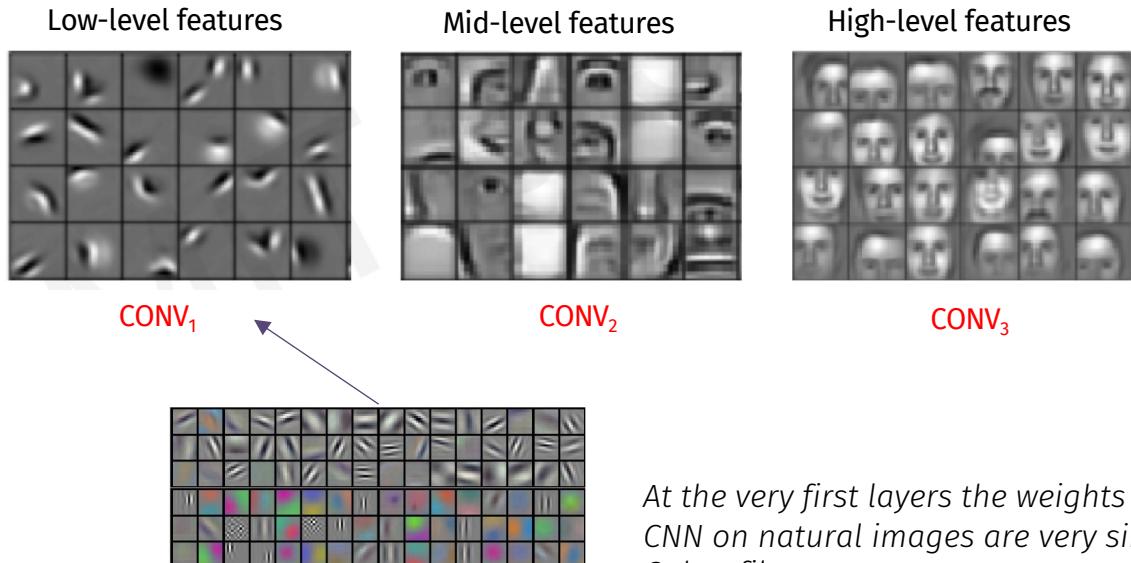
No parameter to be estimated here!

A sequence of Conv layers



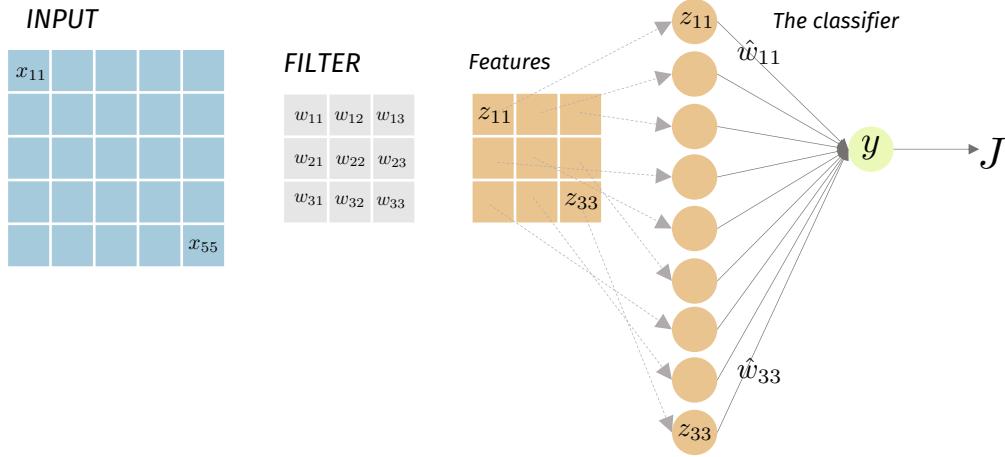
Towards image classification

So... Is it possible to learn the most appropriate representation, possibly with a hierarchy, directly from the data?

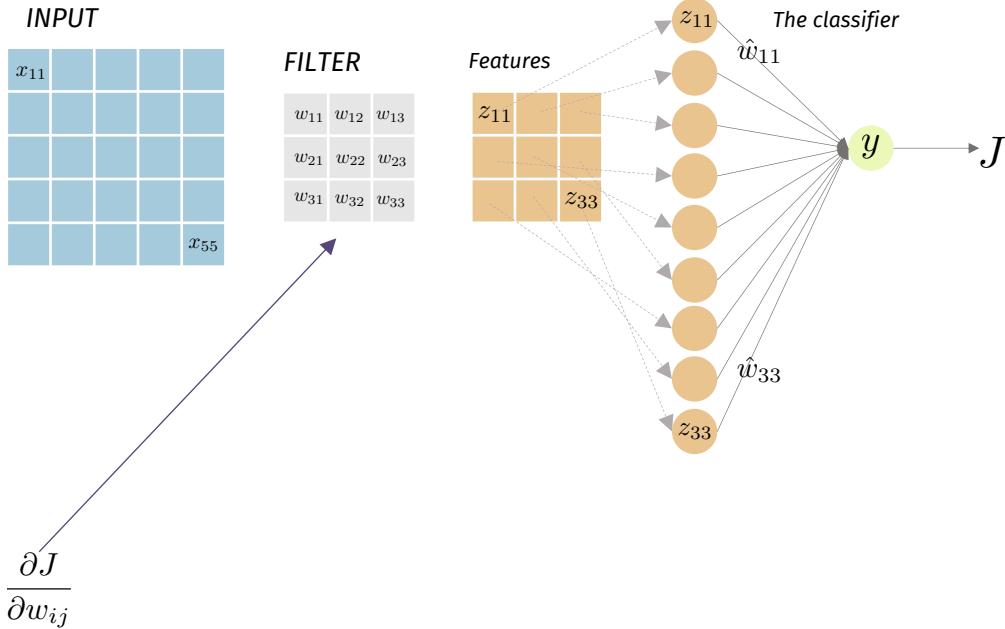


At the very first layers the weights learnt by a CNN on natural images are very similar to Gabor filters

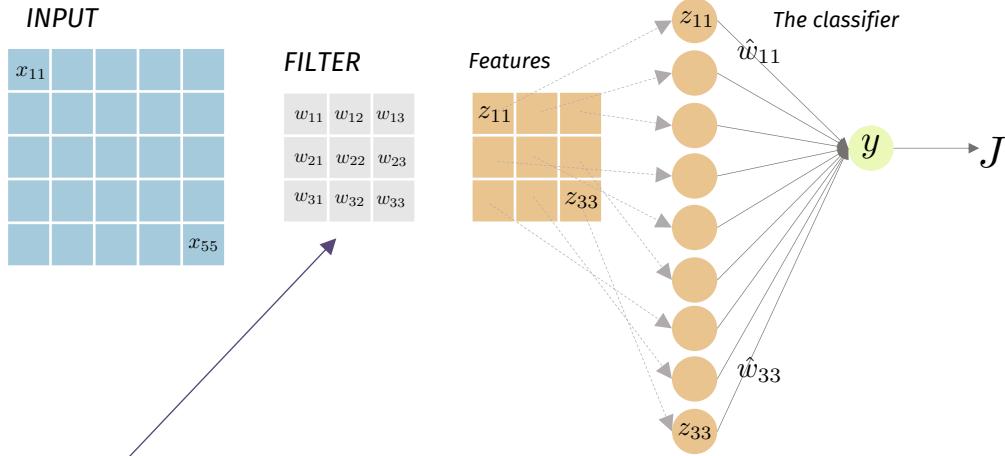
Backpropagation in CNNs (intuition)



Backpropagation in CNNs (intuition)

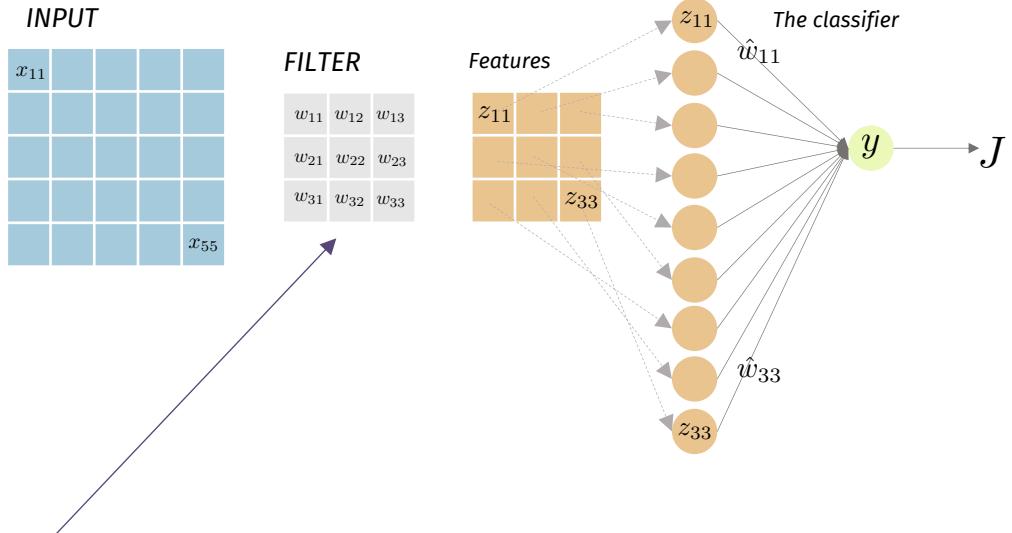


Backpropagation in CNNs (intuition)



$$\frac{\partial J}{\partial w_{ij}} = \frac{\partial J}{\partial f_\sigma(\hat{\mathbf{w}}^\top \mathbf{z})} \frac{\partial f_\sigma(\hat{\mathbf{w}}^\top \mathbf{z})}{\partial \hat{\mathbf{w}}^\top \mathbf{z}}$$

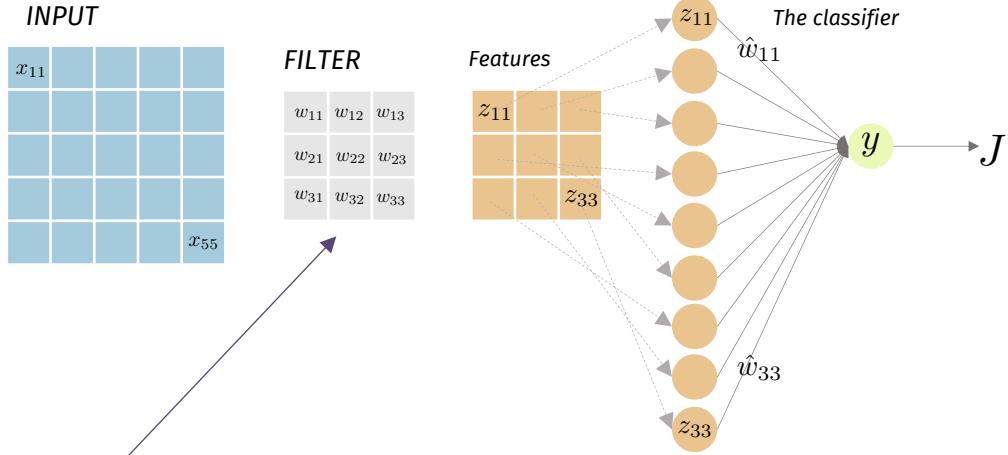
Backpropagation in CNNs (intuition)



$$\frac{\partial J}{\partial w_{ij}} = \frac{\partial J}{\partial f_\sigma(\hat{\mathbf{w}}^\top \mathbf{z})} \frac{\partial f_\sigma(\hat{\mathbf{w}}^\top \mathbf{z})}{\partial \hat{\mathbf{w}}^\top \mathbf{z}} \sum_{k=1}^3 \sum_{h=1}^3 \frac{\partial \hat{\mathbf{w}}^\top \mathbf{z}}{\partial z_{kh}}$$

$$z_{kh} = \sum_{m=-1}^{+1} \sum_{n=-1}^{+1} I(k+m, h+n) K(m, n)$$

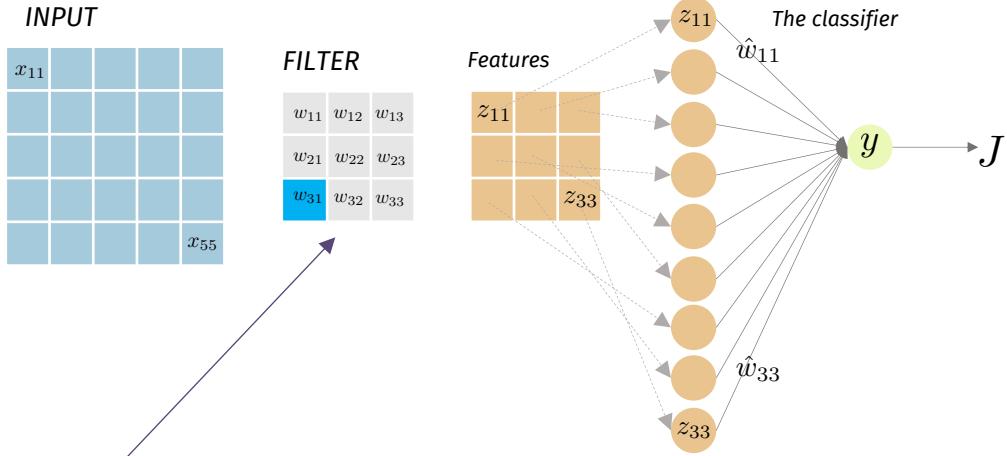
Backpropagation in CNNs (intuition)



$$\frac{\partial J}{\partial w_{ij}} = \frac{\partial J}{\partial f_\sigma(\hat{\mathbf{w}}^\top \mathbf{z})} \frac{\partial f_\sigma(\hat{\mathbf{w}}^\top \mathbf{z})}{\partial \hat{\mathbf{w}}^\top \mathbf{z}} \sum_{k=1}^3 \sum_{h=1}^3 \frac{\partial \hat{\mathbf{w}}^\top \mathbf{z}}{\partial z_{kh}} \frac{\partial z_{kh}}{\partial w_{ij}}$$

$$z_{kh} = \sum_{m=-1}^{+1} \sum_{n=-1}^{+1} I(k+m, h+n) K(m, n)$$

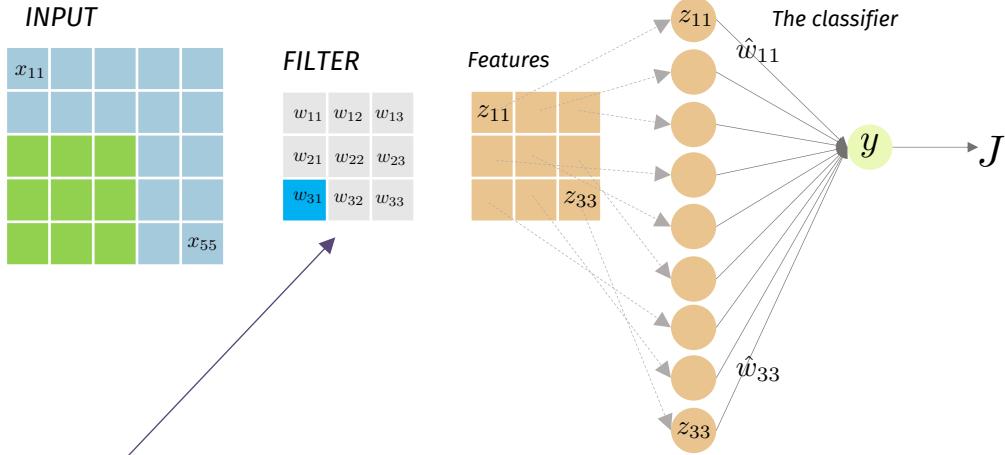
Backpropagation in CNNs (intuition)



$$\frac{\partial J}{\partial w_{31}} = \frac{\partial J}{\partial f_\sigma(\hat{\mathbf{w}}^\top \mathbf{z})} \frac{\partial f_\sigma(\hat{\mathbf{w}}^\top \mathbf{z})}{\partial \hat{\mathbf{w}}^\top \mathbf{z}} \sum_{k=1}^3 \sum_{h=1}^3 \frac{\partial \hat{\mathbf{w}}^\top \mathbf{z}}{\partial z_{kh}} \frac{\partial z_{kh}}{\partial w_{31}}$$

Let's be more specific

Backpropagation in CNNs (intuition)



$$\frac{\partial J}{\partial w_{31}} = \frac{\partial J}{\partial f_\sigma(\hat{\mathbf{w}}^\top \mathbf{z})} \frac{\partial f_\sigma(\hat{\mathbf{w}}^\top \mathbf{z})}{\partial \hat{\mathbf{w}}^\top \mathbf{z}} \sum_{k=1}^3 \sum_{h=1}^3 \frac{\partial \hat{\mathbf{w}}^\top \mathbf{z}}{\partial z_{kh}} \frac{\partial z_{kh}}{\partial w_{31}}$$

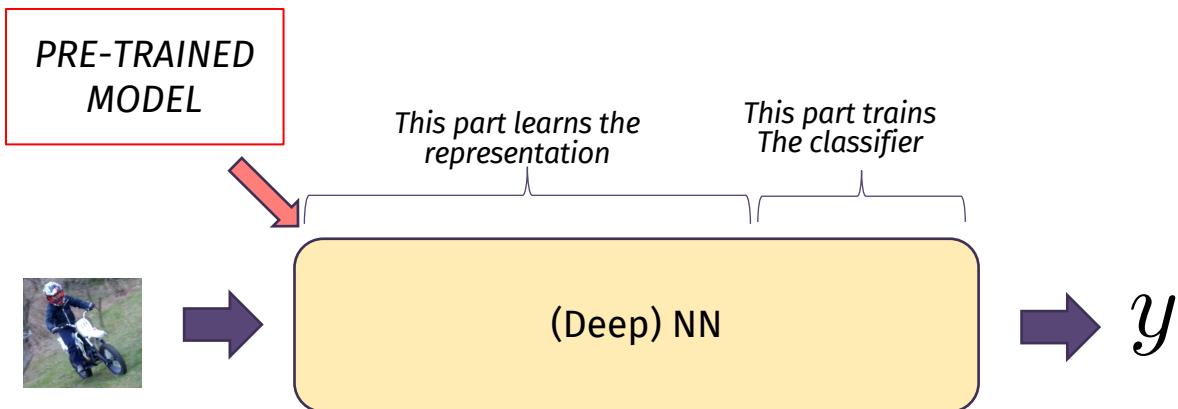
Let's be more specific

... this is still a convolution/correlation

Some further discussions

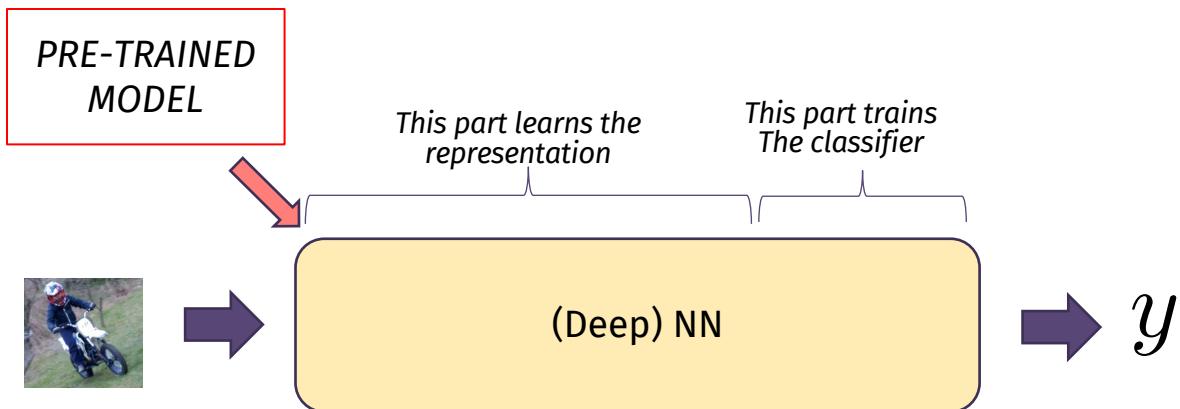
About knowledge reuse

- Efficiency/sustainability are becoming hot keywords...
- Also, in some cases we might not have enough data
- It may be convenient to exploit already existing knowledge!



Using pre-trained models

- Instead of training the network from scratch, one can start from weights learnt on other (possibly very big) datasets, and then fine-tune on the target task



Again on knowledge reuse: the role of data supervision

From unsupervised to self-supervised learning

Sometimes data are available but only few of them have labels

Unlabeled data

Labeled data

Again on knowledge reuse: the role of data supervision

From unsupervised to self-supervised learning

Unlabeled data

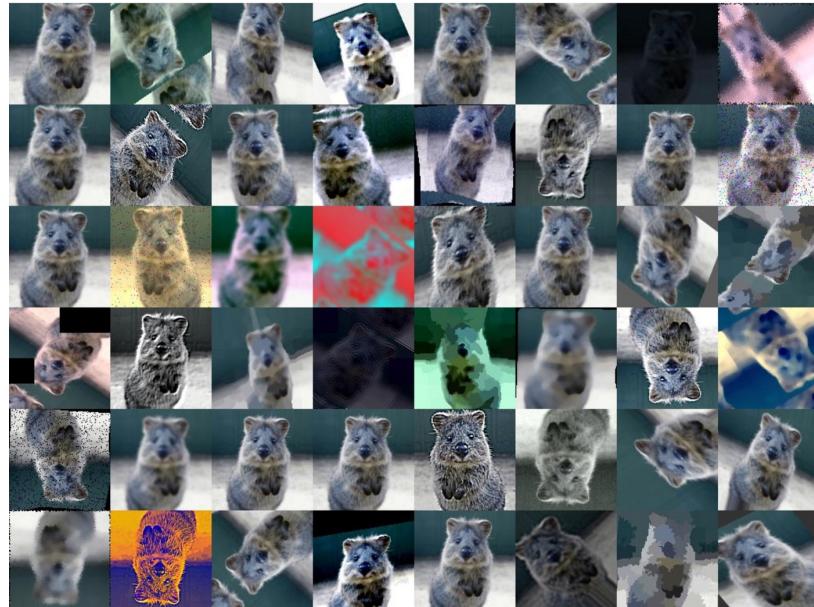
Labeled data

Use the unlabeled data to train a self-supervised problem → The supervision is derived by the data

Fine-tune using the labeled data on the target supervised task

Data Augmentation

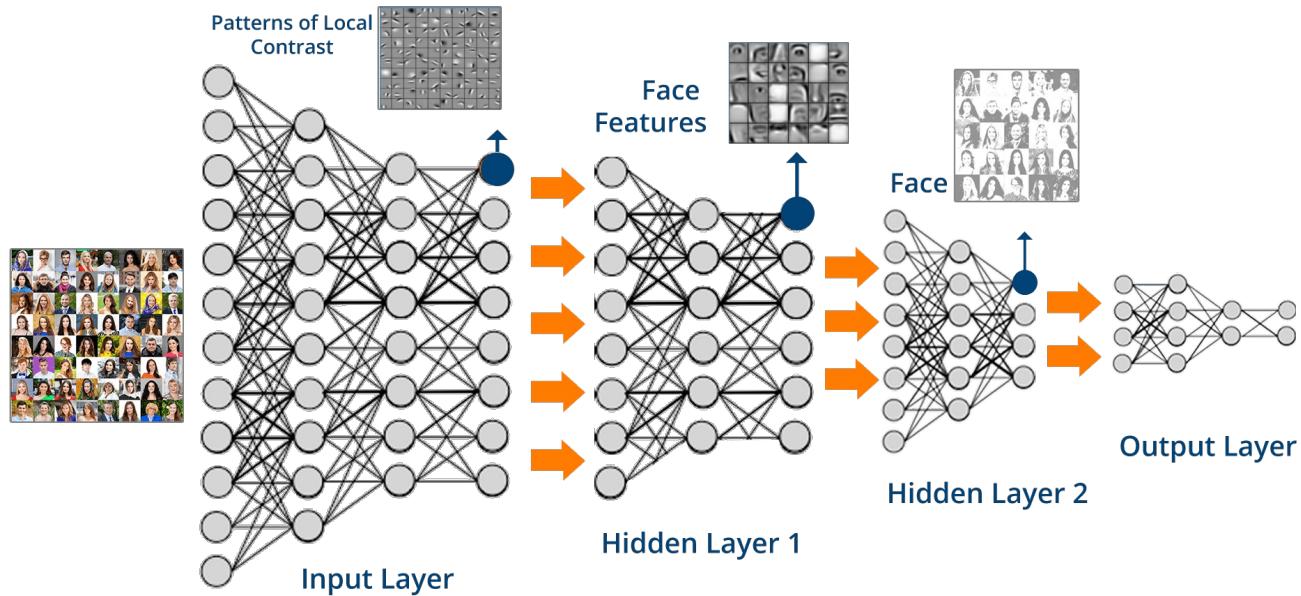
An alternative to enrich the dataset



Is it always possible?

From <https://towardsdatascience.com/image-augmentation-for-deep-learning-using-keras-and-histogram-equalization-9329f6ae5085>

How to interpret a deep architecture?



Interpretable models or interpretable data?

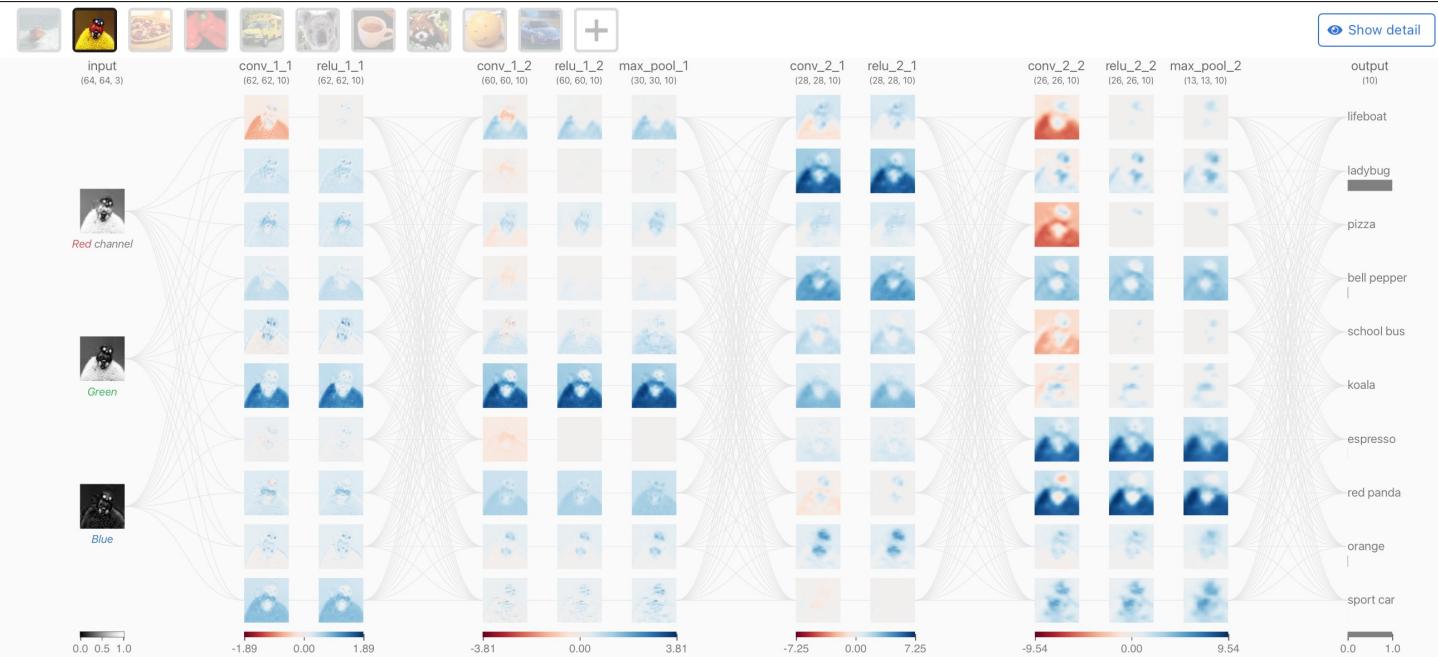


Gradient-weighted Class Activation Mapping (Grad-CAM), uses the gradients of any target concept flowing into the final convolutional layer to produce a coarse localization map highlighting the important regions in the image for predicting the concept

From <https://towardsdatascience.com/understand-your-algorithm-with-grad-cam-d3b62fce353>

A nice visualization

<https://poloclub.github.io/cnn-explainer/>



UniGe

