

Deep Learning - Foundations and Concepts

Chapter 10. Convolutional Networks

nonlineark@github

March 10, 2025

Outline

- 1 Computer Vision
- 2 Convolutional Filters
- 3 Visualizing Trained CNNs

Computer vision

- Computer vision was one of the first fields to be transformed by the deep learning revolution, predominantly thanks to the CNN architecture.
- Recently alternative architectures based on transformers have become competitive with convolutional networks in some applications.
- Some applications for machine learning in computer vision: Classification, detection, segmentation, caption generation, synthesis, inpainting, style transfer, super-resolution, depth prediction, scene reconstruction.

Image data

- The structure of an image:
 - An image comprises a rectangular array of pixels.
 - Each pixel has either a grey-scale intensity or a triplet of red, green and blue channels each with its own intensity value.
- Challenges of applying neural networks to image data:
 - Images generally have a high dimensionality.
 - Image data is highly structured.
- Local correlations can be used to encode strong inductive biases into a neural network, leading to models with far fewer parameters and with much better generalization accuracy.

Inductive biases

To exploit the two-dimensional structure of image data to create inductive biases, we can use four interrelated concepts:

- Hierarchy.
- Locality.
- Equivariance.
- Invariance.

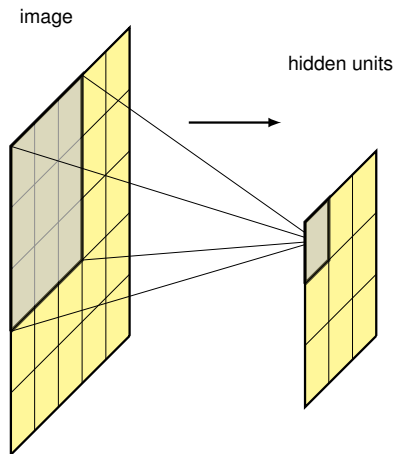
Feature detectors

Let's restrict our attention to grey-scale images first. Consider a single unit in the first layer of a neural network:

- It takes as input just the pixel values from a small rectangular region, or patch, from the image. This patch is referred to as the receptive field of that unit.
- The output of this unit is given by the usual functional form comprising a weighted linear combinations of the input values, which is subsequently transformed using a nonlinear activation function:
$$z = \text{ReLU}(w^T x + w_0).$$
- The weights themselves form a small two-dimensional grid known as a filter, sometimes also called a kernel.
- The unit acts as a feature detector that signals when it finds a sufficiently good match to its kernel (basic properties of inner product).

Feature detectors

Figure: A single unit in the first layer of a neural network



Translation equivariance

- To generalize what our neural network has learned in one location to all possible locations in the image, we can simply replicate the same hidden-unit weight values at multiple locations across the image.
- The units of the hidden layer form a feature map in which all the units share the same weights.

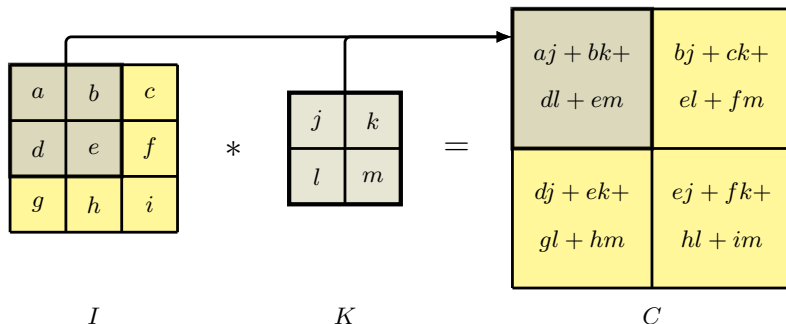
For an image I with pixel intensities $I(j, k)$, and a filter K with pixel values $K(l, m)$, the feature map C has activation values given by:

$$C(j, k) = \sum_l \sum_m I(j + l, k + m) K(l, m)$$

where we have omitted the nonlinear activation function for clarity. This is an example of a convolution and is sometimes expressed as $C = I * K$.

Translation equivariance

Figure: A 3×3 image convolved with a 2×2 filter



Translation equivariance

Comparing this convolutional structure with a standard fully connected network, we see several advantages:

- The connections are sparse, leading to far fewer weights even with large images.
- The weight values are shared, greatly reducing the number of independent parameters and consequently reducing the required size of the training set.
- The same network can be applied to images of different sizes without the need for retraining.
- Convolutions can be implemented very efficiently by exploiting the massive parallelism of GPUs.

Padding

Consider an image of dimensionality $J \times K$ and a kernel of dimensionality $M \times M$:

- The feature map is of dimensionality $(J - M + 1) \times (K - M + 1)$, which is smaller than the original image.
- We can pad the original image with additional pixels around the outside. If we pad with P pixels then the output map has dimensionality $(J + 2P - M + 1) \times (K + 2P - M + 1)$:
 - If $P = 0$, this is called a valid convolution.
 - If $P = \frac{M-1}{2}$, this is called a same convolution.
- Padding values are typically set to zero, after first subtracting the mean from each image so that zero represents the average value of the pixel intensity.

Strided convolutions

Sometimes we wish to use feature maps that are significantly smaller than the original image to provide flexibility in the design of convolutional network architectures. One way to achieve this is to use strided convolutions in which, the filter is moved in larger steps of size S , called the stride. The feature map will have dimensionality:

$$(\lfloor \frac{J + 2P - M}{S} \rfloor + 1) \times (\lfloor \frac{K + 2P - M}{S} \rfloor + 1)$$

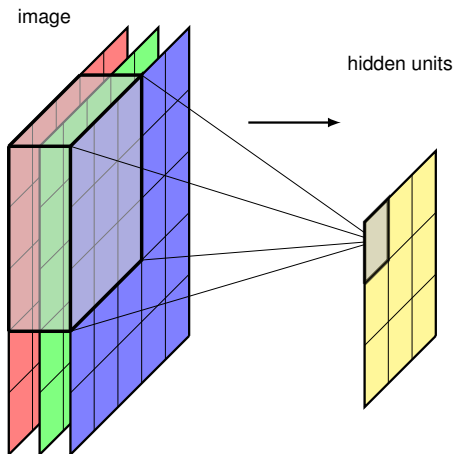
Multi-dimensional convolutions

We can easily extend convolutions to cover multiple channels by extending the dimensionality of the filter:

- An image with $J \times K$ pixels and C channels will be described by a tensor of dimensionality $J \times K \times C$.
- We can introduce a filter described by a tensor of dimensionality $M \times M \times C$.

Multi-dimensional convolutions

Figure: A multi-dimensional filter



Multi-dimensional convolutions

To build more flexible models, we simply include multiple filters, in which each filter has its own independent set of parameters giving rise to its own independent feature map:

- We again refer to these separate feature maps as channels.
- The filter tensor now has dimensionality $M \times M \times C \times C_{\text{OUT}}$ where C is the number of input channels and C_{OUT} is the number of output channels.
- Each output channel will have its own associated bias parameter, so the total number of parameters will be $(M^2C + 1)C_{\text{OUT}}$.

Pooling

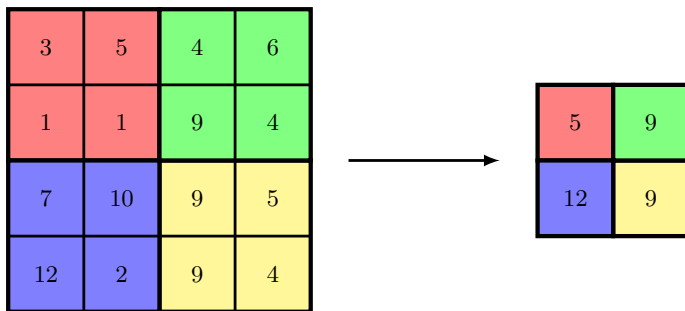
Small changes in the relative locations do not affect the classification, and we want to be invariant to such small translations of individual features.

This can be achieved using pooling:

- A pooling layer is similar to a convolutional layer:
 - Each unit takes input from a receptive field in the previous feature map layer.
 - There is a choice of filter size and of stride length.
- The difference is that the output of a pooling unit is a simple, fixed function of its inputs:
 - Max-pooling.
 - Average pooling.

Pooling

Figure: Illustration of max-pooling

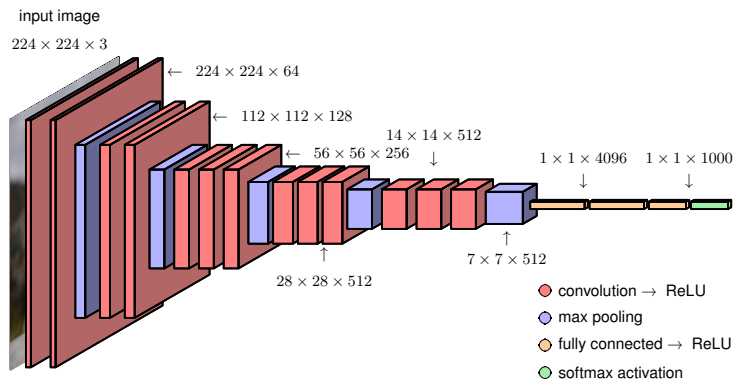


Multilayer convolutions

- To allow the network to discover and represent hierarchical structure in the data, we extend the architecture by considering multiple layers.
- Each convolutional layer is described by a filter tensor of dimensionality $M \times M \times C_{\text{IN}} \times C_{\text{OUT}}$ in which the number of independent weight and bias parameters is $(M^2 C_{\text{IN}} + 1) C_{\text{OUT}}$.
- Each such convolutional layer can optionally be followed by a pooling layer.
- In many applications (e.g., classification tasks), the output units need to combine information from across the whole of the input image. This is typically achieved by introducing one or two standard fully connected layers as the final stages of the network.

Example network architectures

Figure: The architecture of a convolutional network called VGG-16



Visual cortex

Historically, much of the motivation for CNNs came from pioneering research in neuroscience:

- Simple cells: Respond to visual inputs with a simple edge oriented at a particular angle and located at a particular position within the visual field.
- Complex cells: Respond to more complex stimuli and which seem to be derived by combining and processing the output of simple cells.
- Grandmother cells: Respond if, and only if, the visual input corresponds to a person's grandmother, irrespective of location, scale, lighting, or other transformations of the scene.

Visual cortex

The response of the simple cells can be modelled using Gabor filters:

$$G(x, y) = A \exp(-\alpha \tilde{x}^2 - \beta \tilde{y}^2) \sin(\omega \tilde{x} + \phi)$$

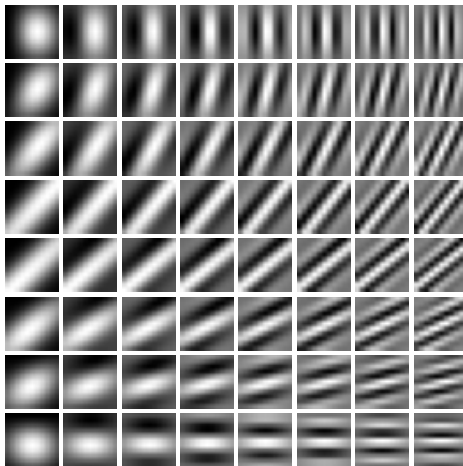
$$\tilde{x} = (x - x_0) \cos \theta + (y - y_0) \sin \theta$$

$$\tilde{y} = -(x - x_0) \sin \theta + (y - y_0) \cos \theta$$

- The \sin term represents a sinusoidal spatial oscillation oriented in a direction defined by the polar angle θ , with frequency ω and phase angle ϕ .
- The exponential factor creates a decay envelope that localizes the filter in the neighborhood of position (x_0, y_0) and with decay rates governed by α and β .

Visual cortex

Figure: Examples of Gabor filters



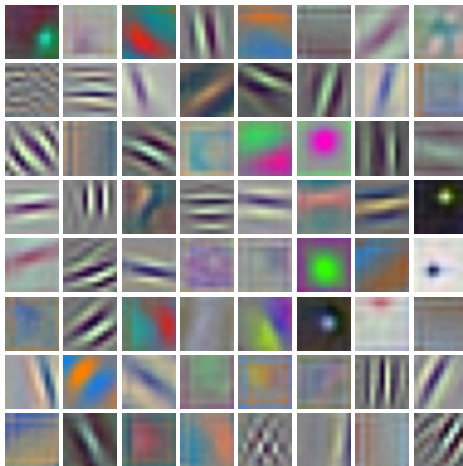
Visualizing trained filters

For the filters in the first convolutional layer, we can visualize the network weights associated with these filters directly as small images (Cauchy-Schwarz inequality).

We see a remarkable similarity between these filters and the Gabor filters. This is because these characteristic filters are a general property of the statistics of natural images and therefore prove useful for image understanding in both natural and artificial systems.

Visualizing trained filters

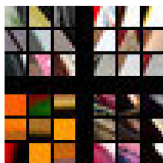
Figure: Examples of learned filters from the first layer of AlexNet



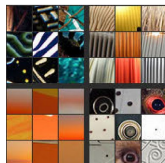
Visualizing trained filters

One approach to visualize the filters in subsequent layers is to present a large number of image patches to the network and see which produce the highest activation value in any particular hidden unit.

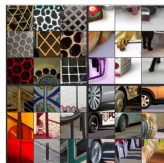
Figure: Examples of image patches that produce the strongest activation in the hidden units in a network



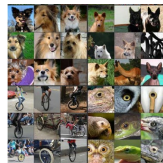
Layer 1



Layer 2



Layer 3

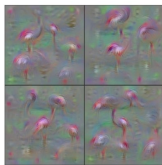


Layer 5

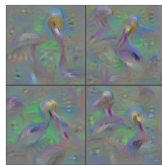
Visualizing trained filters

We can extend this technique to perform a numerical optimization over the input variables to maximize the activation of a particular unit. If we chose the unit to be one of the outputs then we can look for an image that is most representative of the corresponding class label.

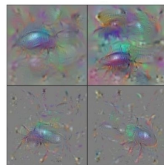
Figure: Examples of synthetic images generated by maximizing the class probability



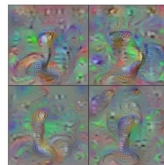
Flamingo



Pelican



Ground Beetle



Indian Cobra

Saliency maps

- Saliency maps aim to identify those regions of an image that are most significant in determining the class label.
- This is best done by investigating the final convolutional layer:
 - This layer still retains spatial localization.
 - This layer has the highest level of semantic representation.

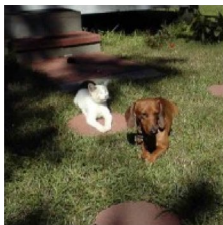
Saliency maps

The Grad-CAM (gradient class activation mapping) method. For a given input image:

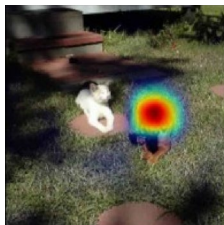
- First compute the derivatives of the output-unit pre-activation a^c for a given class c with respect to the pre-activations $a_{ij}^{(k)}$ of all the units in the final convolutional layer for channel k .
- For each channel in that layer, the average of those derivatives is evaluated to give: $\alpha_k = \frac{1}{M_k} \sum_i \sum_j \frac{\partial a^c}{\partial a_{ij}^{(k)}}$, where M_k is the total number of units in that channel.
- These averages are then used to form a weighted sum of the form: $L = \sum_k \alpha_k A^{(k)}$, where $A^{(k)}$ is a matrix with elements $a_{ij}^{(k)}$.

Saliency maps

Figure: Saliency maps for the VGG-16 network with respect to the dog and cat categories



Original image



Saliency map for 'dog'



Saliency map for 'cat'

Adversarial attacks

- Adversarial attacks involve making very small modifications to an image, at level that is imperceptible to a human, which cause the image to be misclassified by the neural network.
- The fast gradient sign method: Changing each pixel value in an image x by a fixed amount ϵ with a sign determined by the gradient of an error function $E(x, t)$ with respect to the pixel values:
$$x' = x + \epsilon \text{sign}(\nabla_x E(x, t)),$$
 where t is the true label of x .

Synthetic images

DeepDream: Determine which nodes in a particular hidden layer of the network respond strongly to a particular image and then modify the image to amplify those responses:

- We apply an image to the input of the network and forward propagate through to some particular layer.
- We then set the backpropagation δ variables for that layer equal to the pre-activations of the nodes and run backpropagation to the input layer to get a gradient vector over the pixels of the image.
- Finally, we modify the image by taking a small step in the direction of the gradient vector.

Synthetic images

Figure: Examples of DeepDream applied to an image

