

Numerical Simulation Laboratory

LECTURE 12

MACHINE LEARNING

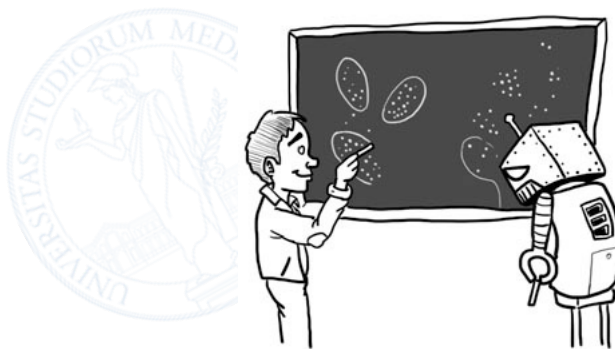


1

Outline

2

- Introduction to Convolutional Neural Networks

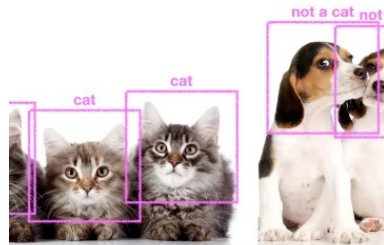


2




Locality and symmetries in the dataset

3

- One of the core lessons of physics is that we should exploit **symmetries** and **invariances** when analyzing physical systems.
- Properties such as **locality and translational invariance** are often **built directly into the physical laws**. For example, we know that in many cases it is sufficient to consider only **local couplings in our Hamiltonians**, or that a model, in which local couplings depend only on the distances among the degrees of freedom, is **translationally invariant** and thus the total momentum is conserved, etc.
- Like physical systems, **many datasets** (especially datasets from Physics) **and supervised learning tasks also possess additional symmetries and structure**.
- For instance, consider a supervised learning task where we want to label images from some dataset as being pictures of cats or not. **Our statistical procedure must first learn features** associated with cats.



3

Cat
Cat
Representation

4

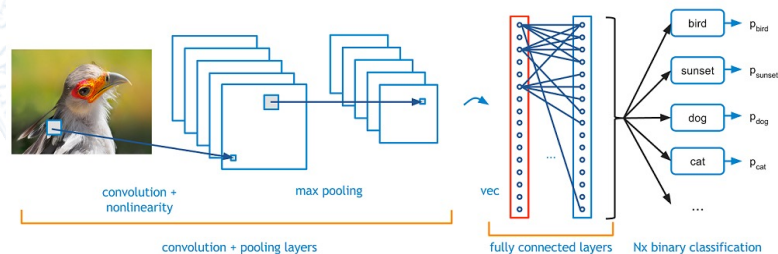
- Because a cat is a physical object, we know that these **features are likely to be local** (whiskers, tails, eyes, etc).
- We also know that the cat can be anywhere in the image. Thus, it does not really matter where in the picture these features occur (though relative positions of features likely do matter). This is a manifestation of **translational invariance** that is built into our supervised learning task.
- This example makes clear that, like many physical systems, many ML tasks (especially in the context of image processing) also possess **additional structure**, such as locality and translation invariance.
- **The all-to-all coupled neural networks in the previous section typically fail to exploit this additional structure**

4

Convolutional Neural Networks (CNNs)

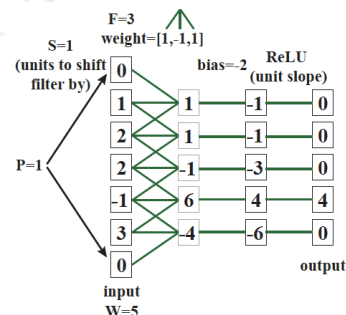
5

- Not surprisingly, the neural networks community realized these problems and designed a class of neural network architectures, **convolutional neural networks or CNNs**, that take advantage of this additional structure (locality and translational invariance)
- A convolutional neural network is a “translationally invariant” neural network that respects locality of the input data.
- There are two kinds of basic layers that make up a CNN: a **convolution layer** that computes the convolution of the input with a series of filters, and **pooling layers** that coarse-grain the input while maintaining locality and spatial structure



5

- For **2D data**, a layer l is characterized by three numbers: **height H_l** , **width W_l** , and the “**number of channels D_l** ”
- The height and width correspond to the sizes of the two-dimensional spatial ($W_l; H_l$)-plane (in neurons) (e.g. if l is the input layer, $W_l \times H_l$ is the number of pixels of the image), and the depth D_l to the “number of channels” in that layer (e.g. 3 channels for an RGB image).
- In general, we will be concerned with **local spatial filters** (often called a receptive field in analogy with neuroscience) that take as inputs a small spatial patch of the previous layer at all depths. **The convolution consists of running this filter over all locations in the spatial plane**
- To demonstrate how this works in practice, let us consider the simple example consisting of a **one-dimensional input of depth 1 (1 channel)**, shown in figure
- In this case, a **filter of size $F \times 1 \times 1$** can be specified by a **vector of weights w of length F**



6

- The **stride, S**, encodes by how many neurons we translate the filter by when performing the convolution.
- In addition, it is common to **pad the input with P zeros**. For an input of width W, the number of neurons (outputs) in the layer is given by $(W-F+2P)/S+1$.
- Below I show you some convolution procedures, for a square input of unit depth:

N.B.: Blue maps are inputs, and cyan maps are outputs.

Diagram illustrating a 1D convolution operation with the following parameters:

- Input width $W=5$
- Filter width $F=3$
- Stride $S=1$ (units to shift filter by)
- Padding $P=1$
- Weight $\text{weight}=[1,-1,1]$
- Bias $\text{bias}=-2$
- ReLU (unit slope)

Input	Weight	Bias	ReLU	Output
0	1	-2	-1	0
1	-1	-2	-3	0
2	1	-2	-1	0
2	-1	-2	-3	0
-1	1	-2	-1	0
3	-1	-2	-3	0
0	1	-2	-1	0

7

- These convolutional layers are interspersed with **pooling layers** that **coarse-grain spatial information by performing a subsampling at each depth**. One common pooling operation is the max pool. A simple example of a **max-pooling** operation is shown in figure
- In a max pool, the spatial dimensions are **coarse-grained** by replacing a small region (say 2×2 neurons) by a single neuron whose output is the maximum value of the output in the region. In physics, this pooling step is **very similar to the decimation step in the Renormalization Group**
- This **reduces the dimension of outputs**. For example, if the region we pool over is 2×2 , then both the height and the width of the output layer will be halved.
- Generally, **pooling operations do not reduce the "number of channels"** of the convolutional layers because pooling is performed separately on each channel.

Diagram illustrating max-pooling operations:

Example 1:

3	0	1	0
0	1	1	1
2	3	2	1
4	1	0	1

channels

3	1
4	2

channels

Example 2:

2	5	2	1
1	1	0	1
1	0	1	0
3	1	0	1

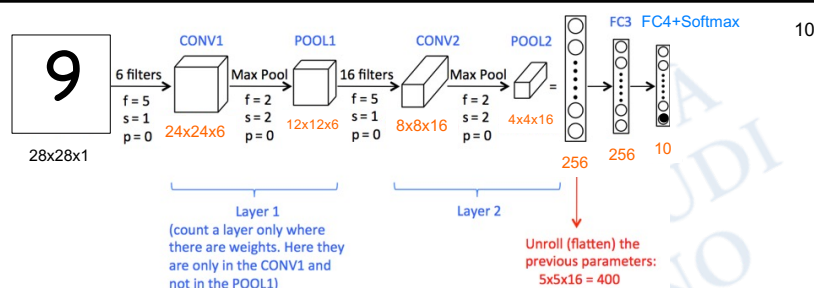
channels

5	2
3	1

channels

8

- In a CNN, the **convolution and max-pool** layers are generally followed by an **all-to-all connected layer** and a classifier such as a soft-max.
- **CNNs are trained using the backpropagation algorithm.** From a backpropagation perspective, **CNNs are almost identical to fully connected NN architectures** except with **tied parameters** ...
- In fact, all **neurons at a given convolutional layer represent the same filter**, and hence can all be **described by a single set of weights and biases**. This reduces the number of free parameters by a factor of $H \times W$ at each layer. For example, for a layer with $D=10^2$ and $H=W=10^2$, this gives a reduction in parameters of nearly 10^6 ... this allows for the training of much larger models!
- The **immense success** of CNNs for image recognition resides in the fact that **filters (receptive fields) learned by the convolution layers should work well on similar tasks**, not just on the ones they were originally trained for. We expect that, since images reflect the natural world, the filters learned by CNNs should transfer over to new tasks with only slight modifications and fine-tuning.



- Let's consider a concrete example of a deep convolutional neural network for digit-image classification where the **input** image size is **28x28x1** (i.e. **grayscale**)
- In the **first layer**, we apply the **convolution** operation with **6 filters** of **5x5** so our output will become **24x24x6**.
- Then we will apply **pooling** with **2x2** filter to reduce the size to **12x12x6**
- In the second layer, we will apply the **convolution** operation with **16 filters** of size **5x5x6**. The output dimensions will become **8x8x16** ...
- on which we will apply **pooling** layer with **2x2** filter and the size will reduce to **4x4x16**.
- Finally, we will pass it through **two fully connected layers** to convert our image matrix **into a classification matrix**.

DNNs vs CNNs in action!

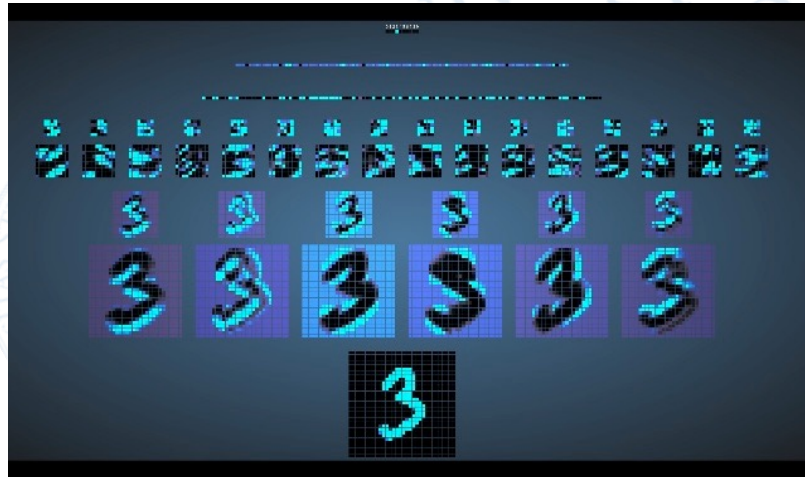
11

- Some fun with CNNs at:

<https://www.cs.cmu.edu/~aharley/vis/fc/flat.html>

Compare with a fully connected DNN at:

<https://www.cs.cmu.edu/~aharley/vis/conv/flat.html>



11

Lecture 12: Suggested books

- P. Mehta, C.H. Wang, A.G.R. Day, and C. Richardson, *A high-bias, low-variance introduction to Machine Learning for physicists*, arXiv:1803.08823 (Feb 2019)
- S. Haykin, *Neural Networks and Learning Machines*, Pearson (2009)
- S. Theodoridis, *Machine Learning. A Bayesian and Optimization Perspective*, Elsevier (2015)
- K.P. Murphy, *Machine Learning. A Probabilistic Perspective*, MIT press (2012)
- J. Kaipio, E. Somersalo, "Statistical and Computational Inverse Problems", Springer

12

12



13