Neural Networks for Data Science Applications

Master's Degree in Data Science

# Lecture 7: Convolutions beyond images

**Lecturer**: S. Scardapane

SAPIENZA
UNIVERSITÀ DI ROMA

Although we focused on images in the previous lecture, many other types of data have a sequential / grid-like structure, albeit with different dimensionality: time-series, audio, videos, text, …

In this lecture we will see how to extend CNNs to these other cases, and also some critical limitations, including their unsuitability in handling long-term or non-regular dependencies. This will provide motivation for the next class of models we will study, **transformers**.

# Convolutions beyond images

## 1D and 3D convolutional layers

CNNs can be extended easily to other domains having grid-like structure of various dimensions.

For example, consider $n$ steps of a time-series $\mathbf{x}_0, \ldots, \mathbf{x}_{n-1} \in \mathbb{R}^c$, each step having $c$ features (e.g., $c$ different readings from different sensors).

We represent it using a matrix $\underset{(n,c)}{\mathbf{X}}$, and define a 1D convolutional layer of receptive field $2k$ as:

$$\underset{(n,c')}{[\mathbf{H}]_{i,d}} = \phi \left( \sum_{i'=-k}^{+k} \sum_{z=1}^{c} \underset{(s,c,c')}{[W]_{i'+k+1,z,d}} [\mathbf{X}]_{i+i',z} \right) .$$

For temporal domains, it is useful to define *causal* (masked) versions of the convolution operation.

For a temporal sequence $\mathbf{X}_{(n,c)}$, a **causal model** $\mathbf{H}_{(n,c')} = f(\mathbf{X})$ is such that $[\mathbf{H}]_i$ depends only on $[\mathbf{X}]_j$ for $j \leq i$.

A causal convolution is defined by removing right neighbours:

$$[\mathbf{H}]_{i,d} \atop (n,c') = \phi \left( \sum_{i'=0}^{k} \sum_{z=1}^{c} [W]_{i'+1,z,d} \atop (k+1,c,c')} [\mathbf{X}]_{i-i',z} \right) .$$
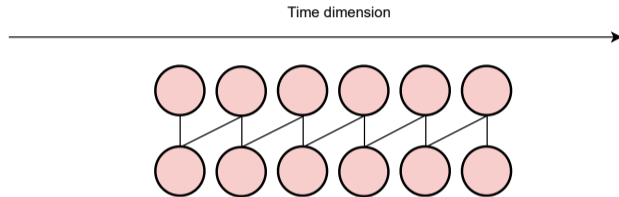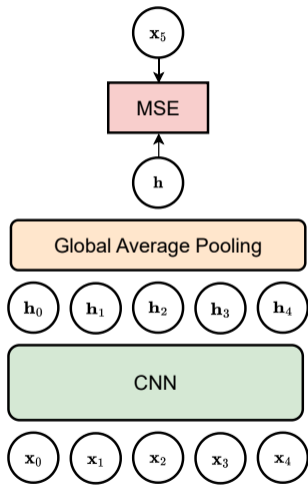
Figure 1: A 1D convolutional layer with causal convolutions. All links going right-to-left are removed. As a result, for each unit the output only depends on the previous time-instants.
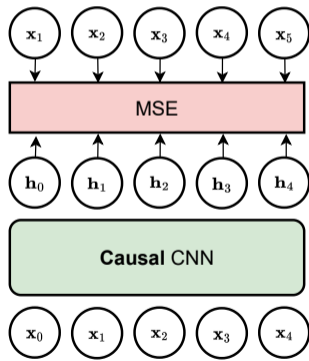
For time-series, a common task is **forecasting**, i.e., predicting the next step in the time-series. With a causal model, we have two options:

▶ Pool the output representation $H$ over $n$, and apply a regressor head to predict $x_n$ (also valid for non-causal models).

▶ Define a *shifted* target $Y = [x_1, \ldots, x_n]$, and train the model such that $H \approx Y$, i.e., at each time step the network predicts the next one. This is only possible with a causal model, otherwise information would 'leak' from the input.

(a) Single-step forecasting     (b) Multi-step forecasting

Models of the second form are useful because they can be used to generate data in an **autoregressive** way. We prompt the model with the beginning of the sequence, and let it forecast the next step. Then, we concatenate it to the input and prompt again the model to generate another step, and so on.

During training, we can feed the network only with true samples (**teacher forcing**), or mix some of the network's own predictions. Extrapolating outside the training data windows may still be difficult.
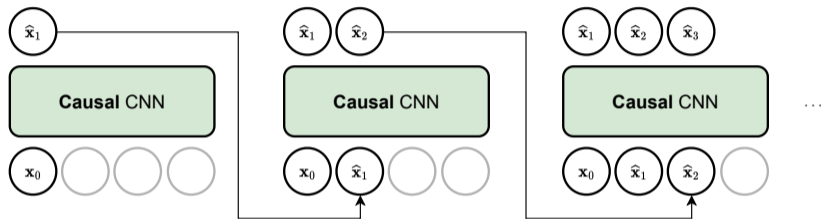
Figure 2: Example of autoregressive generation with a single step in input, and the model being called recursively multiple times. We show in gray padded values.

Now consider a sequence of $n$ video frames, each of shape $(h, w, c)$. The result is a tensor $X_{(n,h,w,c)}$.

A 3D version of the convolutional layer is straightforward:

$$[H]_{t,i,j,d} \atop {(n,h,w,c')}} = \phi \left( \sum_{t'=-k}^{+k} \sum_{i'=-k}^{+k} \sum_{j'=-k}^{+k} \sum_{z=1}^{c} [W]_{t'+k+1,i'+k+1,j'+k+1,z,d} [X]_{t+t',i+i',j+j',z} \right)$$

**Exercise**: define a convolutional operator which is causal in the temporal dimension, and normal in the spatial dimensions.

# Convolutions beyond images

Applying CNNs to audio

Many real-world problems require the classification of **audio samples**, e.g.:

1. Speech / non-speech identification (*is he/she speaking now?*);
2. Language identification (*is it Italian?*);
3. Genre / mood classification (*is it rock?*);
4. Determining the leading instrument;
5. Event recognition (*is someone shooting?*);
6. Scene recognition (*are they in a bus? at a restaurant?*).

An audio is a 1D sequence of **samples**, obtained with a certain **sampling rate**, typically in one or two **channels**.



**Figure 3:** Simple example of an audio **waveform** (image source).

For a given audio, the number of samples can be very high: at a **resolution** of 44kHz, we have almost half a million samples for each 10 seconds.
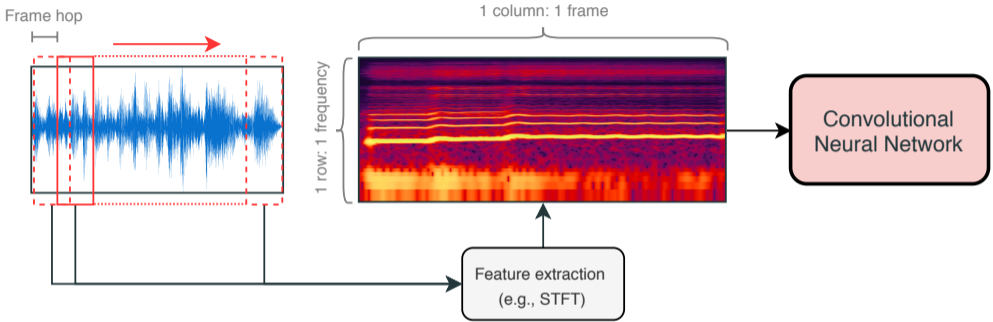
**Figure 4:** To simplify the problem we can sweep a window of fixed size (e.g., 40ms) over the signal, compute a frequency representation (e.g., with a short-term Fourier transform), and merge everything into a single **spectrogram**.

To work with mini-batching, you still need to **zero-pad** all dimensions to a maximum length. For example, you can use `padded_batch()` with `tf.data`:

```
1  im1 = tf.random.normal(shape=(4, 3)) # Image 1 has shape (4, 3)
2  im2 = tf.random.normal(shape=(4, 2)) # Image 2 has shape (4, 2)
3
4  # We need to use from_generator instead of from_tensor_slices
5  data = tf.data.Dataset.from_generator([im1, im2].__iter__,
6  output_types=tf.float64)
7
8  for xb in data.padded_batch(2, padded_shapes=(4, 3)):
9  print(xb) # Zero-padded to (4, 3)
```

By having sufficient computational power, one can also work on the raw audio waveform.

When using convolutions, a key idea is the use of **dilated convolutions** (a.k.a. **atrous** convolutions, from French *à trous*), where neighbours are selected with exponentially increasing steps.

In this way, the receptive field of an item increases *exponentially* with the number of layers.
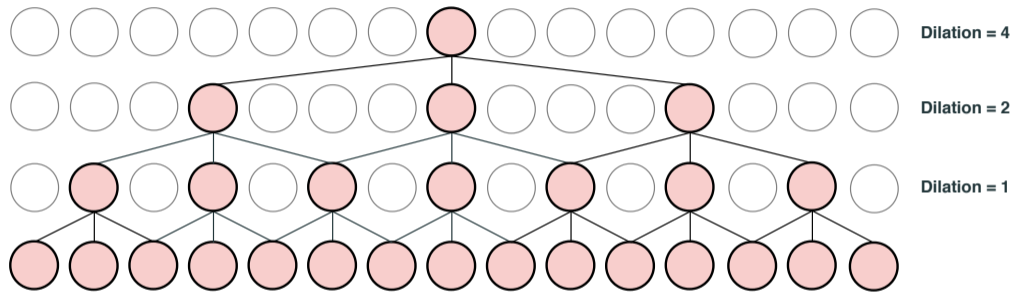
**Figure 5:** Example of dilated (atrous) convolution, with increasing dilation for each layer.

**WaveNet** (2016) was one of the first models to propose a 1D convolutional network dilated and causal convolutions for generating speech. In their design, dilation factors increased from 1 to an upper bound (e.g., 512), before restarting from 1.

Note that autoregressive generation of speech can be very slow, because the model must be called hundreds of thousands of times: although convolutive models have given way to transformers today, this problem is still true (e.g., the recent research on **speculative decoding** in LLMs).

---

Oord, A.V.D., et al., 2016. **WaveNet: A generative model for raw audio**. *arXiv preprint arXiv:1609.03499.*

# Convolutions beyond images

Applying CNNs to text

Text data is another field with a vast range of possible applications, e.g.:

1. Recognition of a topic (*is it talking about soccer?*);
2. Hate speech recognition (*is it respecting our code of conduit?*);
3. Sentiment analysis (*is it a positive review?*);
4. Web page classification (*is it an e-commerce website?*).

Causal models applied to text are the backbone of LLM models such as ChatGPT.

Text must be preprocessed properly to be handled by a neural network. At the very least we need to perform two basic operations:

1. **Tokenize** the text to split it in into basic units (**tokens**) that will form a sequence.
2. **Embed** each token (one or more characters), i.e., convert it into numerical features usable by a neural network.

By the end of this procedure we should have a matrix $\underset{(n,d)}{X}$ of $n$ tokens, each of size $d$.

---

In our discussion, we are ignoring a large number of possible preprocessing operations, e.g., stemming, stop word removal, etc.

# Convolutions beyond images

Text tokenization

['c', 'l', 'a', ..., 't', '!']   ['clas', 'si', ..., 'text']   ['classifying', 'this', 'text']

Character-based model   Sub-word-based model   Word-based model

Character tokenizer   Sub-word tokenizer   Word tokenizer
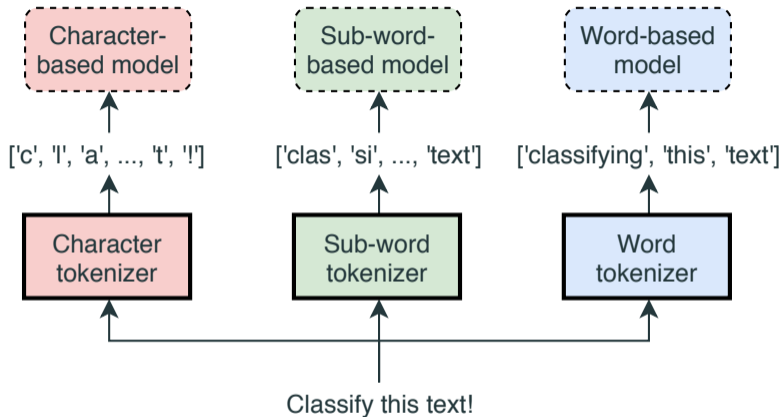
Classify this text!

**Figure 6:** Tokenization can be performed at various levels in the sentence, and it is generally handled by an external library (e.g., Spacy).

Tokenization is crucial because it determines what the network will 'see'. Consider that v15 of Unicode has almost 150k possible symbols, which can be combined in many forms (numbers, dates, function names, …).

**Subword** tokenization is very common, especially when learning the sub-words from large corpora of texts (e.g., **byte pair encoding**, BPE). See `https://huggingface.co/docs/transformers/tokenizer_summary` for a summary.

Also see `https://platform.openai.com/tokenizer` to visualize a tokenization in practice.

# Convolutions beyond images

Token embeddings

The simplest vectorial embedding for text is a **one-hot encoding** according to a predefined **dictionary**:

► Character-level: each character is represented by a 1-of-C binary vector, where *C* is the number of allowable characters.

► Sub-word/word-level: similar, but each word/sub-word is represented with respect to a fixed vocabulary of sub-words / words.

► Sentence-level: each sentence can be represented by summing the one-hot encodings for the single tokens (**bag-of-words**).

One-hot vectors are very simplistic representations of the information contained in text. A more general solution is to *learn* a set of embeddings:

1. For every possible token $c$, initialize randomly a fixed-size vector $\mathbf{v}_c$.
2. During training, substitute each token in the sequence with the corresponding vector (**look-up**).
3. The set of vectors $\mathbf{v}_c$ can be optimized together with the parameters of the neural network by doing gradient descent.
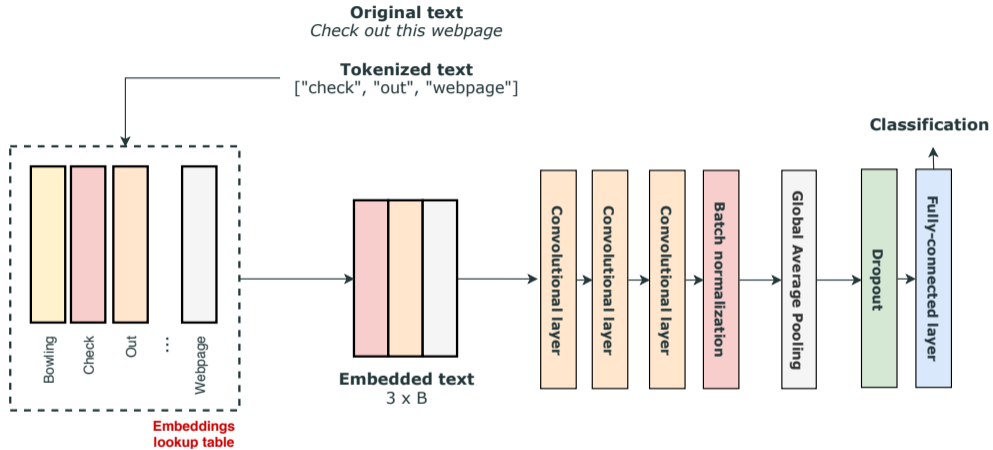
Figure 7: Using custom embeddings for text classification.

Custom embeddings are extremely simple to train within a TensorFlow model:

```python
1  model = Sequential()
2
3  # Get embeddings for each token
4  # Input must be (batch_size, max_sentence_length)
5  # Every element of the input is an index [0, ..., dictionary_size-1]
6  model.add(Embedding(dictionary_size, B, input_length=max_sentence_lengt
7
8  # Optional: Get average embedding for the sentence
9  model.add(GlobalAveragePooling1D())
10
11 # ...
```

One-hot encodings are *sparse*, in the sense that most values are 0. In addition, they are not very informative: for example, the distance between any two encodings $\mathbf{e}_0$ and $\mathbf{e}_1$ is either 0 or $\sqrt{2}$.

Instead, trained embeddings can capture a rich semantic underlying the data, which is reflected in the possibility of doing algebraic manipulations on the embeddings themselves, e.g., see `https://projector.tensorflow.org/`. However, they can also capture **biases** of the data.[1]

---

[1]Bolukbasi, T. et al., 2016. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings. *NeurIPS*.

Text embeddings can also be **pre-trained** using a variety of algorithms:

1. Word2Vec (Mikolov et al., 2013);
2. Global Vectors for Word Representation (GloVe) (Pennington et al., 2014);
3. Embeddings from Language Models (ELMo) (Peters et al., 2018);
4. Generative Pre-Training (GPT) (Radford et al., 2018);
5. Bidirectional Encoder Representations from Transformers (BERT) (Devlin et al., 2018).