

– IN5550 –

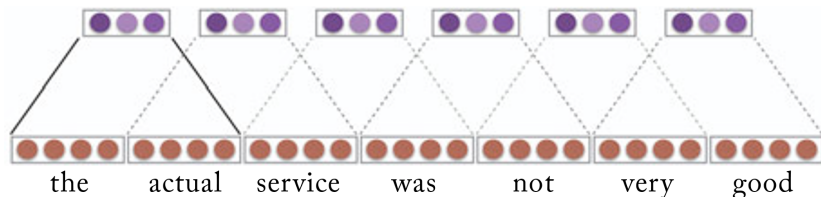
*Neural Methods in Natural Language Processing*

CNNs, Part 4: Additional comments and advanced options

Erik Velldal

Language Technology Group (LTG)  
University of Oslo





- ▶ CNNs improve on CBOW in also capturing **ordered** context.
- ▶ But still rather limited; only relationships **local** to windows of size  $k$ .
- ▶ Due to long-range compositional effects in natural language semantics, we'll often want to model as much context as feasible.
- ▶ One option is to just increase the filter size  $k$ .
- ▶ More powerful: a stack of convolution layers applied one after the other:
- ▶ **Hierarchical convolutions.**



- ▶ Let  $p_{1:m} = \text{CONV}_{U,b}^k(w_{1:n})$  be the result of applying a convolution (with parameters  $U$  and  $b$ ) across  $w_{1:n}$  with window-size  $k$ .
- ▶ Can have a succession of  $r$  layers that feed into each other:

$$p_{1:m_1}^1 = \text{CONV}_{U^1,b^1}^{k_1}(w_{1:n})$$

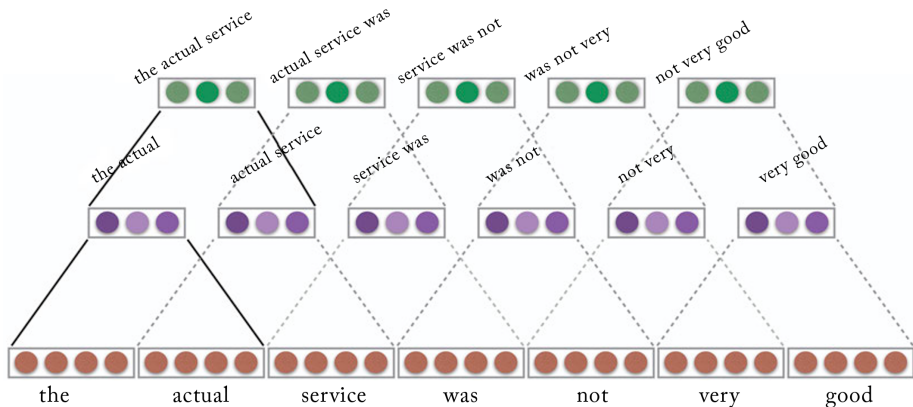
$$p_{1:m_2}^2 = \text{CONV}_{U^2,b^2}^{k_2}(p_{1:m_1}^1)$$

...

$$p_{1:m_r}^r = \text{CONV}_{U^r,b^r}^{k_r}(p_{1:m_{r-1}}^{r-1})$$

- ▶ The vectors  $p_{1:m_r}^r$  capture increasingly larger effective windows.

# Two-layer hierarchical convolution with $k = 2$

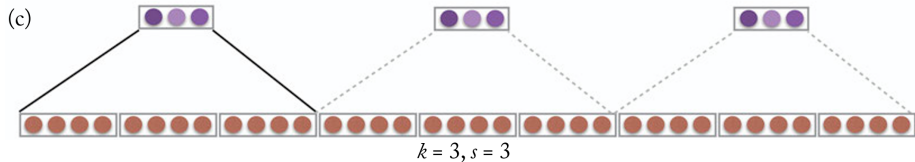
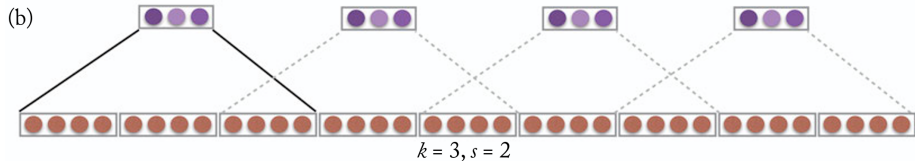
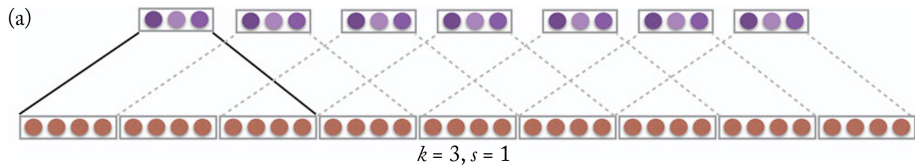


- ▶ Two different but related effects of adding layers:
- ▶ **Larger receptive field** wrt the input at each step: convolutions of successive layers see more of the input.
- ▶ Can learn more **abstract feature combinations**.



- ▶ The **stride size** specifies by how much we shift a filter at each step.
- ▶ So far we've considered convolutions with a stride size of 1: we slide the window by increments of 1 across the word sequence.
- ▶ But using larger strides is possible.
- ▶ Can slide the window with increments of e.g. 2 or 3 words at the time.
- ▶ A larger stride size leads to fewer applications of the filter and a shorter output sequence  $p_{1:m}$ .

# $k = 3$ and stride sizes 1, 2, 3





- ▶ **Dilated convolutions**: skip some of the positions within the filters (or equivalently, introduce zero weights).
- ▶ Gives a wider filter region but with the same number of parameters.



- ▶ **Dilated convolutions**: skip some of the positions within the filters (or equivalently, introduce zero weights).
- ▶ Gives a wider filter region but with the same number of parameters.
- ▶ **Multi-channel CNNs**: each channel providing a different representation of the input. Apply convolutions to each.





- ▶ **Dilated convolutions**: skip some of the positions within the filters (or equivalently, introduce zero weights).
- ▶ Gives a wider filter region but with the same number of parameters.
- ▶ **Multi-channel CNNs**: each channel providing a different representation of the input. Apply convolutions to each.
- ▶ Hierarchical convolutions can be combined with **parameter tying**:
  - ▶ Reusing the same  $U$  and  $b$  across layers.
  - ▶ Allows for using an unbounded number of layers, to extend the receptive field to arbitrary-sized inputs.



- ▶ **Dilated convolutions**: skip some of the positions within the filters (or equivalently, introduce zero weights).
- ▶ Gives a wider filter region but with the same number of parameters.
- ▶ **Multi-channel CNNs**: each channel providing a different representation of the input. Apply convolutions to each.
- ▶ Hierarchical convolutions can be combined with **parameter tying**:
  - ▶ Reusing the same  $U$  and  $b$  across layers.
  - ▶ Allows for using an unbounded number of layers, to extend the receptive field to arbitrary-sized inputs.
- ▶ **Skip-connections** can be useful for deep CNNs:
  - ▶ The output from one layer is passed to not only the next but also subsequent layers (ResNets, Highway nets, DenseNets, ...)



- ▶ While hugely successful in **image processing**, CNNs have had less impact in **NLP** (typically also much more shallow networks).
- ▶ Main use; **document** and **sentence classification** (e.g. for topic or polarity classification).
- ▶ Although they have also been applied to more 'structured' tasks like aspect-based SA and relation extraction.
- ▶ As of today, CNNs often applied at the **character-level**, to generate more robust word representations (typically concatenated with word embeddings before being passed to an RNN).



- ▶ **Convolutional networks** can learn to represent large  $n$ -grams efficiently...
- ▶ ...without blowing up the parameter space and without having to represent the whole vocabulary.
- ▶ **Parameter sharing** (shared weights in all applications of a given filter)
- ▶ Multiple filters; act as specialized **feature extractors**.



- ▶ **Convolutional networks** can learn to represent large  $n$ -grams efficiently...
- ▶ ...without blowing up the parameter space and without having to represent the whole vocabulary.
- ▶ **Parameter sharing** (shared weights in all applications of a given filter)
- ▶ Multiple filters; act as specialized **feature extractors**.
- ▶ But not designed for modeling **sequential** language data:
- ▶ do not offer a very natural way of modeling **long-range and structured dependencies**.



- ▶ Lends itself well to GPU computations; optimized for matrix convolutions.
- ▶ **Easily parallelized**: each 'region' that a convolutional filter operates on is independent of the others; the entire input can be processed concurrently.
- ▶ Each filter also independent.



- ▶ Lends itself well to GPU computations; optimized for matrix convolutions.
- ▶ **Easily parallelized**: each 'region' that a convolutional filter operates on is independent of the others; the entire input can be processed concurrently.
- ▶ Each filter also independent.
- ▶ The cost of this is that we have to **stack** convolutions into deep layers in order to 'view' the entire input
- ▶ Those layers are in fact calculated sequentially...
- ▶ But the **calculations at each layer happen concurrently** and each individual computation is small.
- ▶ In practice, CNNs are much faster than RNNs.



- ▶ Lends itself well to GPU computations; optimized for matrix convolutions.
- ▶ **Easily parallelized**: each 'region' that a convolutional filter operates on is independent of the others; the entire input can be processed concurrently.
- ▶ Each filter also independent.
- ▶ The cost of this is that we have to **stack** convolutions into deep layers in order to 'view' the entire input
- ▶ Those layers are in fact calculated sequentially...
- ▶ But the **calculations at each layer happen concurrently** and each individual computation is small.
- ▶ In practice, CNNs are much faster than RNNs.