

Convolution and Pooling

Benjamin Roth, Nina Poerner

CIS LMU München

November 28, 2018

Gliederung

1 Convolutional Neural Networks

- Convolution with one filter
- Convolution with N filters
- What does convolution do?

2 Pooling

3 Application to NLP

4 Comparison: RNN vs. CNN

Outline

1 Convolutional Neural Networks

- Convolution with one filter
- Convolution with N filters
- What does convolution do?

2 Pooling

3 Application to NLP

4 Comparison: RNN vs. CNN

Convolutional Neural Networks (CNNs)

- Technique from Computer Vision (e.g., object recognition in images)
- Alternative to RNNs for many (not all) NLP tasks
- General idea: Filter bank with N learnable filters

Outline

1 Convolutional Neural Networks

- Convolution with one filter
- Convolution with N filters
- What does convolution do?

2 Pooling

3 Application to NLP

4 Comparison: RNN vs. CNN

Convolution with one filter

- For now, assume we have only one filter
- Move filter over input with step size (stride) s (here: 1)
- At every position, multiply filter and input entries together (elementwise), and sum the results into a single value

Input					Filter			Output	
0	0	-1			2	1			
1	-1	0		*	0	-1	=		
-2	0	0							
-1	-3	1							

- Filter size: 2×2
- Input size: 3×4

Convolution with one filter

- For now, assume we have only one filter
- Move filter over input with step size (stride) s (here: 1)
- At every position, multiply filter and input entries together (elementwise), and sum the results into a single value

Input

0	0	-1
1	-1	0
-2	0	0
-1	-3	1

Filter

2	1
0	-1

Output

The diagram illustrates a 1D convolution operation. The input is a 4x3 grid of numbers. A 2x2 red box highlights the top-left 2x2 subgrid of the input, which corresponds to the first output value. The filter is a 2x2 grid of numbers. The output is a 3x2 grid of empty boxes, representing the result of the convolution operation.

- Filter size: 2×2
- Input size: 3×4

Convolution with one filter

- For now, assume we have only one filter
- Move filter over input with step size (stride) s (here: 1)
- At every position, multiply filter and input entries together (elementwise), and sum the results into a single value

Input

0	0	-1
1	-1	0
-2	0	0
-1	-3	1

Filter

2	1
0	-1

Output

1	

The diagram illustrates a 1D convolution operation. The input is a 4x3 grid of numbers. A 2x2 red box highlights the top-left 2x2 subgrid of the input, which corresponds to the values 0, 0, 1, and -1. This subgrid is multiplied elementwise by the 2x2 filter, which contains the values 2, 1, 0, and -1. The result of this multiplication is the top-left element of the output grid, which is 1. The output grid is a 3x2 grid of numbers.

- Filter size: 2×2
- Input size: 3×4

Convolution with one filter

- For now, assume we have only one filter
- Move filter over input with step size (stride) s (here: 1)
- At every position, multiply filter and input entries together (elementwise), and sum the results into a single value

Input

0	0	-1
1	-1	0
-2	0	0
-1	-3	1

Filter

2	1
0	-1

Output

1	

The diagram illustrates a 1D convolution operation. The input is a 4x3 grid of numbers. A 2x2 sub-region of the input is highlighted with red and green borders, representing the current filter position. The filter is a 2x2 grid of numbers. The output is a 3x2 grid, with the top-left cell containing the value 1, which is the result of the convolution at the current position.

- Filter size: 2×2
- Input size: 3×4

Convolution with one filter

- For now, assume we have only one filter
- Move filter over input with step size (stride) s (here: 1)
- At every position, multiply filter and input entries together (elementwise), and sum the results into a single value

Input

0	0	-1
1	-1	0
-2	0	0
-1	-3	1

Filter

2	1
0	-1

Output

1	-1

The diagram illustrates a 1D convolution operation. The input is a 4x3 grid of numbers. A 2x2 sub-region of the input is highlighted with red and green borders, representing the current filter position. The filter is a 2x2 grid of numbers. The output is a 3x2 grid, with the top row containing the results of the convolution at the first two positions.

- Filter size: 2×2
- Input size: 3×4

Convolution with one filter

- For now, assume we have only one filter
- Move filter over input with step size (stride) s (here: 1)
- At every position, multiply filter and input entries together (elementwise), and sum the results into a single value

Input

0	0	-1
1	-1	0
-2	0	0
-1	-3	1

Filter

2	1
0	-1

Output

1	-1

- Filter size: 2×2
- Input size: 3×4

Convolution with one filter

- For now, assume we have only one filter
- Move filter over input with step size (stride) s (here: 1)
- At every position, multiply filter and input entries together (elementwise), and sum the results into a single value

Input

0	0	-1
1	-1	0
-2	0	0
-1	-3	1

Filter

2	1
0	-1

Output

1	-1
1	

- Filter size: 2×2
- Input size: 3×4

Convolution with one filter

- For now, assume we have only one filter
- Move filter over input with step size (stride) s (here: 1)
- At every position, multiply filter and input entries together (elementwise), and sum the results into a single value

Input

0	0	-1
1	-1	0
-2	0	0
-1	-3	1

Filter

2	1
0	-1

Output

1	-1
1	

- Filter size: 2×2
- Input size: 3×4

Convolution with one filter

- For now, assume we have only one filter
- Move filter over input with step size (stride) s (here: 1)
- At every position, multiply filter and input entries together (elementwise), and sum the results into a single value

Input

0	0	-1
1	-1	0
-2	0	0
-1	-3	1

Filter

2	1
0	-1

Output

1	-1
1	-2

The diagram illustrates a 1D convolution operation. The input is a 4x3 grid of numbers. A 2x2 filter is applied to the input, moving it across the input grid with a step size of 1. The output is a 3x2 grid of numbers, where each cell contains the result of the convolution operation at that position. The output grid shows the results for the first two positions: 1 and -1 in the first row, and 1 and -2 in the second row. The third row is empty, indicating that the filter has moved beyond the input grid.

- Filter size: 2×2
- Input size: 3×4

Convolution with one filter

- For now, assume we have only one filter
- Move filter over input with step size (stride) s (here: 1)
- At every position, multiply filter and input entries together (elementwise), and sum the results into a single value

Input

0	0	-1
1	-1	0
-2	0	0
-1	-3	1

Filter

2	1
0	-1

Output

1	-1
1	-2

The diagram illustrates a 1D convolution operation. The input is a 3x4 grid of numbers. The filter is a 2x2 grid of numbers. The output is a 3x2 grid of numbers. The operation is represented by the equation: Input * Filter = Output. The output values are calculated by elementwise multiplication of the filter with the input, followed by summing the results. The output values are 1, -1, 1, -2, and the remaining cells are empty.

- Filter size: 2×2
- Input size: 3×4

Convolution with one filter

- For now, assume we have only one filter
- Move filter over input with step size (stride) s (here: 1)
- At every position, multiply filter and input entries together (elementwise), and sum the results into a single value

Input

0	0	-1
1	-1	0
-2	0	0
-1	-3	1

Filter

2	1
0	-1

Output

1	-1
1	-2
-1	

The diagram illustrates a 1D convolution operation. The input is a 3x4 grid of numbers. The filter is a 2x2 grid of numbers. The output is a 3x2 grid of numbers. The operation is represented by the equation: Input * Filter = Output. The output values are calculated by multiplying the filter values with the corresponding input values and summing the results. For example, the first output value is 1, which is the sum of 0*2 + 0*1.

- Filter size: 2×2
- Input size: 3×4

Convolution with one filter

- For now, assume we have only one filter
- Move filter over input with step size (stride) s (here: 1)
- At every position, multiply filter and input entries together (elementwise), and sum the results into a single value

Input

0	0	-1
1	-1	0
-2	0	0
-1	-3	1

Filter

2	1
0	-1

Output

1	-1
1	-2
-1	

The diagram illustrates a 1D convolution operation. The input is a 4x3 grid of numbers. The filter is a 2x2 grid of numbers. The output is a 3x2 grid of numbers. The operation is represented by the equation: Input * Filter = Output. The output values are calculated by elementwise multiplication of the filter with the input, followed by summing the results. The output values are: 1, -1, 1, -2, -1.

- Filter size: 2×2
- Input size: 3×4

Convolution with one filter

- For now, assume we have only one filter
- Move filter over input with step size (stride) s (here: 1)
- At every position, multiply filter and input entries together (elementwise), and sum the results into a single value

Input

0	0	-1
1	-1	0
-2	0	0
-1	-3	1

Filter

2	1
0	-1

Output

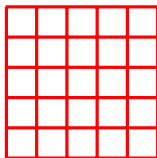
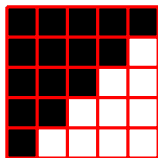
1	-1
1	-2
-1	-1

The diagram illustrates a 1D convolution operation. The input is a 4x3 grid of numbers. The filter is a 2x2 grid of numbers. The output is a 3x2 grid of numbers. The input grid is: $\begin{bmatrix} 0 & 0 & -1 \\ 1 & -1 & 0 \\ -2 & 0 & 0 \\ -1 & -3 & 1 \end{bmatrix}$. The filter grid is: $\begin{bmatrix} 2 & 1 \\ 0 & -1 \end{bmatrix}$. The output grid is: $\begin{bmatrix} 1 & -1 \\ 1 & -2 \\ -1 & -1 \end{bmatrix}$. The operation is represented by an asterisk (*) between the input and filter grids, and an equals sign (=) between the result and the output grid.

- Filter size: 2×2
- Input size: 3×4

Building an edge detector filter

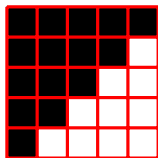
- Assume that -1 means black and $+1$ means white
- We want to build a filter that can detect diagonal edges where the upper left side is dark and the lower right side is bright
- = a filter that calculates a high positive number on windows that look like this:



- In CNNs, the filters are not manually chosen, but learned with gradient descent

Building an edge detector filter

- Assume that -1 means black and +1 means white
- We want to build a filter that can detect diagonal edges where the upper left side is dark and the lower right side is bright
- = a filter that calculates a high positive number on windows that look like this:

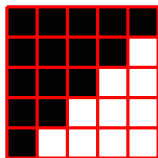


-3	-3	-3	-3	-3
-3	-3	-3	-3	3
-3	-3	-3	3	3
-3	-3	3	3	3
-3	3	3	3	3

- In CNNs, the filters are not manually chosen, but learned with gradient descent

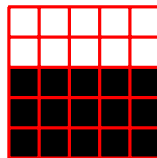
Building an edge detector filter

- Assume that -1 means black and $+1$ means white
- We want to build a filter that can detect diagonal edges where the upper left side is dark and the lower right side is bright
- = a filter that calculates a high positive number on windows that look like this:



-3	-3	-3	-3	-3
-3	-3	-3	-3	3
-3	-3	-3	3	3
-3	-3	3	3	3
-3	3	3	3	3

How would the filter
react to this window?



- In CNNs, the filters are not manually chosen, but learned with gradient descent

Convolution with one filter: Tensor sizes

- Most images are not 2D but 3D
 - ▶ 3rd dimension is # channels, e.g., RGB values
 - ▶ image height \times image width \times # channels
- As a consequence, each filter is also 3D
 - ▶ filter height \times filter width \times # channels
- The operation stays the same, with an additional summation over the channel dimension

Outline

1 Convolutional Neural Networks

- Convolution with one filter
- Convolution with N filters
- What does convolution do?

2 Pooling

3 Application to NLP

4 Comparison: RNN vs. CNN

Convolution with N filters

- Apply N different filters of the same size $\rightarrow N$ matrices with the same size
- Stack the N matrices on top of each other \rightarrow 3D tensor, where the last dimension is N
- Also known as a feature map
- Feature map is slightly smaller than input (why?)

Convolution with N filters

- Apply N different filters of the same size $\rightarrow N$ matrices with the same size
- Stack the N matrices on top of each other \rightarrow 3D tensor, where the last dimension is N
- Also known as a feature map
- Feature map is slightly smaller than input (why?)
- Because a filter of size k fits into an input of size h only $h - k + 1$ times
- ... unless we pad the input with zeros

Convolution with N filters: Tensor sizes

- Tensor sizes:

- ▶ **Input 3D**: input height \times input width \times # channels (if this is the first layer, otherwise # filters of previous layer)
- ▶ **Parameter tensor 4D**: filter height \times filter width \times # channels \times #filters
- ▶ **Output 3D**: input height* \times input width* \times #filters
- ▶ *height and width are slightly reduced by convolution unless we do padding

Outline

1 Convolutional Neural Networks

- Convolution with one filter
- Convolution with N filters
- What does convolution do?

2 Pooling

3 Application to NLP

4 Comparison: RNN vs. CNN

What does convolution do?

- Contextualization: Feature vector computed for position (i, j) contains info from $(i - k, j - k)$ to $(i + k, j + k)$ (where k is filter size).

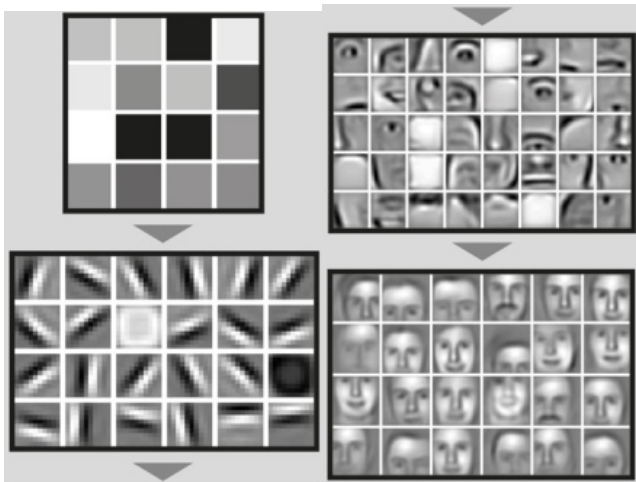
What does convolution do?

- Contextualization: Feature vector computed for position (i, j) contains info from $(i - k, j - k)$ to $(i + k, j + k)$ (where k is filter size).
- Locality-preserving: In one convolution layer, info can travel no further than k positions

What does convolution do?

- Contextualization: Feature vector computed for position (i, j) contains info from $(i - k, j - k)$ to $(i + k, j + k)$ (where k is filter size).
- Locality-preserving: In one convolution layer, info can travel no further than k positions
- Computer Vision: Many convolutional layers applied one after another
- Typical nonlinearity between convolution layers: ReLU
- With every layer, feature maps become more complex
- Pixels \rightarrow edges \rightarrow shapes \rightarrow small objects \rightarrow bigger, compositional

Convolution



Source: Computer science: The learning machines. Nature (2014).

Outline

1 Convolutional Neural Networks

- Convolution with one filter
- Convolution with N filters
- What does convolution do?

2 Pooling

3 Application to NLP

4 Comparison: RNN vs. CNN

Pooling

- Often applied between convolution steps
- Divide feature map into “grid”
- Combine vectors inside the same grid cell with some operator
- Most popular: Average pooling, Max pooling
- Max pooling: only select maximum value for each dimension
- “Feature detector”, “Cat neuron fires”

2	1	2	0
1	0	2	0
0	4	2	3
0	5	1	0

Max pooled

Pooling

- Often applied between convolution steps
- Divide feature map into “grid”
- Combine vectors inside the same grid cell with some operator
- Most popular: Average pooling, Max pooling
- Max pooling: only select maximum value for each dimension
- “Feature detector”, “Cat neuron fires”

2	1	2	0
1	0	2	0
0	4	2	3
0	5	1	0

Max pooled

Pooling

- Often applied between convolution steps
- Divide feature map into “grid”
- Combine vectors inside the same grid cell with some operator
- Most popular: Average pooling, Max pooling
- Max pooling: only select maximum value for each dimension
- “Feature detector”, “Cat neuron fires”

2	1	2	0
1	0	2	0
0	4	2	3
0	5	1	0

Max pooled

2	

Pooling

- Often applied between convolution steps
- Divide feature map into “grid”
- Combine vectors inside the same grid cell with some operator
- Most popular: Average pooling, Max pooling
- Max pooling: only select maximum value for each dimension
- “Feature detector”, “Cat neuron fires”

2	1	2	0
1	0	2	0
0	4	2	3
0	5	1	0

Max pooled

2	

Pooling

- Often applied between convolution steps
- Divide feature map into “grid”
- Combine vectors inside the same grid cell with some operator
- Most popular: Average pooling, Max pooling
- Max pooling: only select maximum value for each dimension
- “Feature detector”, “Cat neuron fires”

2	1	2	0
1	0	2	0
0	4	2	3
0	5	1	0

Max pooled

2	2

Pooling

- Often applied between convolution steps
- Divide feature map into “grid”
- Combine vectors inside the same grid cell with some operator
- Most popular: Average pooling, Max pooling
- Max pooling: only select maximum value for each dimension
- “Feature detector”, “Cat neuron fires”

2	1	2	0
1	0	2	0
0	4	2	3
0	5	1	0

Max pooled

2	2

Pooling

- Often applied between convolution steps
- Divide feature map into “grid”
- Combine vectors inside the same grid cell with some operator
- Most popular: Average pooling, Max pooling
- Max pooling: only select maximum value for each dimension
- “Feature detector”, “Cat neuron fires”

2	1	2	0
1	0	2	0
0	4	2	3
0	5	1	0

Max pooled

2	2
5	

Pooling

- Often applied between convolution steps
- Divide feature map into “grid”
- Combine vectors inside the same grid cell with some operator
- Most popular: Average pooling, Max pooling
- Max pooling: only select maximum value for each dimension
- “Feature detector”, “Cat neuron fires”

2	1	2	0
1	0	2	0
0	4	2	3
0	5	1	0

Max pooled

2	2
5	

Pooling

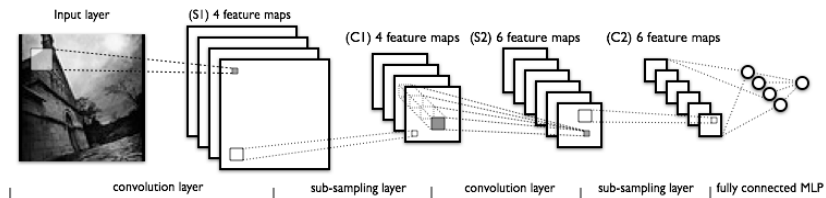
- Often applied between convolution steps
- Divide feature map into “grid”
- Combine vectors inside the same grid cell with some operator
- Most popular: Average pooling, Max pooling
- Max pooling: only select maximum value for each dimension
- “Feature detector”, “Cat neuron fires”

2	1	2	0
1	0	2	0
0	4	2	3
0	5	1	0

Max pooled

2	2
5	3

Convolution and Pooling: LeNet



LeCun et al. (1998). Gradient-based learning applied to document recognition.

Outline

- 1 Convolutional Neural Networks
 - Convolution with one filter
 - Convolution with N filters
 - What does convolution do?
- 2 Pooling
- 3 Application to NLP
- 4 Comparison: RNN vs. CNN

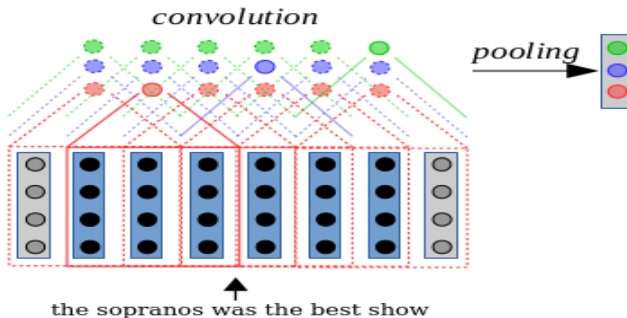
Convolution for NLP

- Images have width and height, but text only has “width” (length)
- → We can discard the “height” dimension from our filters
- Tensor sizes (in NLP):
 - ▶ **Input 2D**: sentence length \times # channels (word embedding size, or # filters of previous convolution)
 - ▶ **Parameter tensor 3D**: filter length \times # channels \times #filters
 - ▶ **Output 2D**: sentence length* \times #filters
 - ▶ *length slightly reduced unless we do padding
- Computer vision: 2D convolution (over height and width)
- NLP: 1D convolution (over length)
- Typically fewer convolutional layers than Computer Vision

Pooling for NLP

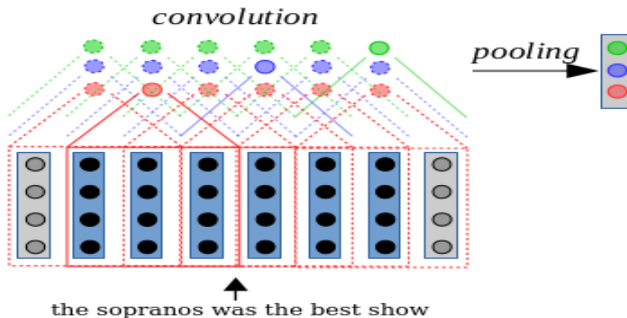
- Pooling between convolutional layers less frequently used than in Computer Vision
- After last convolutional layer: “global” pooling step
- Calculate max/average over the entire sequence (“pooling over time”)

Convolution and Pooling for NLP



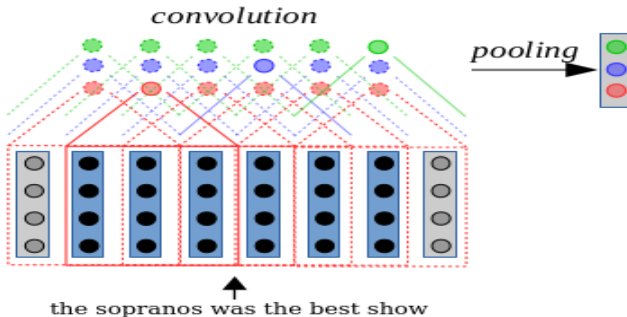
- What is the unpadded input size (=length)?

Convolution and Pooling for NLP



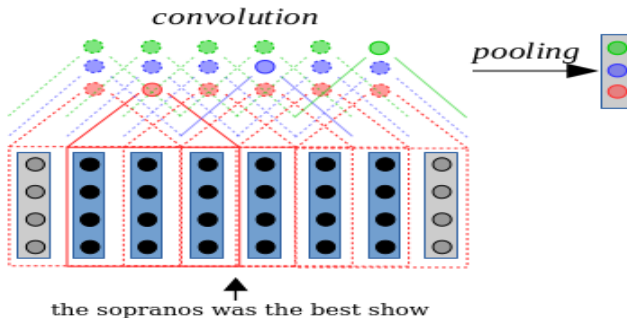
- ▶ What is the unpadded input size (=length)? 6
- ▶ What is the padded input size?

Convolution and Pooling for NLP



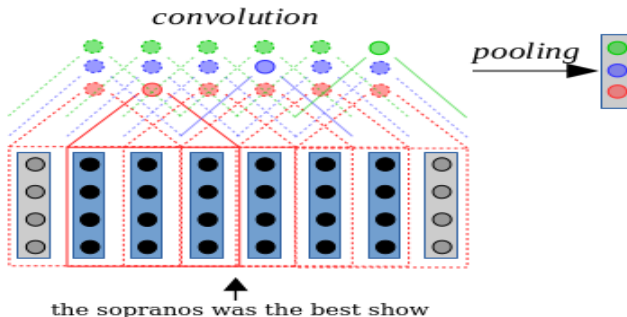
- ▶ What is the unpadded input size (=length)? 6
- ▶ What is the padded input size? 8
- ▶ How many filters?

Convolution and Pooling for NLP



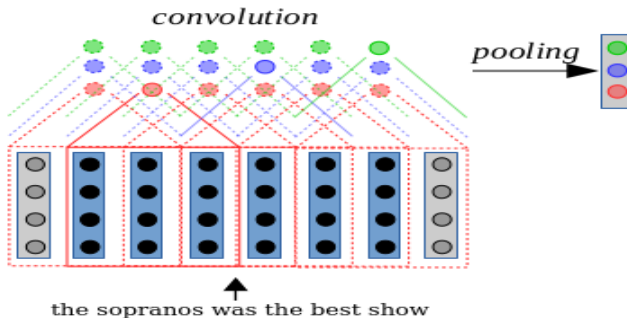
- ▶ What is the unpadded input size (=length)? 6
- ▶ What is the padded input size? 8
- ▶ How many filters? 3
- ▶ How many input channels (=word vector dimensions)?

Convolution and Pooling for NLP



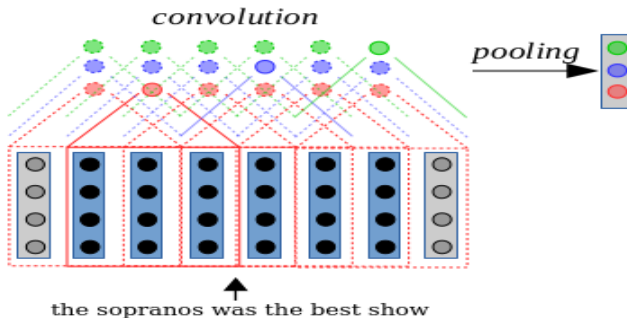
- ▶ What is the unpadded input size (=length)? 6
- ▶ What is the padded input size? 8
- ▶ How many filters? 3
- ▶ How many input channels (=word vector dimensions)? 4
- ▶ What is the filter size (=filter width)?

Convolution and Pooling for NLP



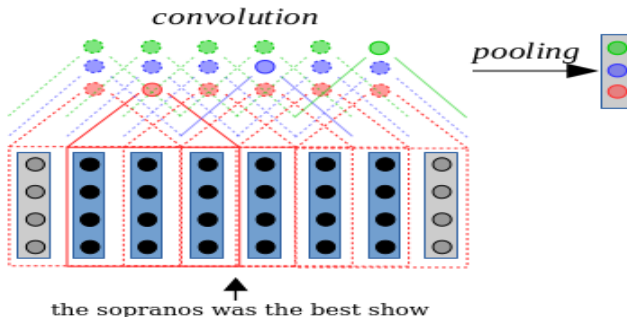
- ▶ What is the unpadded input size (=length)? 6
- ▶ What is the padded input size? 8
- ▶ How many filters? 3
- ▶ How many input channels (=word vector dimensions)? 4
- ▶ What is the filter size (=filter width)? 3
- ▶ What stride (=step size)?

Convolution and Pooling for NLP



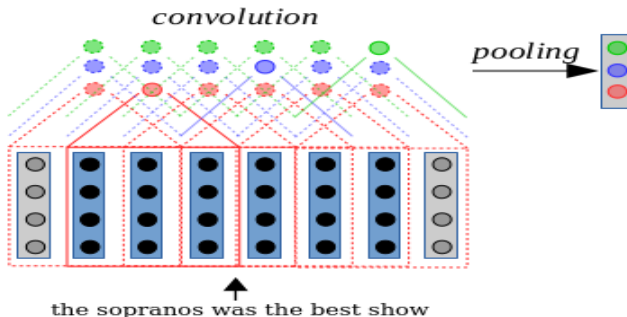
- ▶ What is the unpadded input size (=length)? 6
- ▶ What is the padded input size? 8
- ▶ How many filters? 3
- ▶ How many input channels (=word vector dimensions)? 4
- ▶ What is the filter size (=filter width)? 3
- ▶ What stride (=step size)? 1
- ▶ What is the output size of the convolution operation?

Convolution and Pooling for NLP



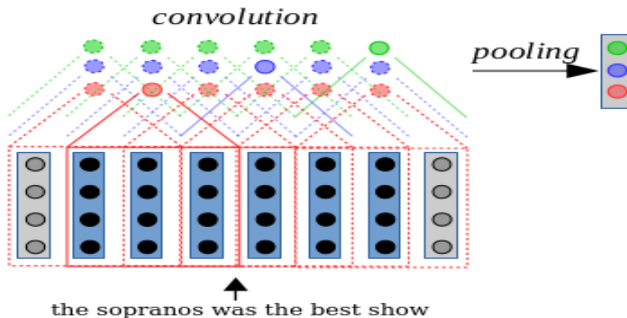
- ▶ What is the unpadded input size (=length)? 6
- ▶ What is the padded input size? 8
- ▶ How many filters? 3
- ▶ How many input channels (=word vector dimensions)? 4
- ▶ What is the filter size (=filter width)? 3
- ▶ What stride (=step size)? 1
- ▶ What is the output size of the convolution operation? 6×3
- ▶ What is the output size of the pooling operation?

Convolution and Pooling for NLP



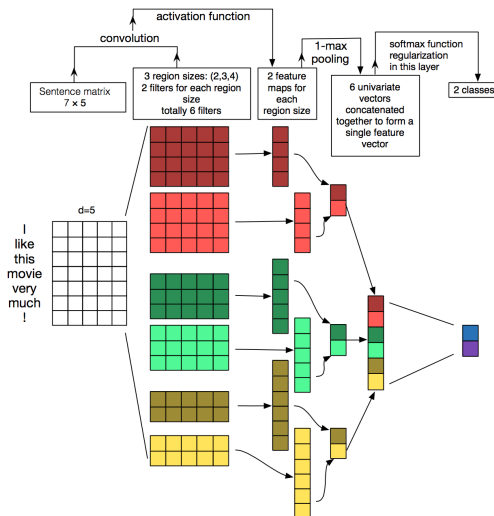
- ▶ What is the unpadded input size (=length)? 6
- ▶ What is the padded input size? 8
- ▶ How many filters? 3
- ▶ How many input channels (=word vector dimensions)? 4
- ▶ What is the filter size (=filter width)? 3
- ▶ What stride (=step size)? 1
- ▶ What is the output size of the convolution operation? 6×3
- ▶ What is the output size of the pooling operation? 3
- ▶ How many parameters have to be learned?

Convolution and Pooling for NLP



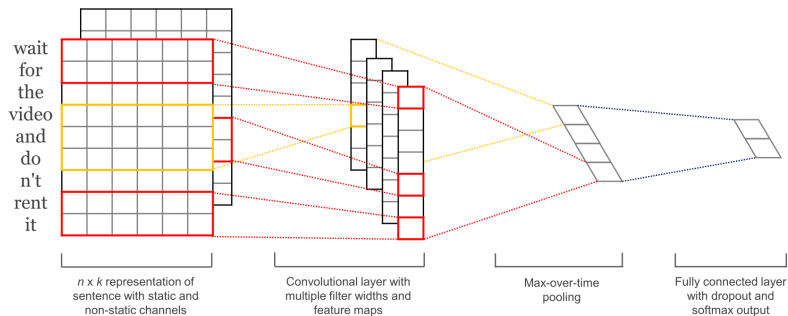
- ▶ What is the unpadded input size (=length)? 6
- ▶ What is the padded input size? 8
- ▶ How many filters? 3
- ▶ How many input channels (=word vector dimensions)? 4
- ▶ What is the filter size (=filter width)? 3
- ▶ What stride (=step size)? 1
- ▶ What is the output size of the convolution operation? 6×3
- ▶ What is the output size of the pooling operation? 3
- ▶ How many parameters have to be learned? $3 \times 3 \times 4 = 36$

Convolution and Pooling for NLP



Source: Zhang, Y., & Wallace, B. (2015). A Sensitivity Analysis of ConvNets for Sentence Classification.

Convolution and Pooling for NLP



Source: Kim, Y. (2014). Convolutional Neural Networks for Sentence Classification.

binary classifier, e.g., sentiment polarity

```
from keras.layers import Embedding, Conv1D, GlobalMaxPooling1D, Dense
from keras.models import Sequential
```

```
embedding = Embedding(input_dim = VOCAB_SIZE, output_dim = EMB_DIM)
conv_layer = Conv1D(filters = NUM_FILTERS, kernel_size = FILTER_WIDTH,
                    activation = "relu")
pool_layer = GlobalMaxPooling1D()
dense_layer = Dense(units = 1, activation = "sigmoid")
```

```
model = Sequential(layers = [emb_layer, conv_layer, pool_layer, dense_layer])
model.compile(loss = "binary_crossentropy", optimizer = "sgd")
X, Y = # load_data()
model.fit(X, Y)
```

Outline

- 1 Convolutional Neural Networks
 - Convolution with one filter
 - Convolution with N filters
 - What does convolution do?
- 2 Pooling
- 3 Application to NLP
- 4 Comparison: RNN vs. CNN

RNN vs. CNN

- Range
 - ▶ CNN: Cannot capture dependencies with range above $k \times L$ (where k is filter width and L is the number of layers)
 - ▶ RNN: Can capture long-range dependencies
- Information transport
 - ▶ RNN: Must learn to “transport” salient information across many time steps.
 - ▶ CNN: No information transport across time, salient information “fast-tracked” by global max pooling
- Efficiency
 - ▶ RNN: Sequential data processing → not parallelizable over time
 - ▶ CNN: Input windows are independent from one another → highly parallelizable over time

RNN vs. CNN: Quiz

- Given a task description, choose appropriate architecture!
 - ▶ Task: predict the number of the main verb (sleep or sleeps)
 - ★ *The cats, who were sitting on the map inside the house, [sleep/sleeps?]*
 - ▶ Which architecture should we use?

RNN vs. CNN: Quiz

- Given a task description, choose appropriate architecture!
 - ▶ Task: predict the number of the main verb (sleep or sleeps)
 - ★ *The cats, who were sitting on the map inside the house, [sleep/sleeps?]*
 - ▶ Which architecture should we use? RNN

RNN vs. CNN: Quiz

- Given a task description, choose appropriate architecture!
 - ▶ Task: predict the number of the main verb (sleep or sleeps)
 - ★ *The cats, who were sitting on the map inside the house, [sleep/sleeps?]*
 - ▶ Which architecture should we use? RNN
 - ▶ Task: predict the polarity of the review:
 - ★ *[... many useless sentences ...] best book ever [... many useless sentences ...]*
 - ▶ Which architecture should we use?

RNN vs. CNN: Quiz

- Given a task description, choose appropriate architecture!
 - ▶ Task: predict the number of the main verb (sleep or sleeps)
 - ★ *The cats, who were sitting on the map inside the house, [sleep/sleeps?]*
 - ▶ Which architecture should we use? RNN
 - ▶ Task: predict the polarity of the review:
 - ★ *[... many useless sentences ...] best book ever [... many useless sentences ...]*
 - ▶ Which architecture should we use? CNN

RNN vs. CNN: Quiz

- Given a task description, choose appropriate architecture!
 - ▶ Task: predict the number of the main verb (sleep or sleeps)
 - ★ *The cats, who were sitting on the map inside the house, [sleep/sleeps?]*
 - ▶ Which architecture should we use? RNN
 - ▶ Task: predict the polarity of the review:
 - ★ *[... many useless sentences ...] best book ever [... many useless sentences ...]*
 - ▶ Which architecture should we use? CNN
- Task: Machine Translation
- Which architecture should we use?

RNN vs. CNN: Quiz

- Given a task description, choose appropriate architecture!
 - ▶ Task: predict the number of the main verb (sleep or sleeps)
 - ★ *The cats, who were sitting on the map inside the house, [sleep/sleeps?]*
 - ▶ Which architecture should we use? RNN
 - ▶ Task: predict the polarity of the review:
 - ★ *[... many useless sentences ...] best book ever [... many useless sentences ...]*
 - ▶ Which architecture should we use? CNN
- Task: Machine Translation
- Which architecture should we use?
 - ▶ Intuitively RNN (because MT is all about long-range dependencies), but ...

RNN vs. CNN: Quiz

- Given a task description, choose appropriate architecture!
 - ▶ Task: predict the number of the main verb (sleep or sleeps)
 - ★ *The cats, who were sitting on the map inside the house, [sleep/sleeps?]*
 - ▶ Which architecture should we use? RNN
 - ▶ Task: predict the polarity of the review:
 - ★ *[... many useless sentences ...] best book ever [... many useless sentences ...]*
 - ▶ Which architecture should we use? CNN
- Task: Machine Translation
- Which architecture should we use?
 - ▶ Intuitively RNN (because MT is all about long-range dependencies), but ...
 - ▶ **Attention** gives CNNs the ability to capture long-range dependencies, while maintaining parallel processing (Gehring et al.)
 - ▶ More about attention: Next week

Next week?

- Attention
- Keras advanced features
- Adversarial training
- Generative pre-training (a.k.a. Elmo & Bert)