

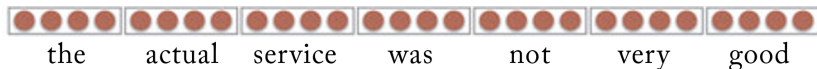
– IN5550 –
Neural Methods in Natural Language Processing
CNNs, Part 2: Convolutions

Erik Velldal

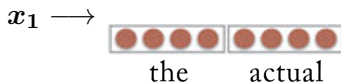
Language Technology Group (LTG)
University of Oslo



- ▶ Consider a sequence of words $w_{1:n} = w_1, \dots, w_n$.
- ▶ Each word is represented by a d dimensional embedding $E_{[w_i]} = \mathbf{w}_i$.



- ▶ A **convolution** corresponds to 'sliding' a **window of size k** across the sequence and applying a **filter** to each.
- ▶ Let $\oplus(\mathbf{w}_{i:i+k-1}) = [\mathbf{w}_i; \mathbf{w}_{i+1}; \dots; \mathbf{w}_{i+k-1}]$ be the concatenation of the embeddings $\mathbf{w}_i, \dots, \mathbf{w}_{i+k-1}$.
- ▶ The vector for the i th window is $\mathbf{x}_i = \oplus(\mathbf{w}_{i:i+k-1})$, where $\mathbf{x}_i \in \mathbb{R}^{kd}$.





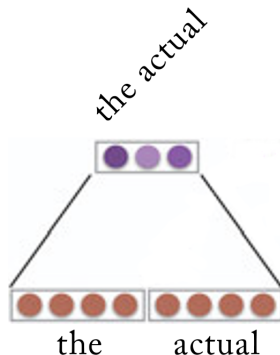
To apply a filter to a window x_i :

- ▶ compute its **dot-product** with a weight vector $u \in \mathbb{R}^{kd}$
- ▶ and then apply a **non-linear activation** g ,
- ▶ resulting in a **scalar** value $p_i = g(x_i \cdot u)$

To apply a filter to a window x_i :

- ▶ compute its **dot-product** with a weight vector $u \in \mathbb{R}^{kd}$
- ▶ and then apply a **non-linear activation** g ,
- ▶ resulting in a **scalar** value $p_i = g(x_i \cdot u)$

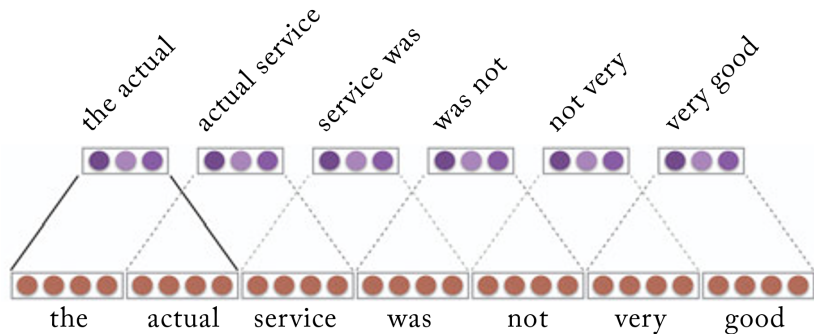
- ▶ Typically use ℓ different filters, u_1, \dots, u_ℓ .
- ▶ Can be arranged in a matrix $U \in \mathbb{R}^{kd \times \ell}$.
- ▶ Also include a bias vector $b \in \mathbb{R}^\ell$.
- ▶ Gives an **ℓ -dimensional vector** p_i summarizing the i th window: $p_i = g(x_i \cdot U + b)$
- ▶ Ideally different dimensions captures different indicative information.



Convolutions on sequences



- ▶ Applying the convolutions over the text results in m vectors $p_{1:m}$.
- ▶ Each $p_i \in \mathbb{R}^{\ell}$ represents a particular k -gram in the input.
- ▶ Sensitive to the identity and order of tokens within the sub-sequence,
- ▶ but independent of its particular position within the sequence.





- ▶ What is m in $p_{1:m}$?
- ▶ For a given window size k and a sequence w_1, \dots, w_n , how many vectors p_i will be extracted?
- ▶ There are $m = n - k + 1$ possible positions for the window.
- ▶ This is called a **narrow convolution**.

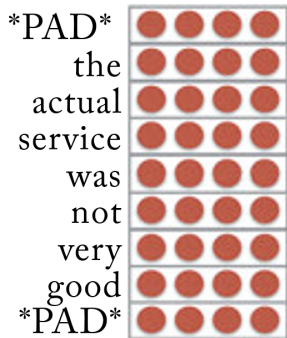
Narrow vs. wide convolutions



- ▶ What is m in $p_{1:m}$?
 - ▶ For a given window size k and a sequence w_1, \dots, w_n , how many vectors p_i will be extracted?
 - ▶ There are $m = n - k + 1$ possible positions for the window.
 - ▶ This is called a **narrow convolution**.
-
- ▶ Another strategy: pad with $k - 1$ extra dummy-tokens on each side.
 - ▶ Let's us slide the window beyond the boundaries of the sequence.
 - ▶ We then get $m = n + k - 1$ vectors p_i .
 - ▶ Called a **wide convolution**.
 - ▶ Necessary when using window-sizes that might be wider than the input.

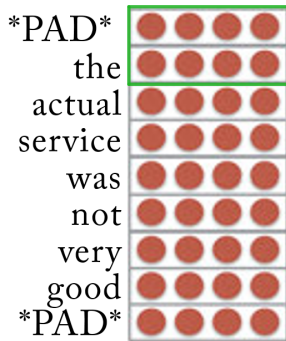


- ▶ So far we've visualized inputs, filters, and filter outputs as sequences:
- ▶ What Goldberg (2017) calls the 'concatenation notation'.



- ▶ An alternative (and perhaps more common) view: 'stacking notation'.
- ▶ Imagine the n input embeddings stacked on top of each other, resulting in an $n \times d$ sentence matrix.

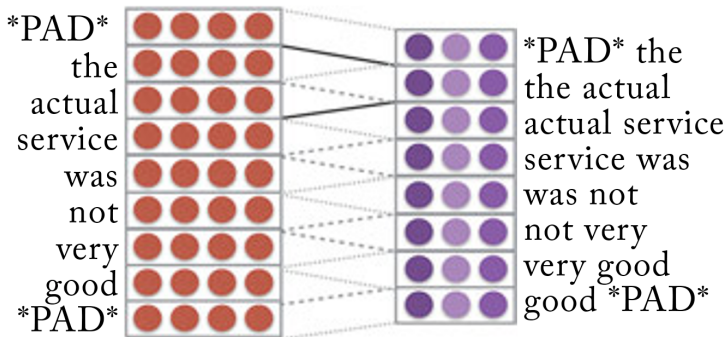
Stacking view (2:3)



- ▶ Correspondingly, imagine each column u in the matrix $U \in \mathbb{R}^{kd \times \ell}$ be arranged as a $k \times d$ matrix.
- ▶ We can then slide ℓ different $k \times d$ filter matrices down the sentence matrix, computing **matrix convolutions**:
- ▶ Sum of element-wise multiplications.

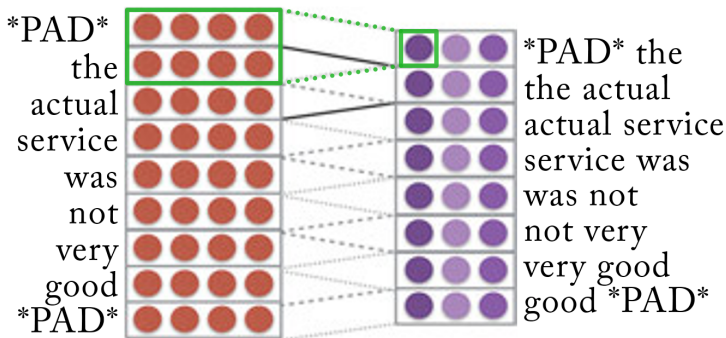
Stacking view (3:3)

- ▶ Now imagine the output vectors $p_{1:m}$ stacked in a matrix $P \in \mathbb{R}^{m \times \ell}$.
- ▶ Each ℓ -dimensional **row** of P holds the features extracted for a given k -gram by different filters.
- ▶ Each m -dimensional **column** of P holds the features extracted across the sequence for a given filter.
- ▶ These columns are sometimes referred to as **feature maps**.



Stacking view (3:3)

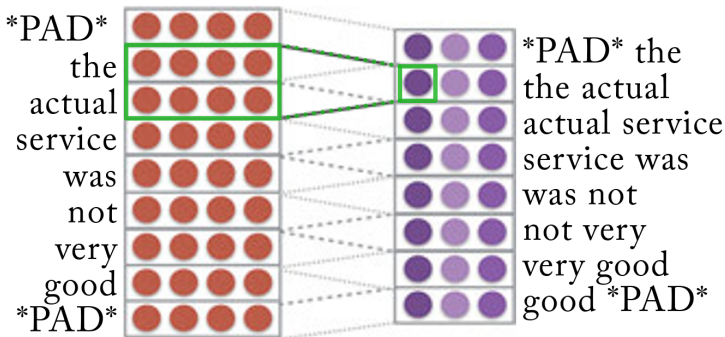
- ▶ Now imagine the output vectors $p_{1:m}$ stacked in a matrix $P \in \mathbb{R}^{m \times \ell}$.
- ▶ Each ℓ -dimensional **row** of P holds the features extracted for a given k -gram by different filters.
- ▶ Each m -dimensional **column** of P holds the features extracted across the sequence for a given filter.
- ▶ These columns are sometimes referred to as **feature maps**.



Stacking view (3:3)



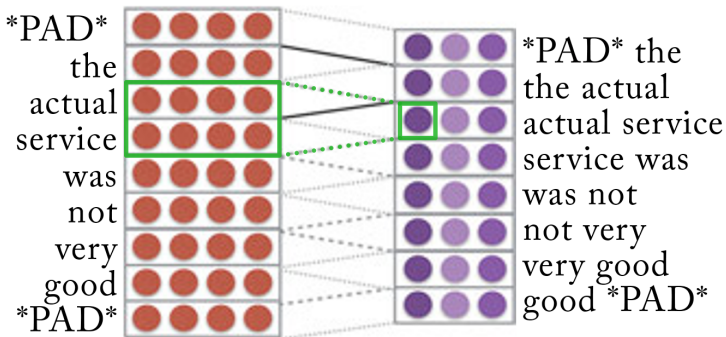
- ▶ Now imagine the output vectors $p_{1:m}$ stacked in a matrix $P \in \mathbb{R}^{m \times \ell}$.
- ▶ Each ℓ -dimensional **row** of P holds the features extracted for a given k -gram by different filters.
- ▶ Each m -dimensional **column** of P holds the features extracted across the sequence for a given filter.
- ▶ These columns are sometimes referred to as **feature maps**.



Stacking view (3:3)

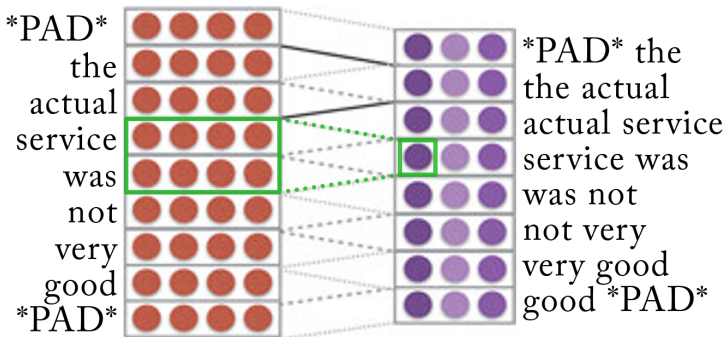


- ▶ Now imagine the output vectors $p_{1:m}$ stacked in a matrix $P \in \mathbb{R}^{m \times \ell}$.
- ▶ Each ℓ -dimensional **row** of P holds the features extracted for a given k -gram by different filters.
- ▶ Each m -dimensional **column** of P holds the features extracted across the sequence for a given filter.
- ▶ These columns are sometimes referred to as **feature maps**.



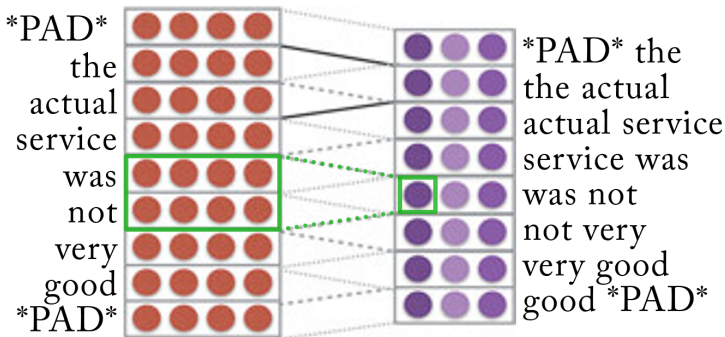
Stacking view (3:3)

- ▶ Now imagine the output vectors $p_{1:m}$ stacked in a matrix $P \in \mathbb{R}^{m \times \ell}$.
- ▶ Each ℓ -dimensional **row** of P holds the features extracted for a given k -gram by different filters.
- ▶ Each m -dimensional **column** of P holds the features extracted across the sequence for a given filter.
- ▶ These columns are sometimes referred to as **feature maps**.



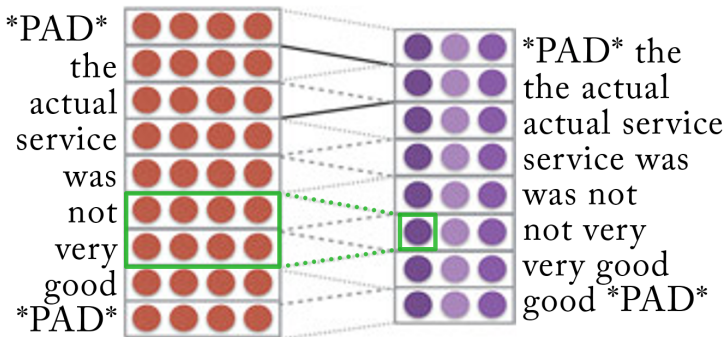
Stacking view (3:3)

- ▶ Now imagine the output vectors $p_{1:m}$ stacked in a matrix $P \in \mathbb{R}^{m \times \ell}$.
- ▶ Each ℓ -dimensional **row** of P holds the features extracted for a given k -gram by different filters.
- ▶ Each m -dimensional **column** of P holds the features extracted across the sequence for a given filter.
- ▶ These columns are sometimes referred to as **feature maps**.



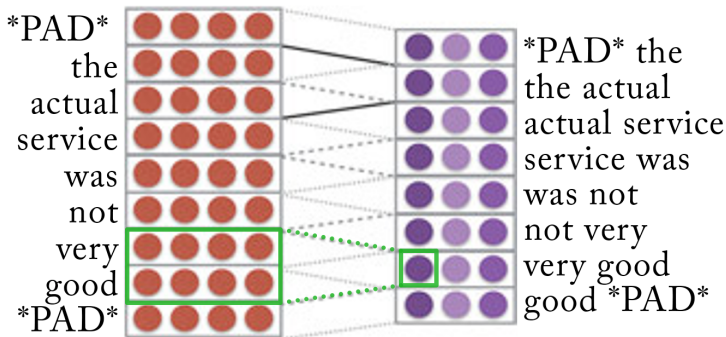
Stacking view (3:3)

- ▶ Now imagine the output vectors $p_{1:m}$ stacked in a matrix $P \in \mathbb{R}^{m \times \ell}$.
- ▶ Each ℓ -dimensional **row** of P holds the features extracted for a given k -gram by different filters.
- ▶ Each m -dimensional **column** of P holds the features extracted across the sequence for a given filter.
- ▶ These columns are sometimes referred to as **feature maps**.



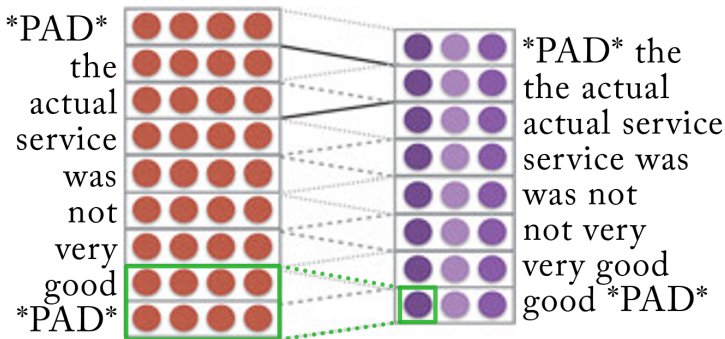
Stacking view (3:3)

- ▶ Now imagine the output vectors $p_{1:m}$ stacked in a matrix $P \in \mathbb{R}^{m \times \ell}$.
- ▶ Each ℓ -dimensional **row** of P holds the features extracted for a given k -gram by different filters.
- ▶ Each m -dimensional **column** of P holds the features extracted across the sequence for a given filter.
- ▶ These columns are sometimes referred to as **feature maps**.



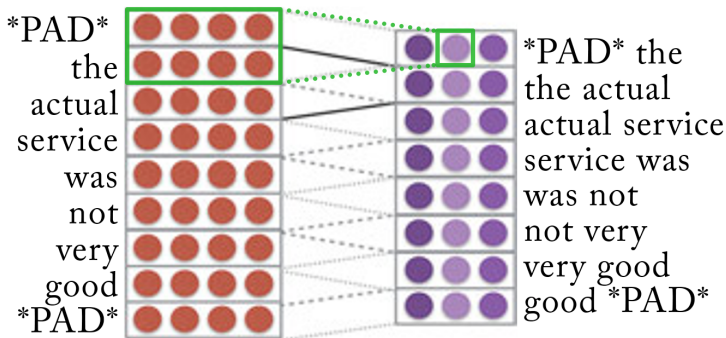
Stacking view (3:3)

- ▶ Now imagine the output vectors $p_{1:m}$ stacked in a matrix $P \in \mathbb{R}^{m \times \ell}$.
- ▶ Each ℓ -dimensional **row** of P holds the features extracted for a given k -gram by different filters.
- ▶ Each m -dimensional **column** of P holds the features extracted across the sequence for a given filter.
- ▶ These columns are sometimes referred to as **feature maps**.



Stacking view (3:3)

- ▶ Now imagine the output vectors $p_{1:m}$ stacked in a matrix $P \in \mathbb{R}^{m \times \ell}$.
- ▶ Each ℓ -dimensional **row** of P holds the features extracted for a given k -gram by different filters.
- ▶ Each m -dimensional **column** of P holds the features extracted across the sequence for a given filter.
- ▶ These columns are sometimes referred to as **feature maps**.





- ▶ Next, in Part 3 of the CNN lecture we cover *pooling*.