

Fully Connected/Dense

A Fully Connected/Dense layer is insensitive to the order of features.

This is just a property of the dot product

$$\Theta^T \cdot \mathbf{x} = \Theta[\text{perm}]^T \cdot \mathbf{x}[\text{perm}]$$

where $\Theta[\text{perm}]^T$ and $\mathbf{x}[\text{perm}]$ are permutations of Θ, \mathbf{x} .

$$\begin{aligned} \sum \begin{Bmatrix} \text{Machine} & \text{Learning} & \text{is} & \text{easy} & \text{not} & \text{hard} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \Theta_{\text{Machine}} & \Theta_{\text{Learning}} & \Theta_{\text{is}} & \Theta_{\text{easy}} & \Theta_{\text{not}} & \Theta_{\text{hard}} \end{Bmatrix} \\ = \\ \sum \begin{Bmatrix} \text{Machine} & \text{Learning} & \text{is} & \text{hard} & \text{not} & \text{easy} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \Theta_{\text{Machine}} & \Theta_{\text{Learning}} & \Theta_{\text{is}} & \Theta_{\text{hard}} & \Theta_{\text{not}} & \Theta_{\text{easy}} \end{Bmatrix} \end{aligned}$$

Convolution

We have spoken about convolutions as

- Identifying the presence/absence of a feature
- At a *spatial* location

The one-dimensional convolution, when applied to a sequence of tokens

- Identifies the presence/absence of a feature
- At a *temporal* location (index within the sequence)

One dimensional convolution
Slide blue kernel over input

Using one dimensional convolution with kernel size n

- The convolution creates an n -gram feature
- At each (temporal) location in the sequence

As with any other CNN, we can apply multiple kernels

- Each matching a different pattern
- To identify a different feature (n-gram)
- At each location in the sequence

One dimensional convolution multiple kernels

Recurrent

Recurrent layers take *sequences* of vectors as input

$$\mathbf{x}_0, \mathbf{x}_1, \dots \mathbf{x}_t \dots \mathbf{x}_T$$

RNN

$\mathbf{h}_{(t)}$ is a **fixed length** vector that "summarizes" the prefix of sequence \mathbf{x} up to element t .

The sequence is processed element by element, so order matters.

$$\mathbf{h}_{(0)} = \text{summary}([\text{Machine}])$$

$$\mathbf{h}_{(1)} = \text{summary}([\text{Machine}, \text{Learning}])$$

$$\vdots$$

$$\mathbf{h}_{(t)} = \text{summary}([\mathbf{x}_{(0)}, \dots, \mathbf{x}_{(t)}])$$

$$\vdots$$

$$\mathbf{h}_{(5)} = \text{summary}([\text{Machine}, \text{Learning}, \text{is}, \text{easy}, \text{not}, \text{hard}])$$

Dropout

- Regularization
- Prevents over-fitting

Normalization

- For use in very deep networks
- Keeps distribution of layer outputs "normalized"