

Fully recurrent network

Activation function or non-linearity

Forward pass of FRNN:

Input weight matrix. It's trans-on between input space R^D & hidden space R^I , i.e. maps input features to hidden repr-ns

$$\vec{s}(t) = \vec{W}^T \vec{x}(t) + \vec{R}^T \vec{a}(t-1)$$

$$\vec{a}(t) = \phi(\vec{s}(t))$$

$$\vec{y}(t) = \vec{V}^T \vec{a}(t)$$

where $\vec{W} \in R^{D \times I}$, $\vec{R} \in R^{I \times I}$, $\vec{V} \in R^{I \times K}$

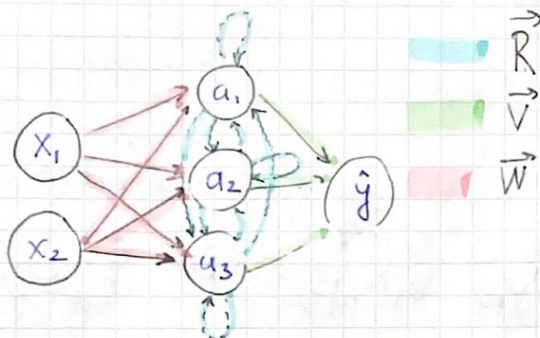
Output activation function & mostly depends on the task & not on the architecture of the network

Recurrent weight matrix. Holds the weights of the loop connections & determines time dependant behaviour of the system.

Output weight matrix. Maps the network's hidden repr-s to the output space R^K

In contrast to the Elman network, we'll now treat the recurrent weights \vec{R} as trainable \Rightarrow weight updates on \vec{R}

Also, in contrast to the Elman networks, a hidden unit may now not only depend on a former version of itself, but also of all neighboring neurons.



Before we look at the derivatives with respect to $r_{ij} = \vec{R}$, $v_{ij} = \vec{V}$, $w_{id} = \vec{W}$ we introduce an ancillary concept: derivatives of the loss with respect to pre-activations, which are called deltas.

Deltas are helpful in backpropagation because they enable a recursive calculation of the gradients through time.

We define:

$$\begin{aligned} \vec{\delta}(t)^T &= \frac{\partial L}{\partial \vec{s}(t)} = \frac{\partial L}{\partial \vec{a}(t)} \frac{\partial \vec{a}(t)}{\partial \vec{s}(t)} = \\ &= \left(\frac{\partial L(\vec{y}(t), \hat{\vec{y}}(t))}{\partial \vec{a}(t)} + \frac{\partial L}{\partial \vec{s}(t+1)} \cdot \frac{\partial \vec{s}(t+1)}{\partial \vec{a}(t)} \right) \frac{\partial \vec{a}(t)}{\partial \vec{s}(t)} = \\ &= \left(\frac{\partial L}{\partial \vec{y}(t)} \cdot \frac{\partial \vec{y}(t)}{\partial \vec{a}(t)} + \frac{\partial L}{\partial \vec{s}(t+1)} \cdot \frac{\partial \vec{s}(t+1)}{\partial \vec{a}(t)} \right) \frac{\partial \vec{a}(t)}{\partial \vec{s}(t)} = \\ &= (\vec{e}(t)^T \text{diag}(\phi'(\vec{V}^T \vec{a}(t))) \vec{V}^T + \vec{\delta}(t+1)^T \vec{R}^T) \text{diag}(\phi'(\vec{W}^T \vec{x}(t))) \end{aligned}$$

derivative of the loss w.r.t. model output

derivative of the model output w.r.t. output weights

Now, let's express derivatives w.r.t. \vec{R} and \vec{W} :

$$\textcircled{I} \quad \frac{\partial L}{\partial \vec{R}} = \sum_{t=1}^T \frac{\partial L}{\partial \vec{s}(t)} \cdot \frac{\partial \vec{s}}{\partial \vec{R}} = \sum_{t=1}^T \vec{\delta}(t) \cdot \frac{\partial (\vec{W}(t)^T \vec{z}(t) + \vec{R}(t)^T \vec{a}(t-1))}{\partial \vec{R}} =$$

$$= \sum_{t=1}^T \vec{\delta}(t) \cdot \vec{a}(t-1)$$

$$\textcircled{II} \quad \frac{\partial L}{\partial \vec{W}} = \sum_{t=1}^T \frac{\partial L}{\partial \vec{s}(t)} \cdot \frac{\partial \vec{s}}{\partial \vec{W}} = \sum_{t=1}^T \vec{\delta}(t) \cdot \frac{\partial (\vec{W}(t)^T \vec{z}(t) + \vec{R}(t)^T \vec{a}(t-1))}{\partial \vec{W}} =$$

$$= \sum_{t=1}^T \vec{\delta}(t) \cdot \vec{z}(t)$$

In order to derive wrt \vec{V} , we have to include additional definition:

$$\vec{\psi}(t)^T = \frac{\partial L}{\partial \vec{V}^T \vec{a}(t)} = \frac{\partial L}{\partial \hat{y}(t)} \cdot \frac{\partial \hat{y}(t)}{\partial \vec{V}^T \vec{a}(t)} = \vec{e}(t)^T \cdot \frac{\partial (\varphi(\vec{V}^T \vec{a}(t)))}{\partial (\vec{V}^T \vec{a}(t))} =$$

$$= \vec{e}(t)^T \circ \varphi'(\vec{a}(t)^T \vec{V})$$

Hence $\vec{a}(t)$ can be multiplied with $\psi(t)$ and we can get $\frac{\partial L}{\partial \vec{V}}$

$$\textcircled{III} \quad \frac{\partial L}{\partial \vec{V}} = \vec{\psi}(t) \vec{a}(t)^T = (\vec{e}(t)^T \circ \varphi'(\vec{a}(t)^T \vec{V})) \vec{a}(t)^T$$