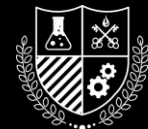


Recurrent Neural Networks



Mass Street
University
PER EDUCATIONEM PROGRESSUS

Shortcomings of Neural Networks

- Neural Networks struggle to deal with spatial and sequential data.
- They have no memory of the past. Only impressions on the network that develop an intuition.
- Fixed inputs and outputs. Weights are not shared and can not store information.



Recurrent Neural Networks

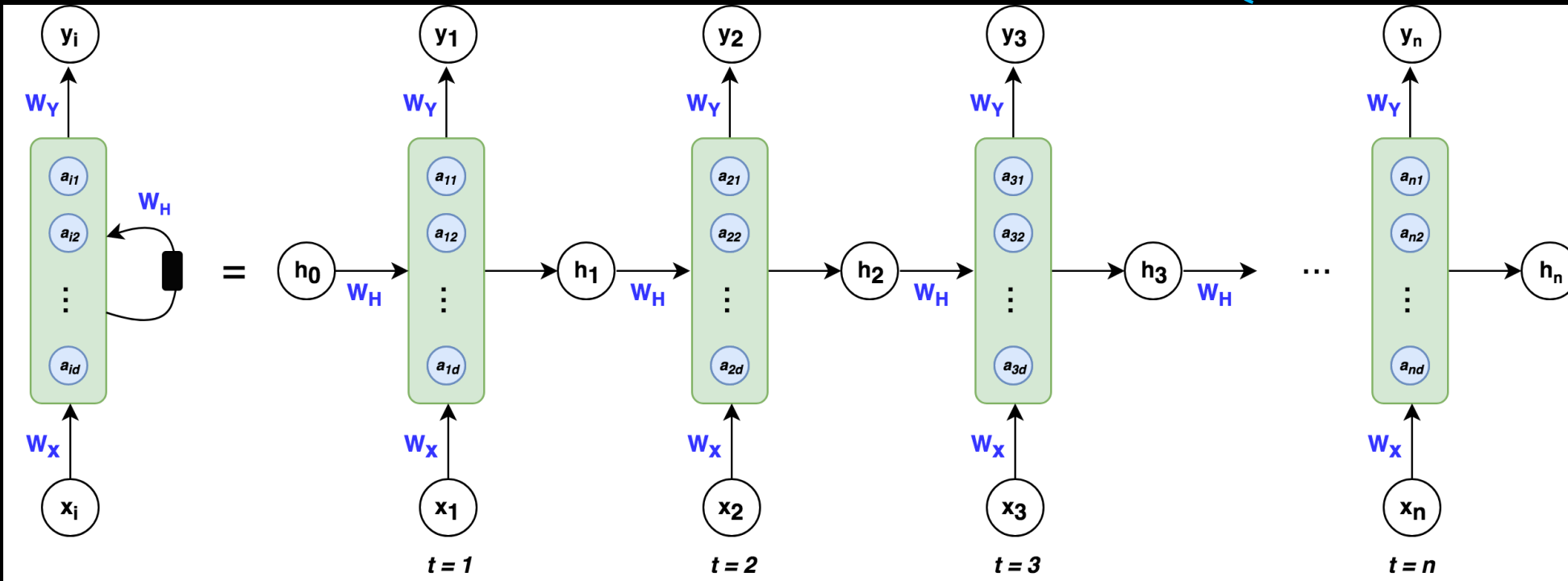
- RNNs are networks that reuse previous weights and information to continuously process information.
- This allows it to process sequential data and work with data that needs long term information, such as text data.



How Recurrent Neural Networks Work

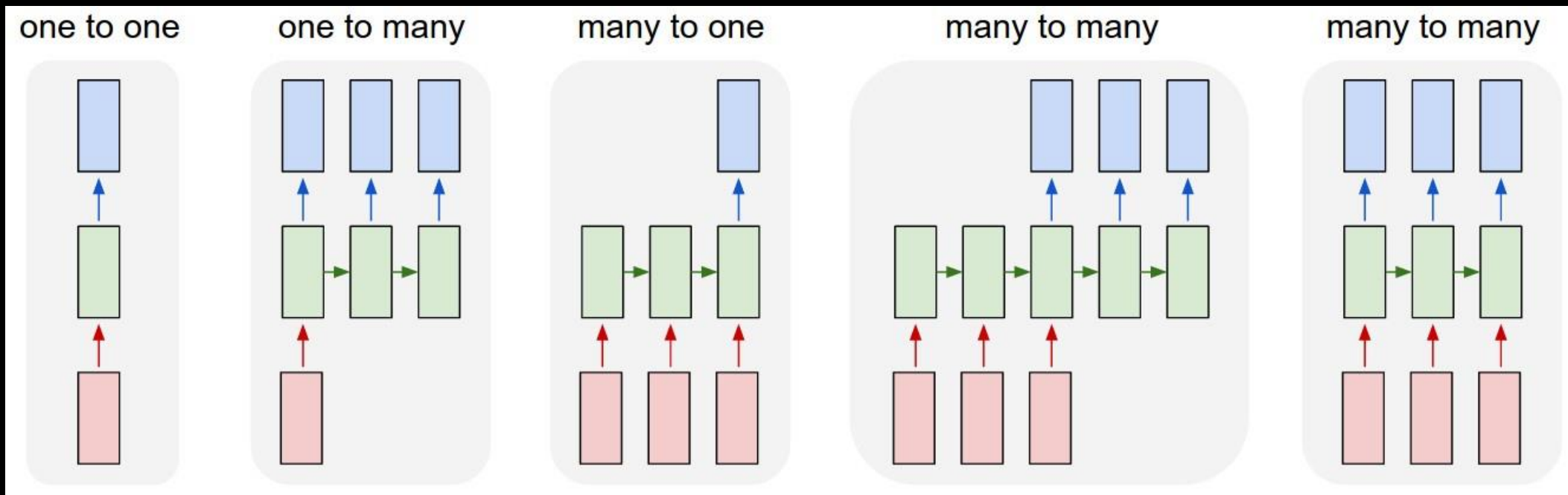
- They have three unique parts. Inputs, Outputs, and State Vectors.
- The state vector and memory loop control what information is stored and the current information that is remembered.

$$H = \tanh(W_h * H + W_x * X)$$



RNN Uses

- Sequence data including One to One, One to Many, Many to One, and Many to Many.



RNN uses

- Sequence data
- One to many operations
 - Single input many outputs
- Many to one operations
 - Single output classification on the sequence
- Many to many operations
 - Synced – an output at every step
 - Un-synced – delay between first input and output



Cont.

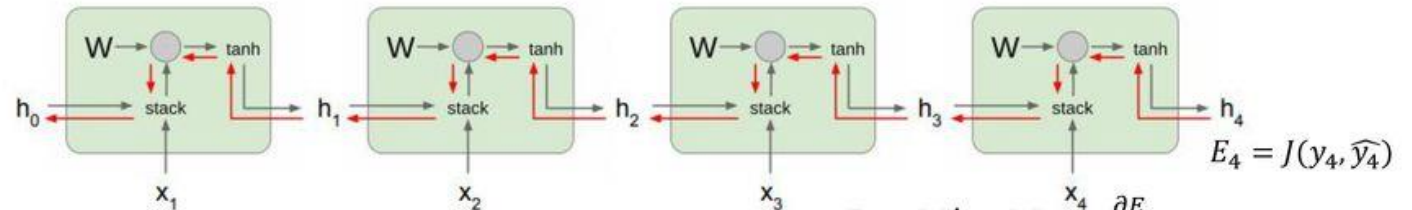
- RNNs can process non-sequential data to keep spatial information intact like a CNN.
- Applications include document generation, text translation, chat bots, law searches for legal firms, lead generation, recommendation machines, and series forecasting.



Limitations

- Exploding and Vanishing Gradients.
- The length of the sequence (amount of memory) we can process.

Exploding and Vanishing Gradients



- Backpropagation algorithms correct parameters with the error E by $W' = W - \eta \frac{\partial E}{\partial W}$.

$$\frac{\partial E_k}{\partial W} = \sum_{i=0}^k \frac{\partial E_k}{\partial \hat{y}_k} \left(\frac{\partial \hat{y}_k}{\partial h_k} \right) \left(\frac{\partial h_k}{\partial W} \right)$$

$$\frac{\partial h_k}{\partial h_i} = \frac{\partial h_k}{\partial h_{k-1}} \frac{\partial h_{k-1}}{\partial h_{k-2}} \dots \frac{\partial h_{i+1}}{\partial h_i}$$

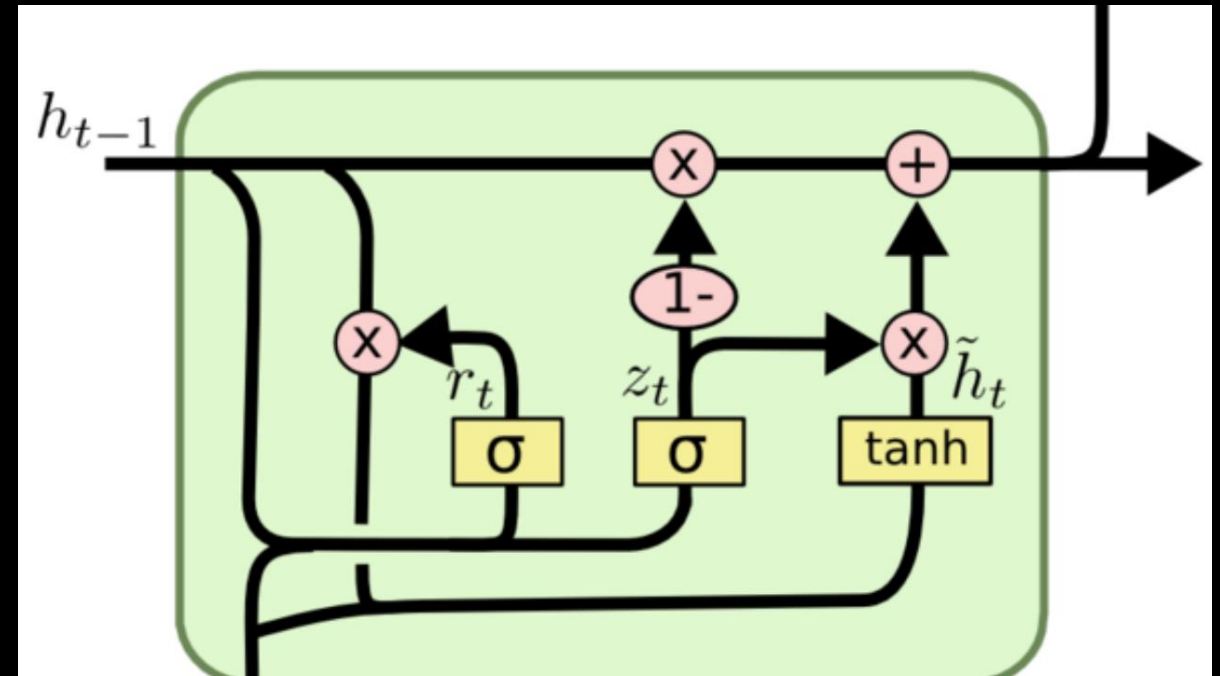
$$h_t = \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

- However, in vanilla RNNs, computing this gradient involves many factors of W_{hh} (and repeated \tanh)*. If we decompose the singular values of the gradient multiplication matrix,
 - Largest singular value $> 1 \rightarrow$ **Exploding gradients**
 - Slight error in the late time steps causes drastic updates in the early time steps \rightarrow Unstable learning
 - Largest singular value $< 1 \rightarrow$ **Vanishing gradients**
 - Gradients passed to the early time steps is close to 0. \rightarrow Uninformed correction

* Refer to Bengio et al. (1994) or Goodfellow et al. (2016) for a complete derivation

Gated Recurrent Networks.

- GRUs use two types of gates: Update and Reset gates. These gates are vectors that are multiplied by the input vectors.
- The update gate (Z_t) keeps information and the reset gate (R_t) forgets information.



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

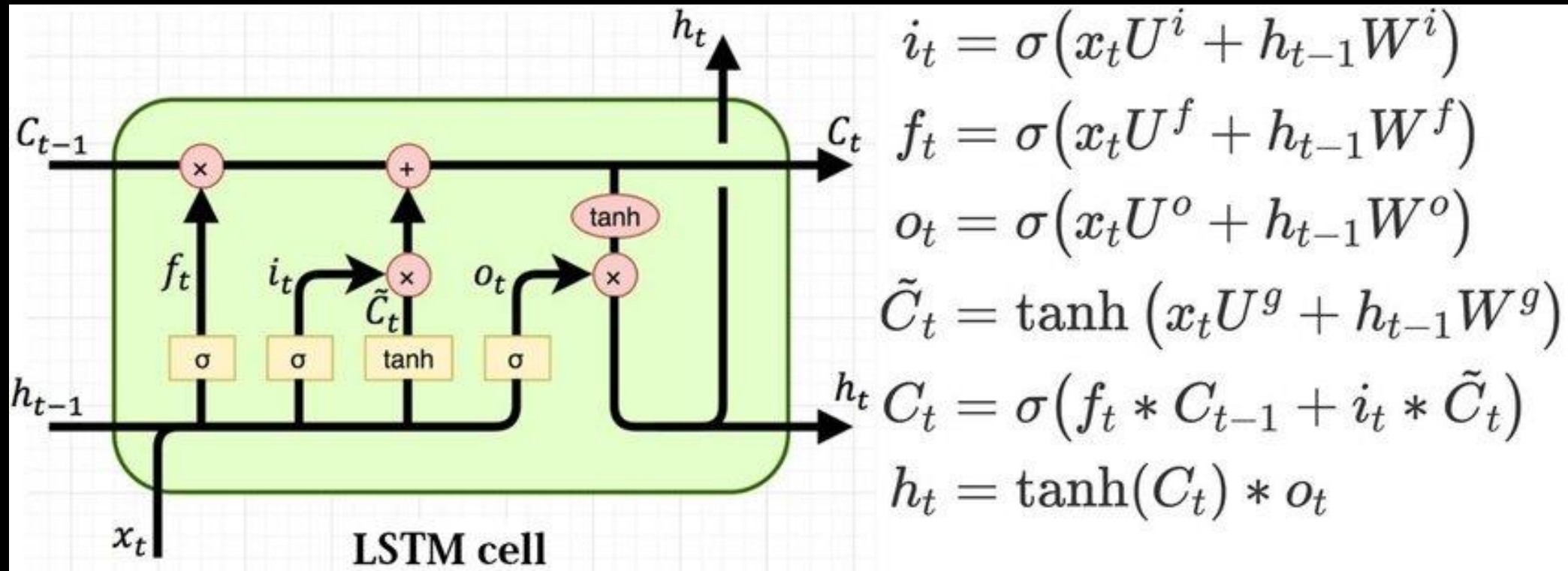


Long Short Term Memory Networks

- LSTMs use three different gates: Forget Gate, Input Gate, Output Gate. It also has a cell state that retains the past values.
- The gates regulate the flow of information in the RNN much like the GRU.



LSTM Cont.



Bi-Directional RNNs

- They process information forward and backwards to better understand context and long-term dependencies.

