

# Recurrent Neural Networks (RNNs) and the Exploding & Vanishing Gradient Problems

Le Gia Khang	21522189
Bui Manh Hung	21522110
Nguyen Hoang Tan	21521413
Pham Tram Anh	21520587
Le Thi Kim Yen	21521695

**University of Information Technology**

December 27, 2022

# References

The contents of this document are taken mainly from the follow sources:

- Robin M.Schmidt, Recurrent Neural Networks (RNNs): a gentle introductionand overview<sup>1</sup>;
- Razvan Pascanu, Tomas Mikolov, Yoshua Bengio, On the difficulty of training RNNs<sup>2</sup>
- Yoshua Bengio, Patrice Simard, Paolo Frasconi, Learning long-term dependencies with gradients descent is difficult<sup>3</sup>

---

<sup>1</sup><https://arxiv.org/pdf/1912.05911.pdf>

<sup>2</sup><https://arxiv.org/pdf/1211.5063.pdf>

<sup>3</sup>[https://www.researchgate.net/profile/Y-Bengio/publication/5583935\\_Learning\\_long-term\\_dependencies\\_with\\_gradient\\_descent\\_is\\_difficult/links/546b702d0cf2397f7831c03d/Learning-long-term-dependencies-with-gradient-descent-is-difficult.pdf](https://www.researchgate.net/profile/Y-Bengio/publication/5583935_Learning_long-term_dependencies_with_gradient_descent_is_difficult/links/546b702d0cf2397f7831c03d/Learning-long-term-dependencies-with-gradient-descent-is-difficult.pdf)

# Table of Contents

- ① Introduction & Notation
- ② Training Recurrent Neural Networks
- ③ Exploding and Vanishing Gradient

# Table of Contents

① Introduction & Notation

② Training Recurrent Neural Networks

③ Exploding and Vanishing Gradient

# Sequential Data

- Sequential Data refers to any data that contain elements that are ordered into sequences.
- Examples include time series, DNA sequences (see biomedical informatics) and sequences of user actions.

## Question ?

**Named-entity Recognition for person names:**

Harry Potter and Hermione Granger invented a new spell.

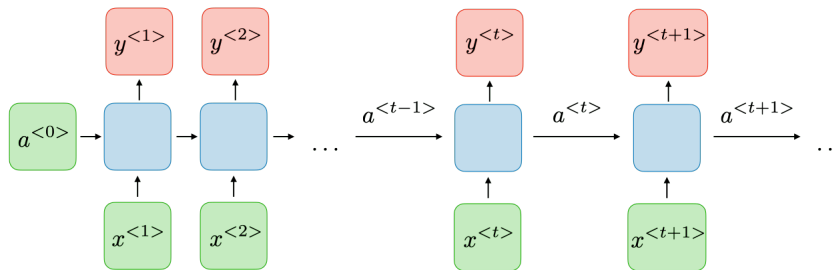
- Machine learning models that input or output data sequences are known as sequence models.

# What's Recurrent Neural Networks (RNNs)

- Recurrent Neural Networks (RNNs) are a type of neural network architecture which is mainly used to detect patterns in a sequence of data.
- However, they are also applicable to images if these get respectively decomposed into a series of patches and treated as a sequence.
- What differentiates Recurrent Neural Networks from Feedforward Neural Networks also known as Multi-Layer Perceptrons (MLPs) is how information gets passed through the network.

# Model Architecture

- The RNN has cycles and transmits information back into itself. This enables them to extend the functionality of Feedforward Networks to also take into account previous inputs  $X_{0:t-1}$  and not only the current input  $X_t$ .



# Model Architecture

- We denote the hidden state and the input at time step  $t$  respectively as  $\mathbf{H}_t$  and  $X_t$ . Further, we use  $\mathbf{W}_{xh}$ ,  $\mathbf{W}_{hh}$ ,  $\mathbf{W}_{ho}$  and a bias  $b_h$  as shared parameters. Lastly, all these informations get passed to a activation function  $\phi$  which is usually a logistic sigmoid or tanh function.

Hidden state at the time step  $t$

$$\mathbf{H}_t = \phi_h(\mathbf{X}_t \mathbf{W}_{xh} + \mathbf{H}_{t-1} \mathbf{W}_{hh} + \mathbf{b}_h)$$

Output at time step  $t$

$$\mathbf{O}_t = \phi_o(\mathbf{H}_t \mathbf{W}_{ho} + \mathbf{b}_o)$$



# Loss Function

- Loss function evaluates performance of the network by comparing the output  $\mathbf{y}_t$  with the corresponding target  $\mathbf{z}_t$  defined as:

## Loss Function

$$L(\mathbf{y}, \mathbf{z}) = \sum_{t=1}^T L_t(\mathbf{y}_t, \mathbf{z}_t)$$

- Selection of the loss function is problem dependent. Some popular loss function are Euclidean distance and Hamming distance and cross-entropy.

## Benefits

- Possibility of processing input of any length;
- Model size not increasing with size of input;
- Computation takes into account historical information;
- Weights are shared across time

## Costs

- Computation being slow;
- Difficulty of accessing information from a long time ago;
- Cannot consider any future input for the current state;

# Various usage of RNN

## One to One RNN

- This type of neural network is known as the Vanilla Neural Network. It's used for general machine learning problems, which has a single input and a single output.

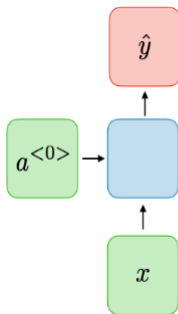


Figure: Vanilla Neural Network

## One to Many RNN

- A type of RNN that gives multiple outputs when given a single input. It takes a fixed input size and gives a sequence of data outputs.
- Its applications can be found in Music Generation and Image Captioning.

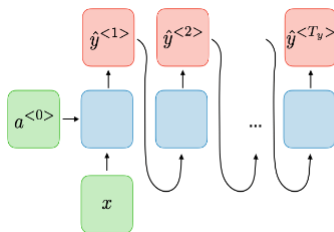


Figure: One-to-Many mode

## Many to One RNN

- Many-to-One is used when a single output is required from multiple input units or a sequence of them. It takes a sequence of inputs to display a fixed output.
- Sentiment Analysis is a common example of this type of Recurrent Neural Network.

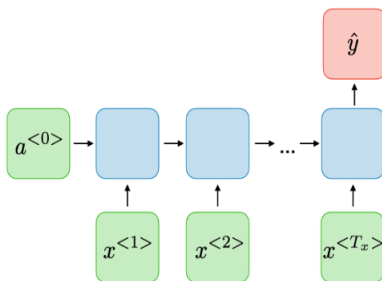
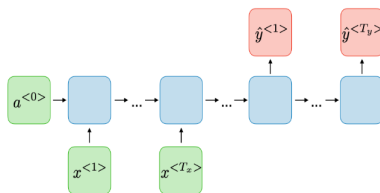
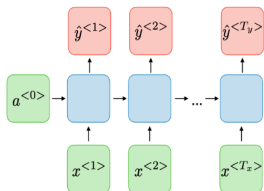


Figure: Many to One mode

# Various usage of RNN

## Many to Many RNN

- Many-to-Many is used to generate a sequence of output data from a sequence of input units.
- This type of RNN is further divided into the following two subcategories:
  - Equal Unit Size:** In this case, the number of both the input and output units is the same. A common application can be found in Name-Entity Recognition.
  - Unequal Unit Size:** In this case, inputs and outputs have different numbers of units. Its application can be found in Machine Translation.



# Table of Contents

- ① Introduction & Notation
- ② Training Recurrent Neural Networks
- ③ Exploding and Vanishing Gradient

- A general rule is to assign small values to the weights. A Gaussian draw with a standard deviation of 0.001 or 0.01 is a reasonable choice.
- The biases are usually set to zero, but the output bias can also be set to a very small value.
- However, the initialization of parameters is dependent on the task and properties of the input data such as dimensionality.



# Gradient-based Learning Methods

- Gradient descent (GD) adjusts the weights of the model by finding the error function derivatives with respect to each member of the weight matrices.

## Batch Gradient Descent

Computing the gradient for the whole dataset in each optimization iteration to perform a single update as

$$\theta_{t+1} = \theta_t - \frac{\lambda}{U} \sum_{k=1}^U \frac{\partial \ell_k}{\partial \theta}$$

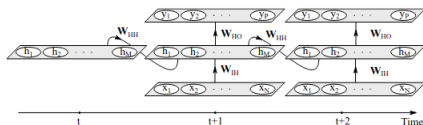
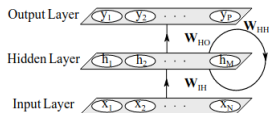
- However, computing error-derivatives through time is difficult. This is mostly due to the relationship among the parameters and the dynamics of the RNN, that is highly unstable and makes GD ineffective.

# Gradient-based Learning Methods

- Gradient-based algorithms have difficulty in capturing dependencies as the duration of dependencies increases.
- The derivatives of the loss function with respect to the weights only consider the current state, without using the history information for weights updating.
- RNNs cannot learn long-range temporal dependencies when GD is used for training.

# Back-propagation through time (BPTT)

- Backpropagation Through Time (BPTT) is the adaption of the backpropagation algorithm for RNNs.
- BPTT unfolds the RNN to construct a traditional Feedforward Neural Network where we can apply backpropagation.



# Back-propagation through time (BPTT)

Since we have three weight matrices  $\mathbf{W}_{xh}$ ,  $\mathbf{W}_{hh}$  and  $\mathbf{W}_{ho}$  we need to compute the partial derivative to each of these weight matrices using the chain rule.

Partial derivative with respect to  $\mathbf{W}_{ho}$

$$\frac{\partial L}{\partial \mathbf{W}_{ho}} = \sum_{t=1}^T \frac{\partial \ell_t}{\partial \mathbf{O}_t} \cdot \frac{\partial \mathbf{O}_t}{\partial \phi_o} \cdot \frac{\partial \phi_o}{\partial \mathbf{W}_{ho}} = \sum_{t=1}^T \frac{\partial \ell_t}{\partial \mathbf{O}_t} \cdot \frac{\partial \mathbf{O}_t}{\partial \phi_o} \cdot \mathbf{H}_t$$

Partial derivative with respect to  $\mathbf{W}_{hh}$

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{W}_{hh}} &= \sum_{t=1}^T \frac{\partial \ell_t}{\partial \mathbf{O}_t} \cdot \frac{\partial \mathbf{O}_t}{\partial \phi_o} \cdot \frac{\partial \phi_o}{\partial \mathbf{H}_t} \cdot \frac{\partial \mathbf{H}_t}{\partial \phi_h} \cdot \frac{\partial \phi_h}{\partial \mathbf{W}_{hh}} \\ &= \sum_{t=1}^T \frac{\partial \ell_t}{\partial \mathbf{O}_t} \cdot \frac{\partial \mathbf{O}_t}{\partial \phi_o} \cdot \mathbf{W}_{ho} \cdot \frac{\partial \mathbf{H}_t}{\partial \phi_h} \cdot \frac{\partial \phi_h}{\partial \mathbf{W}_{hh}} \end{aligned}$$

# Back-propagation through time (BPTT)

- Since each  $\mathbf{H}_t$  depends on the previous time step we can substitute the last part from above equation

Partial derivative with respect to  $\mathbf{W}_{hh}$

$$\frac{\partial L}{\partial \mathbf{W}_{hh}} = \sum_{t=1}^T \frac{\partial \ell_t}{\partial \mathbf{O}_t} \cdot \frac{\partial \mathbf{O}_t}{\partial \phi_o} \cdot \mathbf{W}_{ho} \sum_{k=1}^t \frac{\partial \mathbf{H}_t}{\partial \mathbf{H}_k} \cdot \frac{\partial \mathbf{H}_k}{\partial \mathbf{W}_{hh}}$$

- Similarly, we have the partial derivative with respect to  $\mathbf{W}_{xh}$  as below

Partial derivative with respect to  $\mathbf{W}_{xh}$

$$\frac{\partial L}{\partial \mathbf{W}_{xh}} = \sum_{t=1}^T \frac{\partial \ell_t}{\partial \mathbf{O}_t} \cdot \frac{\partial \mathbf{O}_t}{\partial \phi_o} \cdot \mathbf{W}_{ho} \sum_{k=1}^t \frac{\partial \mathbf{H}_t}{\partial \mathbf{H}_k} \cdot \frac{\partial \mathbf{H}_k}{\partial \mathbf{W}_{xh}}$$

# Back-propagation through time (BPTT)

- In order to transport the error through time from timestep  $t$  back to timestep  $k$  we can have

$$\frac{\partial \mathbf{H}_t}{\partial \mathbf{H}_k} = \prod_{i=k+1}^t \frac{\partial \mathbf{H}_i}{\partial \mathbf{H}_{i-1}}$$

- We can consider previous equation as a Jacobian matrix for the hidden state parameter as

$$\prod_{i=k+1}^t \frac{\partial \mathbf{H}_i}{\partial \mathbf{H}_{i-1}} = \prod_{i=k+1}^t \mathbf{W}_{hh}^T \text{diag}(\phi'_h(\mathbf{H}_{i-1}))$$

# Table of Contents

- ① Introduction & Notation
- ② Training Recurrent Neural Networks
- ③ Exploding and Vanishing Gradient

# Vanishing Gradient

- This problem refers to the exponential shrinking of gradient magnitudes as they are propagated back through time.
- This phenomena causes memory of the network to ignore long term dependencies and hardly learn the correlation between temporally distant events
- There are two reasons for that:
  - ① Standard nonlinear functions have a gradient which is almost everywhere close to zero;
  - ② The magnitude of gradient is multiplied over and over by the recurrent matrix as it is back-propagated through time.



# Vanishing Gradient

- It is possible to use singular values to generalize it to the non-linear function  $\phi'_h$  by bounding it with  $\gamma \in R$  such as

$$\|diag(\phi'_h(\mathbf{H}_k))\| \leq \gamma$$

- The Jacobian matrix  $\frac{\partial \mathbf{H}_{k+1}}{\partial \mathbf{H}_k}$  and the bound in previous equation, we can have

$$\forall k, \left\| \frac{\partial \mathbf{H}_{k+1}}{\partial \mathbf{H}_k} \right\| \leq \|\mathbf{W}_{hh}^T\| \|diag(\sigma'(\mathbf{H}_k))\| < 1$$

# Vanishing Gradient

- We can consider  $\left\| \frac{\partial \mathbf{H}_{k+1}}{\partial \mathbf{H}_k} \right\| \leq \eta < 1$  such as  $\eta \in R$  for each step  $k$ . By continuing it over different timesteps and adding the loss function component we can have

$$\left\| \frac{\partial \ell_t}{\partial \mathbf{H}_t} \left( \prod_{i=k}^{t-1} \frac{\partial \mathbf{H}_{i+1}}{\partial \mathbf{H}_i} \right) \right\| \leq \eta^{t-k} \left\| \frac{\partial \ell_t}{\partial \mathbf{H}_t} \right\|$$

- Finally, we can see that the sufficient condition for the gradient vanishing problem to appear is that the largest singular value of the recurrent weights matrix  $\mathbf{W}_{hh}$  satisfies  $\lambda_1 < \frac{1}{\gamma}$

# Exploding Gradient

- Gradients in training RNNs on long sequences may explode as the weights become larger and the norm of the gradient during training largely increases.
- As it is stated before, the necessary condition for this situation to happen is  $\lambda > \frac{1}{\gamma}$ .

# Conclusion for the 2 problems

## Vanishing Gradient

$$\left\| \frac{\partial \mathbf{H}_{i+i}}{\partial \mathbf{H}_i} \right\| < 1$$

## Exploding Gradient

$$\left\| \frac{\partial \mathbf{H}_{i+1}}{\partial \mathbf{H}_i} \right\| > 1$$

# Dealing with the exploding and vanishing gradient

- To deal with the exploding gradients problem, we propose a solution that involves clipping the norm of the exploded gradients when it is too large.

---

$$\begin{aligned} \hat{\mathbf{g}} &\leftarrow \frac{\partial \mathcal{E}}{\partial \theta} \\ \text{if } \|\hat{\mathbf{g}}\| &\geq \text{threshold} \text{ then} \\ &\hat{\mathbf{g}} \leftarrow \frac{\text{threshold}}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}} \\ \text{end if} \end{aligned}$$

---

- To deal with the vanishing gradients problem, we propose using LSTMs layers rather than normal Recurrent layers.

