

---

# Report for the Deep Learning Course Assignment 2

---

Martin de la Riva  
11403799  
martin7557@gmail.com

## 1 Vanilla RNN

### Question 1.1

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{oh}} &= - \sum_k y_k \frac{\partial \log(\hat{y}_k)}{\partial \mathbf{W}_{oh}} = - \sum_k y_k \frac{\partial \log(\hat{y}_k)}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial \mathbf{W}_{oh}} \\ &= - \sum_k y_k \frac{1}{\hat{y}_k} \frac{\partial \hat{y}_k}{\partial \mathbf{W}_{oh}} = (\hat{y}_k - y_k) \frac{\partial p_k}{\partial \mathbf{W}_{oh}} = (\hat{y}_k - y_k) h^{(t)}\end{aligned}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{hh}} = (\hat{y}_k - y_k) \frac{\partial p_k}{\partial \mathbf{W}_{hh}} = (\hat{y}_k - y_k) \mathbf{W}_{oh} (1 - \tanh^2(h^{(t)})) h^{(t-1)}$$

### Question 1.3

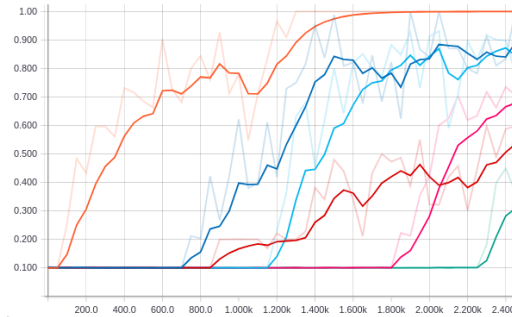


Figure 1: Accuracies for Vanilla RNN. Input lengths from 5 to 30.

### Question 1.4

**RMSProp:** In Adagrad the gradient is monotonically increasing. This could lead the model to stop learning, which could be problematic. RMSProp decay the past accumulated gradient (adaptive learning rate) in an attempt to reduce its aggressively decreasing learning rate.

**Adam:** Adam implements momentum on RMSProp in order not to get stuck in local optima.

**Momentum:** Adding momentum to the update optimizer is used to not stay in local optima and pass through them. It is clear why adding this improves from the classic vanilla SGD.

## 2 Long-Short Term Network (LSTM)

### Question 1.5

**modulation gate:** Tanh layer. Creates a vector of candidate features that could be inserted to the state. If you use only sigmoids, the result values will be between  $[0,1]$ . Then, no value will be zero and will therefore be forgotten. That is why in this gate tanh is used, so it can have values between  $[-1,1]$  and therefore be able to forget memory.

**input gate:** A sigmoid layer that decides which features will be updated in the cell state. It is combined with the modulation gate to create a state update.

**forget gate:** Sigmoid layer in charge of deciding what features's information is going to be not kept from the cell state.

**output gate:** Sigmoid layer. Defines how much of the cell state is exposed to the external network (higher layers and next time step).

The sigmoid gates are used so the value of the gates are between 0 and 1 at each feature; 0 meaning not to take any information from that feature and 1 to completely take it.

$$\mathbf{W}_{\mu x} \in \mathbb{R}^{d \times n} \quad \forall \mu \in g, i, f, o,$$

$$\mathbf{W}_{\mu h} \in \mathbb{R}^{n \times n} \quad \forall \mu \in g, i, f, o,$$

$$\mathbf{W}_{out} \in \mathbb{R}^{d \times n}$$

$$\mathbf{b}_{\mu} \in \mathbb{R}^n \quad \forall \mu \in g, i, f, o,$$

$$\mathbf{b}_{out} \in \mathbb{R}^n$$

Given these definitions, the total number of trainable parameter is:

$$d * n * 5 + n * n * 4 + n * 5$$

### Question 1.6

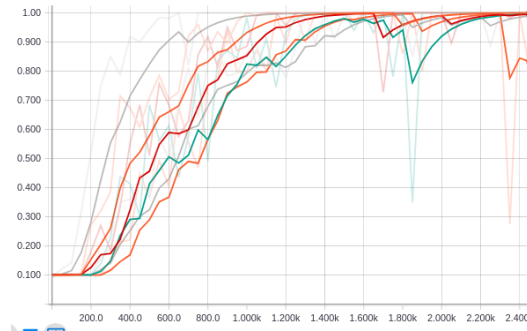


Figure 2: Accuracies for LSTM. Input lengths from 5 to 45.

## 3 Recurrent Nets as Generative Model

### Question 2.1

I trained different models for 20,000 steps on the *democracy in the US*. I inputted the sentence *in the year* and got the following results:

- 1 step: in the yearzz,./zJMMM2222AAAUuuussssggg
- 750 steps: in the yeare of the porest of the porest
- 1250 steps: in the years of the properal the power th
- 5000 steps: in the year the political power of the Un

