

# Assignment 2 - Implementation of Recurrent Perceptron

Jimut Bahan Pal, 22D1594  
Shambhavi Pandey, 23D1145  
Saikat Dutta, 23D2031

IIT Bombay

CS-772

31 March, 2024

# Problem Statement

- **Input:** POS-tagged input tokens
- **Output:**
  - Noun chunk labels on tokens .
  - The beginning of the chunk will be labeled 1 and the rest of the words in the chunk will be labeled 0.
  - All other words are labeled 1.

# Implementation Details

$$s_t = \tanh(Wx_t + W_0s_{t-1} + Vp_{t-1})$$

$$o_t = \text{sigmoid}(s_t)$$

Here  $p_{t-1} = x_{t-1}$  and  $x_t = x_t[1 : l]$

$$x_t \in \mathbb{R}^5 \iff x \in \mathbb{R}^{l \times 5}$$

$$o_t \in \mathbb{R}^5 \iff x \in \mathbb{R}^{l \times 5}$$

$$s_t \in \mathbb{R}^1 \iff x \in \mathbb{R}^{l \times 1}$$

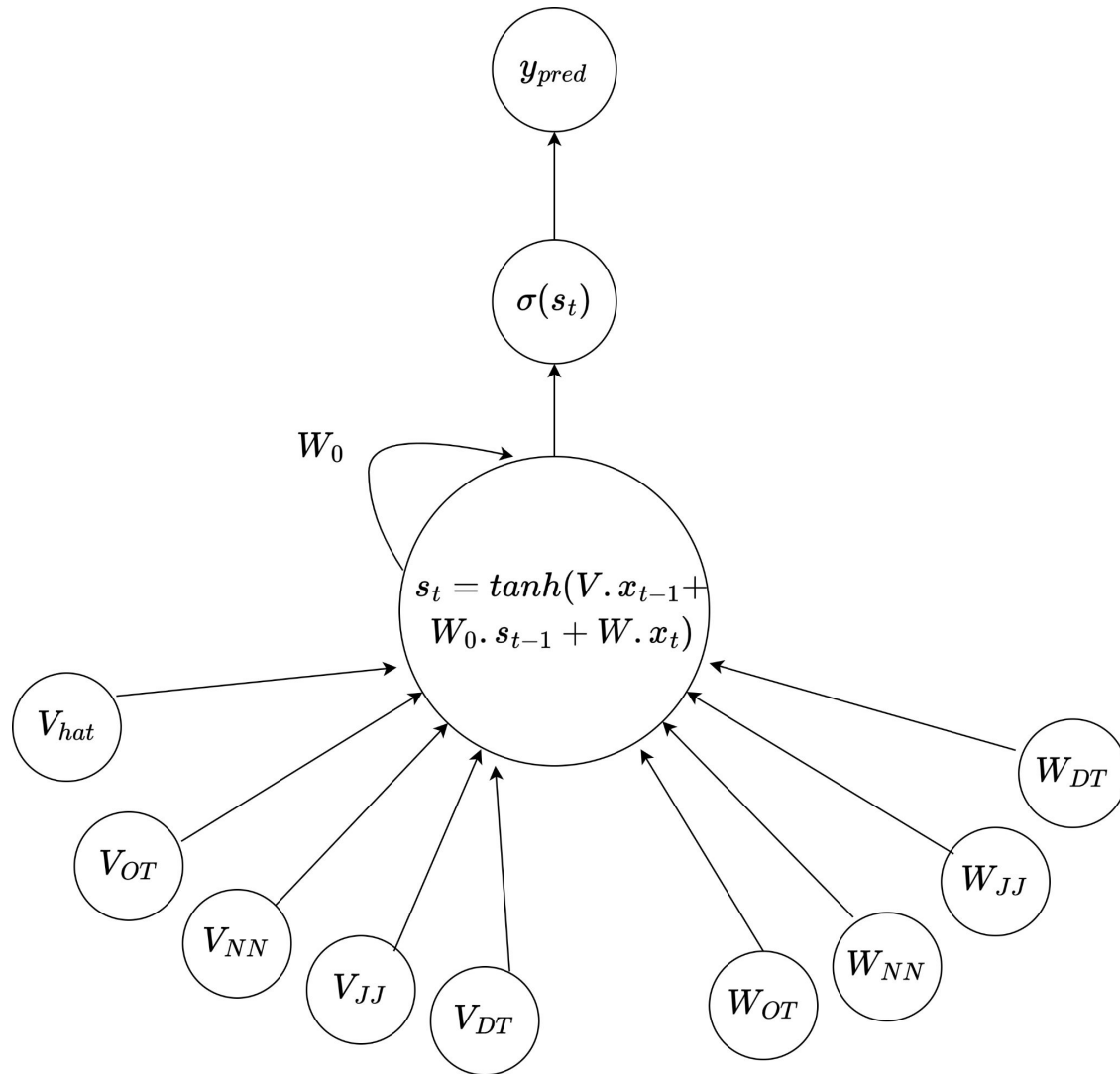
$$W \in \mathbb{R}^{4 \times 1}$$

$$V \in \mathbb{R}^{5 \times 1}$$

$$W_0 \in \mathbb{R}^{1 \times 1}$$

Here  $l$  is the length of the sentence.

# Implementation Details



# Implementation Details

- The BPTT equations with respect to weights  $W, W_0, V$  are given as follows at 3rd time step-

$$\begin{aligned}\frac{\partial E_3}{\partial W} &= \frac{\partial E_3}{\partial \hat{y}^3} \frac{\partial \hat{y}^3}{\partial s_3} \frac{\partial s_3}{\partial W} \\ \frac{\partial E_3}{\partial W_0} &= \frac{\partial E_3}{\partial \hat{y}^3} \frac{\partial \hat{y}^3}{\partial s_3} \frac{\partial s_3}{\partial W_0} \\ \frac{\partial E_3}{\partial V} &= \frac{\partial E_3}{\partial \hat{y}^3} \frac{\partial \hat{y}^3}{\partial s_3} \frac{\partial s_3}{\partial V}\end{aligned}$$

- The last term of each equation which will incorporate propagation towards back upto initial time.

# Implementation Details

- We have used cross-entropy loss
- $Lr = 0.001$
- Experimented with uniform initialization between  $[-5, 5]$  and  $[-1/\sqrt{4}, 1/\sqrt{4}]$ .
- Experimented with adaptive learning rate
- Experimented with linear activation function versus non-linear activation function.

# Overall performance

- **Quantitative Results:**
  - We have used 4 metrics for evaluation using 5-fold cross validation and full dataset training.
  - Upper row reports metrics on the held out fold and the below row on the full test dataset.
  - We see that initialization plays an important role in performance.

LeCun Uniform Initialization	Dataset	Accuracy	Precision	Recall	F1
	Fold-1	37.1 38.9	8.2 8.0	88.3 81.8	0.15 0.15
	Fold-2	37.8 38.9	8.3 8.0	87.7 81.8	0.15 0.15
	Fold-3	37.6 38.9	8.2 8.0	86.8 81.8	0.15 0.15
	Fold-4	50.7 50.2	45.4 40.9	71.4 70.3	0.55 0.52
	Fold-5	61.1 59.7	88.7 88.3	65.5 63.8	0.75 0.74
	Full-train	50.2	40.9	70.3	0.52
Uniform Initialization [-5, 5]	Fold-1	65.0 62.9	78.4 78.7	72.3 68.8	0.75 0.73
	Fold-2	64.4 62.9	78.4 78.7	71.3 68.8	0.75 0.73
	Fold-3	67.0 <b>65.2</b>	100.0 <b>100.0</b>	67.0 <b>65.2</b>	0.80 <b>0.79</b>
	Fold-4	64.7 62.9	78.5 78.7	71.9 69.8	0.75 0.73
	Fold-5	64.3 62.9	78.0 78.7	71.4 69.8	0.74 0.73
	Full-train	<b>62.9</b>	<b>78.7</b>	<b>68.8</b>	<b>0.73</b>

# Language constraint table

- Learnt weight values:

V_hat	16627.57	W_0	-10772.25
V_NN	16626.59	W_NN	-10796.90
V_DT	16634.67	W_DT	-10798.23
V_JJ	16632.02	W_JJ	2.31
V_OT	16633.09	W_OT	0.98
theta	4885.66		



# Language constraint table

- Constraints table (1/2):

Current(W) / Prev (V)	DT	JJ	NN	OT
hat	$f(V_{\text{hat}} + W_{\text{DT}} + \theta) > 0.5$ <b>YES</b>	$f(V_{\text{hat}} + W_{\text{JJ}} + \theta) > 0.5$ <b>YES</b>	$f(V_{\text{hat}} + W_{\text{NN}} + \theta) > 0.5$ <b>YES</b>	$f(V_{\text{hat}} + W_{\text{OT}} + \theta) > 0.5$ <b>YES</b>
DT	x	$f(W_0 + V_{\text{DT}} + W_{\text{JJ}} + \theta) < 0.5$ <b>NO</b>	$f(W_0 + V_{\text{DT}} + W_{\text{NN}} + \theta) < 0.5$ <b>YES</b>	x
JJ	x	$f(V_{\text{JJ}} + W_{\text{JJ}} + \theta) < 0.5$ <b>NO</b>	$f(V_{\text{JJ}} + W_{\text{NN}} + \theta) < 0.5$ <b>NO</b>	x
	x	$f(W_0 + V_{\text{JJ}} + W_{\text{JJ}} + \theta) < 0.5$ <b>NO</b>	$f(W_0 + V_{\text{JJ}} + W_{\text{NN}} + \theta) < 0.5$ <b>YES</b>	x

# Language constraint table

- Constraints table (2/2):

Current(W) / Prev (V)	DT	JJ	NN	OT
NN	x	x	x	$f(V\_NN + W\_OT + \theta) > 0.5$ <b>YES</b>
	x	x	x	$f(W\_0 + V\_NN + W\_OT + \theta) > 0.5$ <b>YES</b>
OT	$f(W\_0 + V\_OT + W\_DT + \theta) > 0.5$ <b>NO</b>	$f(W\_0 + V\_OT + W\_JJ + \theta) > 0.5$ <b>YES</b>	$f(W\_0 + V\_OT + W\_NN + \theta) > 0.5$ <b>NO</b>	$f(W\_0 + V\_OT + W\_OT + \theta) > 0.5$ <b>YES</b>

- In total, 10 out of 16 constraints are satisfied.

# Error Analysis - Qualitative

- Sentences from test set which received **over 85% accuracy**:
  - **Example: 1**
    - Tokens: ['RKC', 'Waalwijk', '1', 'Willem', 'II', 'Tilburg', '2']
    - POS Tags: [1, 1, 4, 1, 1, 1, 4]
    - Chunk Tags: [1, 0, 0, 0, 0, 0, 0]
    - Predicted Chunk Tags: [0 0 0 0 0 0 0]
    - Accuracy: **85.71%**
  - **Example: 2**
    - Tokens: ['Ironi', 'Rishon', 'Lezion', '1', 'Maccabi', 'Herzliya', '0']
    - POS Tags: [1, 1, 1, 4, 1, 1, 4]
    - Chunk Tags: [1, 0, 0, 0, 0, 0, 0]
    - Predicted Chunk Tags: [0 0 0 0 0 0 0]
    - Accuracy: **85.71%**

# Error Analysis - Qualitative

- Sentences from test set which received **between 60%-70% accuracy**:
- **Example: 1**
  - Tokens: ['UAE', '-', 'Hassan', 'Ahmed', '53', ',', 'Adnan', 'Al', 'Talyani', '55', ',', 'Bakhit', 'Saad', '80']
  - POS Tags: [1, 4, 1, 1, 4, 4, 1, 1, 1, 4, 4, 1, 1, 4]
  - Chunk Tags: [1, 1, 1, 0, 0, 1, 1, 0, 0, 0, 1, 1, 0, 0]
  - Predicted Chunk Tags: [0 0 0 0 0 1 0 0 0 0 1 0 0 0]
  - Accuracy: **66.67%**
- **Example: 2**
  - Tokens: ['Tasmania', '481', 'for', 'eight', 'declared', '(', 'Michael', 'DiVenuto', '119', ',', 'David', 'Boon', '118', ',', 'Shaun', 'Young', '113', ')', ';', 'Victoria', '220', 'for', 'three', '(', 'Dean', 'Jones', '130', 'not', 'out', ')', '.']
  - POS Tags: [1, 4, 4, 4, 1, 4, 1, 1, 4, 4, 1, 1, 4, 4, 1, 1, 4, 4, 4, 4, 1, 1, 4, 4, 4, 4]
  - Chunk Tags: [1, 0, 1, 1, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1]
  - Predicted Chunk Tags: [0 0 1 0 0 0 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 1 0 0 0 1 0 1 0]
  - Accuracy: **66.67%**

# Error Analysis - Qualitative

- Sentences from test set which received between **40%-60%** accuracy:
- **Example: 1**
  - Tokens: ['The', 'lanky', 'former', 'Leeds', 'United', 'defender', 'did', 'not', 'make', 'his', 'England', 'debut', 'until', 'the', 'age', 'of', '30', 'but', 'eventually', 'won', '35', 'caps', 'and', 'was', 'a', 'key', 'member', 'of', 'the', '1966', 'World', 'Cup', 'winning', 'team', 'with', 'his', 'younger', 'brother', ',', 'Bobby', '.']
  - POS Tags: [2, 3, 3, 1, 1, 1, 4, 4, 4, 4, 1, 1, 4, 2, 1, 4, 4, 4, 4, 4, 4, 1, 4, 4, 2, 3, 1, 4, 2, 4, 1, 1, 3, 1, 4, 4, 3, 1, 4, 1, 4]
  - Chunk Tags: [1, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 1, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 1]
  - Predicted Chunk Tags: [0 1 0 0 0 0 0 1 0 1 0 0 0 1 0 0 1 0 1 0 1 0 1 0 0 1 0 1 0 0 1 0 0 1 0 0 0 0 0 1 0 0 0 0 0]
  - Accuracy: **56.09%**

# Error Analysis - Qualitative

- Sentences from test set which received **below 40% accuracy**:
- **Example: 1**
  - Tokens: ['It', 'all', 'culminated', 'in', 'the', 'fact', 'that', 'I', 'now', 'have', 'lots', 'of', 'great', ',', 'great', 'friends', 'in', 'Ireland', '.']
  - POS Tags: [4, 2, 4, 4, 2, 1, 4, 4, 4, 4, 1, 4, 3, 4, 3, 1, 4, 1, 4]
  - Chunk Tags: [1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1]
  - Predicted Chunk Tags: [0 1 0 1 0 1 0 1 0 1 0 0 1 0 1 0 0 0 0]
  - Accuracy: **36.84%**
- **Example: 2**
  - Tokens: ['(', 'tabulate', 'under', 'won', ',', 'lost', ',', 'percentage', ',', 'games', 'behind', ')', ':']
  - POS Tags: [4, 1, 4, 3, 4, 4, 4, 1, 4, 1, 4, 4, 4]
  - Chunk Tags: [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]
  - Predicted Chunk Tags: [0 0 0 1 0 1 0 0 0 0 0 0 1 0]
  - Accuracy: **23.07%**

# Learnings

- Weight initialization plays an important role in convergence of models.
- Even a very small model like recurrent perceptron does pretty good job in classifying noun chunk.
- The simpler the model, the more explainable it is.
- As per our finding large values of weights could be controlled using some regularisation technique or gradient clipping.

**Demo**



**Thank You**

# Metrics used

- Here, TP = True Positive, FP = False Positive.

$$\text{Accuracy (AC)} = \frac{TP + TN}{TP + TN + FP + FN}$$

- For Binary Cross Entropy loss function:

$$\text{Precision (PC)} = \frac{TP}{TP + FP}$$

- $y'$  is the ground truth (actual label) and
- $y$  is the predicted label (neural network's output) of the class.

$$\text{F1-Score} = \frac{2 \times (PC \times SE)}{PC + SE}$$

$$\text{Recall or Sensitivity (SE)} = \frac{TP}{TP + FN}$$

$$L_{y'}(y) := -\frac{1}{N} \sum_{i=1}^N (y'_i \log(y_i) + (1 - y'_i) \log(1 - y_i))$$

# k-Fold

- We have used 4 metrics for evaluation using 5-fold cross validation and full dataset training.
- Upper row reports metrics on the held out fold and the below row on the full test dataset, (mean +- standard deviation)
- We see that initialization plays an important role in performance.

	Dataset	Accuracy	Precision	Recall	F1-score
LeCun Uniform Initialization	Fold-1	44.4 (24.3) 46.7 (25.1)	14.3 (28.5) 13.8 (28.6)	34.9 (46.3) 31.5 (45.1)	0.16 (0.28) 0.15 (0.27)
	Fold-2	45.2 (24.0) 46.7 (25.1)	13.8 (27.5) 13.8 (28.6)	36.4 (46.8) 31.5 (45.1)	0.16 (0.27) 0.15 (0.27)
	Fold-3	44.6 (24.2) 46.7 (25.1)	13.7 (28.0) 13.8 (28.6)	34.6 (46.3) 31.5 (45.1)	0.16 (0.27) 0.15 (0.27)
	Fold-4	53.3 (26.2) 53.0 (26.4)	39.7 (48.9) 34.9 (47.7)	28.0 (36.8) 23.7 (34.9)	0.32 (0.41) 0.27 (0.39)
	Fold-5	60.6 (21.8) 58.6 (22.9)	89.7 (12.6) 90.0 (12.9)	64.3 (24.0) 61.7 (25.9)	0.71 (0.20) 0.69 (0.22)
	Full-train	53.0 (26.4)	34.9 (47.7)	23.7 (34.9)	0.27 (0.39)
Uniform Initialization [-5, 5]	Fold-1	63.8 (19.2) 60.1 (21.5)	84.1 (13.4) 85.5 (13.3)	68.5 (23.1) 63.8 (25.9)	0.72 (0.18) 0.68 (0.20)
	Fold-2	62.9 (19.7) 60.1 (21.5)	84.2 (13.6) 85.5 (13.3)	67.4 (23.6) 63.8 (25.9)	0.71 (0.19) 63.8 (25.9)
	Fold-3	62.7 (22.9) 59.4 (24.9)	100.0 (0.0) 100.0 (0.0)	62.7 (22.9) 59.4 (24.9)	0.74 (0.20) 0.71 (0.23)
	Fold-4	62.8 (19.7) 60.1 (21.5)	84.0 (13.3) 85.5 (13.3)	67.6 (23.7) 63.8 (25.9)	0.71 (0.18) 0.68 (0.20)
	Fold-5	62.5 (19.3) 60.1 (21.5)	83.9 (13.5) 85.5 (13.3)	66.8 (23.7) 63.8 (25.9)	0.71 (0.19) 0.68 (0.20)
	Full-train	<b>60.1 (21.5)</b>	<b>85.5 (13.3)</b>	<b>63.8 (25.9)</b>	<b>0.68 (0.20)</b>