May 3, 2023.        DSML : NLP module.    Word embeddings in a nutshell

# Recurrent Neural Networks.

Class starts @ 9:05

**Recap:**

* Document Vectorization: Bag of Words (BoW) TF/IDF.

* Word Vectorization: Continuous BoW, Skip-gram.

* Language Modeling: Naïve Bayes to predict the next word.

* Topic Modeling: Heuristics, Latent Dirichlet Allocation

# Agenda:

* Recurrent Neural Networks. (RNNs)

* A new type of Architecture.

* Different variants.

* Training - Backpropagation through Time.

But first, business case: Classifying News Articles.

## Sequence Data :

* **I** travelled to **France** last December.

1    2    3    4    5    6

$t=1$    $t=2$    $t=4$

* A   C V   conference was held at Nice.

7    8    9    10    11    12   13

* It lasted for three days.

15    16    17    18    19.

# Approaches:

**1]** Doc → TF/IDF → <u>Multiclass classification.</u>
SVMs, Decision trees,

**2]** <u>Neural network approaches:</u>

Problem 1] Length of inputs is varying.

Solution: clip the input.

&lt;blank&gt; &lt;blank&gt;

max_length: 12 words.

| title |
|---|
| Ad sales boost Time Warner profit → 7 k |
| Dollar gains on Greenspan speech → S |
| Yukos unit buyer faces loan claim → 6 |
| High fuel prices hit BA's profits |
| Pernod takeover talk lifts Domecq |
| ... |
| BT program to beat dialler scams |
| Spam e-mails tempt net shoppers |
| Be careful how you code |
| US cyber security chief resigns |
| Losing yourself in online gaming |

→ 0
→→ 0
→ 0
− 0
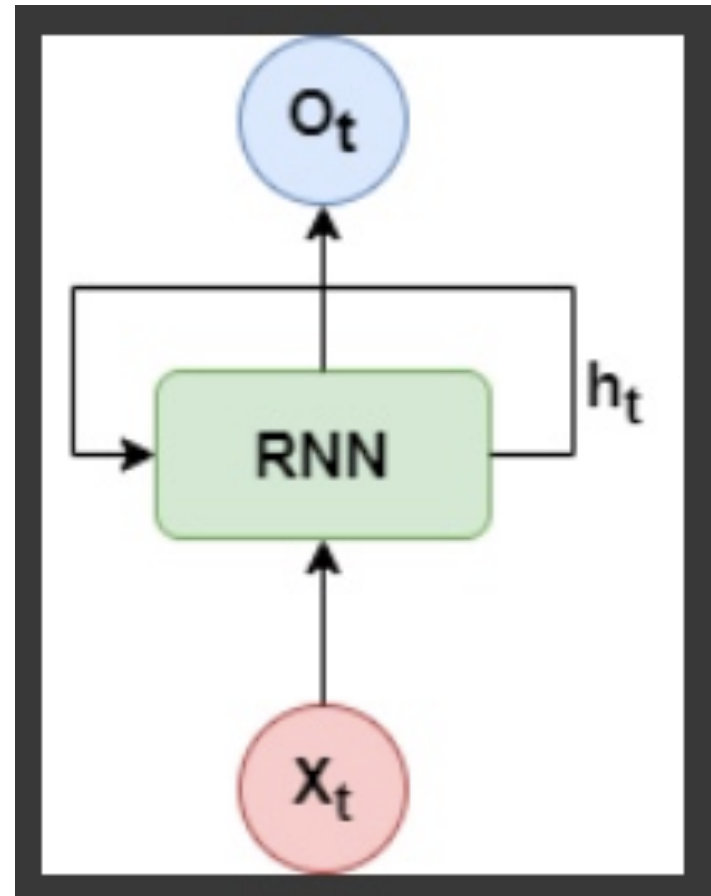⋮
0

Problem 2] Content information.

Solution: N-grams. Force the NN to see N words at a time.

Length of the N-gram is very difficult to optimize.
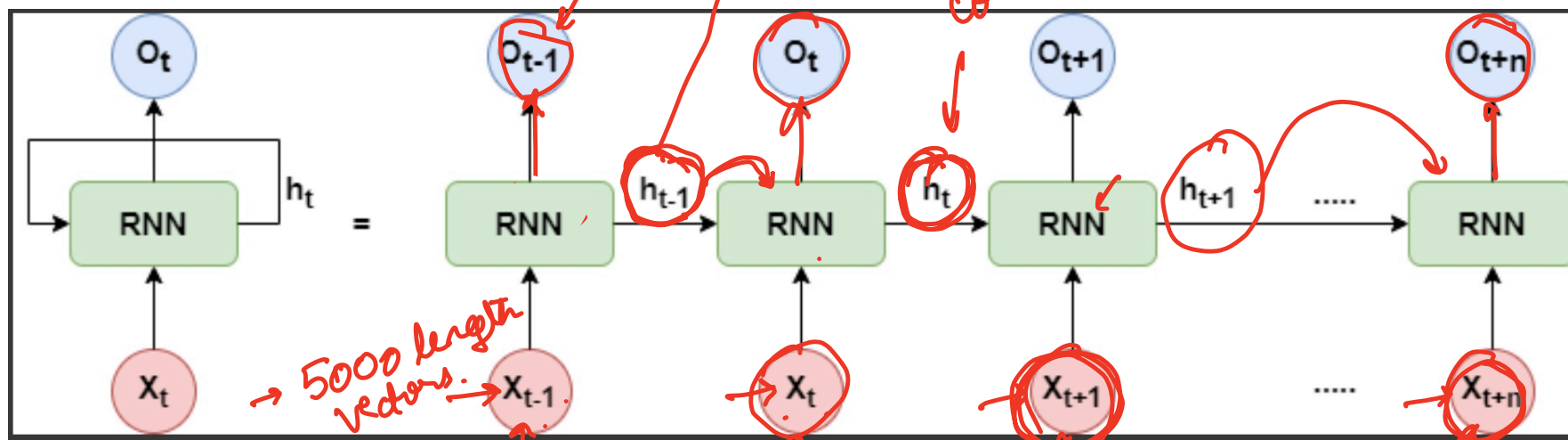
Problem 3]. Time & Space complexity.

Motivation behind this
architecture:

→ NN lack memory.

→ How do we give
   it memory?

# Recurrent Neural Networks.
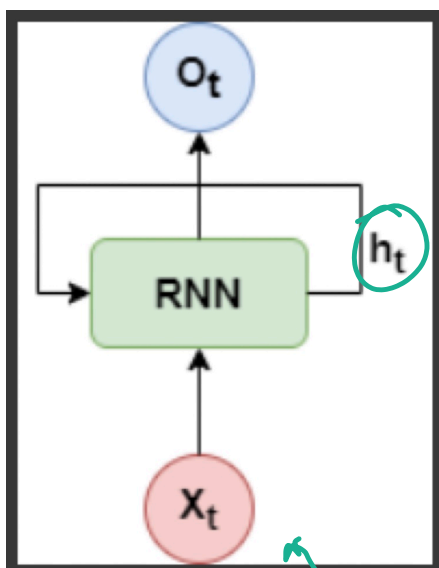
↳ "Recursion"

hidden state.
contextual info.



→ 5000 length vectors.

OHE    OHE    OHE

OHE

I    travelled    to    France    last    December.

|V| = 5000

$O_t$

RNN  $h_t$

$X_t$

— Garfield is a cat
　　　①　　②　③　　　④

He　is　fat.
①　②　③

# RNN architecture deep dive.

$$\frac{\partial L_t}{\partial V} = \frac{\partial L_t}{\partial O_t} * \frac{\partial O_t}{\partial Z_t} * \frac{\partial Z_t}{\partial V}$$

$t = 0$.

① $h_0 = 0$ vector

② $h_0 \to$ random.

③ $h_0 -$ some fixed vector.

$\frac{\partial L}{\partial O_t}$

$100$

Output at time $t$.

$h_{t-1}$  $W_{nh}$  $200$

$W_{hx}$  $+$  $200$

$200$

$100$

$tanh \to$ Activation.

$Z_t$  $200$

$O_t$

$h_t \to$ next hidden state.

$x_t \to$ Vector Representation of our input word at time $t$.

① OHE word.

② Fixed Word2Vec.

③ Trained Word2Vec.

Loss = $L(Y, O_t)$

$O_t$

$Z_{y1}$    $Z_{y2}$    $Z_{y3}$    $Z_{y4}$    $Z_{y5}$

$h_{t-4}$   W    tanh   V    $h_{t-1}$   W   tanh   V    $h_{t-3}$   W   tanh   V    $h_{t+2}$   W   tanh   V    $h_{t-1}$   W   tanh   V   $h_t$

$Z_{h1}$    $Z_{h2}$    $Z_{h3}$    $Z_{h4}$    $Z_{h5}$

U    U    U    U    U

Embed    Embed    Embed    Embed    Embed

$X_1$ t-4    $X_2$ t-3    $X_3$ t-2    $X_4$ t-1    $X_5$ t

Federer    joins    all    time    greats

$Z_{h1} = (X_1 * U) + (h_{t-5} * W) + b_h$
$h_1 = \sigma(Z_{h1})$, where $\sigma$ = tanh
$Z_{y1} = (V * h_{t-4}) + b_y$

$Z_{h2} = (X_2 * U) + (h_{t-4} * W) + b_h$
$h_2 = \sigma(Z_{h2})$, where $\sigma$ = tanh
$Z_{y2} = (V * h_{t-3}) + b_y$

$Z_{h3} = (X_3 * U) + (h_{t-3} * W) + b_h$
$h_3 = \sigma(Z_{h3})$, where $\sigma$ = tanh
$Z_{y3} = (V * h_{t-2}) + b_y$

$Z_{h4} = (X_4 * U) + (h_{t-2} * W) + b_h$
$h_4 = \sigma(Z_{h4})$, where $\sigma$ = tanh
$Z_{y4} = (V * h_{t-1}) + b_y$

$Z_{h5} = (X_5 * U) + (h_{t-1} * W) + b_h$
$h_5 = \sigma(Z_{h5})$, where $\sigma$ = tanh
$Z_{y5} = (V * h_t) + b_y$
$O_{t+4} = \sigma(Z_{y5})$, where $\sigma$ = Softmax
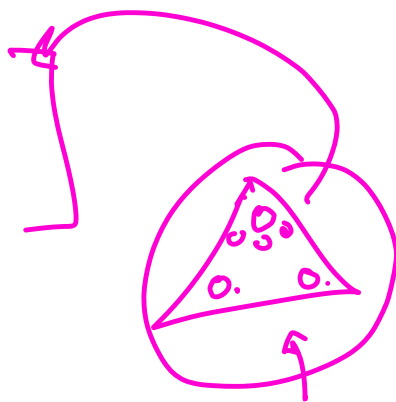
$$\nabla_{O_t} h \qquad \nabla_W \ell = \nabla_{O_t} \ell \times \nabla_{Z_{h5}} O_t \times \nabla_W Z_{h_5}$$

(2$^{nd}$ time)

$$\nabla_W \ell = \nabla_{Z_{h4}} \ell \times \nabla_W Z_{h4}$$

$$Z_{h_4} = (X_4 * U) + \sigma\left((X_5 * U) + (h_{t-1} * W) + b_h\right) * W) + b_h$$

$k = 10.$