

## Αλληλεπίδραση Ανθρώπου-Υπολογιστή

B6. Επεξεργασία φυσικής γλώσσας με ανατροφοδοτούμενα νευρωνικά δίκτυα

#### (2024–25)

Ίων Ανδρουτσόπουλος http://www.aueb.gr/users/ion/

## Contents

- Recurrent neural networks (RNNs), GRUs/LSTMs.
- Applications in token classification (e.g., named entity recognition).
- RNN language models.
- Layer normalization.
- RNNs with self-attention and applications in text classification.
- Bidirectional and stacked RNNs.
- Sequence-to-sequence RNN models with attention, and applications in machine translation.

## Extracting contract elements

THIS AGREEMENT is made the 15th day of October 2009 (The "Effective Date") BETWEEN:

Sugar 13 Inc., a corporation whose office is at James House,
42-50 Bond Street, London, EW2H TL ("Sugar");

(2) **E2 UK Limited**, a limited company whose registered office is at 260 Bathurst Road, Yorkshire, SL3 4SA (**"Provider"**).

#### **RECITALS:**

A. The Parties wish to enter into a framework agreement which will enable Sugar, from time to time, to [...]B. [...]

#### NO THEREFORE IT IS AGREED AS FOLLOWS:

#### **ARTICLE I - DEFINITIONS**

"Sugar" shall mean: Sugar 13 Inc.

"Provider" shall mean: E2 UK Limited

"1933 Act" shall mean: Securities Act of 1933

#### **ARTICLE II - TERMINATION**

The Service Period will be for five (5) years from the Effective Date (The "Initial Term"). The agreement is considered to be terminated in October 16, 2014.

#### **ARTICLE III - PAYMENT - FEES**

During the service period monthly payments should occur. The estimated fees for the Initial Term are £154,800.

#### **ARTICLE IV - GOVERNING LAW**

This agreement shall be governed and construed in accordance with the Laws of England & Wales. Each party hereby irrevocably submits to the exclusive jurisdiction of the courts sitting in Northern London.

**IN WITNESS WHEREOF**, the parties have caused their respective duly authorized officers to execute this Agreement.

**BY:** George Fake Authorized Officer Sugar 13 Inc.

**BY:** Olivier Giroux CEO E2 UK LIMITED

Identify start/end dates, duration, contractors, amount, legislations refs, jurisdiction etc. Similar to Named Entity Recognition (NER).

I. Chalkidis, I. Androutsopoulos, A. Michos, "Extracting Contract Elements", ICAIL 2017, http://nlp.cs.aueb.gr/pubs/icail2017.pdf.

### Window-based token classification



## Window-based token classification



We learn  $W^{(1)}$ ,  $W^{(2)}$  with **backpropagation**. We can also learn (or modify) the **word embeddings** *E* during **backpropagation**! But when we don't have large training datasets (e.g., corpus manually annotated with B-I-O tags), it may be better to use **pre-trained embeddings**, which can be obtained from large non-annotated corpora (e.g., via Word2Vec, GloVe).

We can use the same window-based approach for **POS-tagging**, **chunking**, ...



## **RNN-based** token classification





## RNN language model



### Reminder: LMs as next word predictors

Sequence probability using a bigram LM:

$$P(w_1^k) = P(w_1, \dots, w_k) = P(w_1) \cdot P(w_2 \mid w_1) \cdot P(w_3 \mid w_1, w_2) \cdot P(w_4 \mid w_1^3) \cdots P(w_k \mid w_1^{k-1}) \simeq$$

$$P(w_1 | start) \cdot P(w_2 | w_1) \cdot P(w_3 | w_2) \cdots P(w_k | w_{k-1})$$

- We can think of the LM as a system that provides the probabilities  $P(w_i|w_{i-1})$ , which we then multiply.
  - Or the probabilities  $P(w_i|w_{i-2}, w_{i-1})$  for a trigram LM.
  - Or the probabilities P(w<sub>i</sub>|h) for an LM that considers all the "history" (previous words) h, e.g., in an RNN LM.
- An LM typically provides a distribution P(w|h) showing how probable it is for every word  $w \in V$  to be the next one.



## More about RNNs

- Trained by **backpropagation** (with **unrolled** view).
  - For each sentence (or window), feed it to the unrolled RNN, compute the loss and backpropagate, adding gradients obtained for the same matrix (e.g., same W<sup>(h)</sup> at each cell).
  - GRU or LSTM cells help avoid vanishing gradients.
  - The norms of the **gradients** can be **clipped** (when larger than a max value) to avoid **exploding gradients**.
  - Use layer normalization, not batch normalization in RNNs.
- We can also **learn** the **word embeddings** (*E*) with an RNN LM. Billions of **free training examples**!
  - We can then use the **word embeddings** in **other NLP tasks**.
  - With a **large vocabulary**, **softmax** is too **slow** (alternatives: small vocabulary, hierarchical softmax, negative sampling).

#### What about the right-context of each token?





#### Stacked bidirectional RNN



Each layer revises the word embeddings of the previous (lower) layer. The embeddings become increasingly more context-aware and also increasingly more appropriate for the particular task we address...

#### Token classification with a stacked biRNN

**Compare to the correct** predictions (**sum** the **cross-entropy loss** for **all token positions**) and **backpropagate** to **adjust all the weights**, including the weights of the stacked biRNN.



#### User comment moderation



#### RNN with deep self-attention



J. Pavlopoulos, P. Malakasiotis and I. Androutsopoulos, "Deeper Attention to Abusive User Content Moderation", EMNLP 2017, <u>http://nlp.cs.aueb.gr/pubs/emnlp2017.pdf</u>.

#### RNN with deep self-attention

The entire input text is now represented by the weighted (by *a<sub>i</sub>* scores) sum of the revised embeddings of its words.  $\alpha_k$  $\alpha_2 \mathbf{x}$  $h_2$  + softmax dense  $W^{(o)}$ & softmax rejectior orobabilit acceptance Attention MLP probabilitv RNN

Hello there ... relax

We pass the **weighted sum vector** (point) through another **dense layer and softmax** to obtain a **probability** score for **each class** (here accept, reject).

**Compare to the correct** predictions with a **cross-entropy loss** and **backpropagate** to **adjust the weights** of the **entire neural net**, including the MLP and RNN(s).

The attention scores  $a_i$  can also be used to highlight the words that influence the system's decision most.

Go	and	hang	yourself	!				
You	are	ignorant	and	vandal	!	Stop	it	!
Thanks	•	Please	go	<b>ROK</b>	yourself		ty	!

J. Pavlopoulos, P. Malakasiotis and I. Androutsopoulos, "Deeper Attention to Abusive User Content Moderation", EMNLP 2017, <u>http://nlp.cs.aueb.gr/pubs/emnlp2017.pdf</u>.

#### RNN with deep self-attention



J. Pavlopoulos, P. Malakasiotis and I. Androutsopoulos, "Deeper Attention to Abusive User Content Moderation", EMNLP 2017, <u>http://nlp.cs.aueb.gr/pubs/emnlp2017.pdf</u>.

#### Text classification with stacked biRNN

Compare (via categorical cross entropy) the predicted  $\vec{o}$  to the correct 1-hot distribution and backpropagate to adjust all the weights, including the weights of the stacked biRNN.



### **RNNs for Machine Translation**

From the slides of R. Socher's course "Deep Learning for NLP", 2015. http://cs224d.stanford.edu/ Encoder:  $h_t = \phi(h_{t-1}, x_t) = f\left(W^{(hh)}h_{t-1} + W^{(hx)}x_t\right)$ Decoder:  $h_t = \phi(h_{t-1}) = f\left(W^{(hh)}h_{t-1}\right)$  $y_t = softmax\left(W^{(S)}h_t\right)$ 

Minimize cross entropy error for all target words conditioned on source words







### **RNN-based** Machine Translation







Images from Stephen Merity's http://smerity.com/articles/2016/ google\_nmt\_arch.html



### Basic Encoder-Decoder NMT



**During training**, at each time-step of the **decoder**, we can use the **correct previous word** of the human translation (**teacher forcing**), or we can **randomly use the correct or the predicted** previous word (**scheduled sampling**).

**During testing (inference)**, we always use the **predicted previous word**; and we either **greedily select the most probable next word**, or we use **beam search** to find the translation  $y_1^m$  of  $x_1^n$  with the highest probability:  $p(y_1|x_1^n) p(y_2|y_1, x_1^n) p(y_3|y_1^2, x_1^n) \dots p(y_m|y_1^{m-1}, x_1^n)$ Google's paper: <u>https://arxiv.org/abs/1609.08144</u> Images from Stephen Merity's <u>http://smerity.com/articles/2016/google\_nmt\_arch.html</u>

#### Encoder-Decoder with attention

The source sentence is now represented by the weighted sum of the encoder states:



For each German word, the **attention scores** over the English words **change**! Each "**attention**" weight  $a_j$  is a function of the corresponding encoder state  $\vec{h}_j$ and the **previous state**  $\vec{z}_{i-1}$  of the decoder (memory of translation so far), e.g.:  $\tilde{a}_j = v^T \cdot f(W^{(h)}\vec{h}_j + W^{(z)}\vec{z}_{i-1}) = v^T \cdot f(W[\vec{h}_j; \vec{z}_{i-1}]), \quad a_j = softmax(\tilde{a}_j)$ with a **softmax** to make the  $a_j$  weights sum to 1.

Google's paper: <u>https://arxiv.org/abs/1609.08144</u> Images from Stephen Merity's <u>http://smerity.com/articles/2016/google\_nmt\_arch.html</u>

## Bidirectional LSTM encoder

The encoder is now a **bidirectional LSTM**. The **encoder state** for the *j*-th word of the source sentence is the concatenation of the corresponding states of the forward and backward LSTM.



Google's paper: <u>https://arxiv.org/abs/1609.08144</u> Images from Stephen Merity's <u>http://smerity.com/articles/2016/google\_nmt\_arch.html</u>

## Stacking RNNs and residuals



"Residual" connections (a kind of skip-connections) helps fight vanishing gradients in backpropagation sum-nodes copy the gradients to their inputs). Also allows upper layers to learn only modifications (differences) from representations of lower layers.

Google's paper: <u>https://arxiv.org/abs/1609.08144</u>

Images from Stephen Merity's <u>http://smerity.com/articles/2016/google\_nmt\_arch.html</u>

## **RNN-based** Machine Translation

Softmax

Decoder

Decoder

Decoder

Decoder

Decoder

Decoder

Decoder

Sum

Sum

Sum

Sum

Sum

Sum



Attention based on the previous state of the bottom decoder only, to speed up computations. 29

Google's paper: https://arxiv.org/abs/1609.08144

Images from Stephen Merity's http://smerity.com/articles/2016/ google\_nmt\_arch.html

# Recommended reading

- M. Surdeanu and M.A. Valenzuela-Escarcega, *Deep Learning for Natural Language Processing: A Gentle Introduction*, Cambridge Univ. Press, 2024.
  - Chapters 11, 12, 14. See <u>https://clulab.org/gentlenlp/text.html</u>
  - Also available at AUEB's library.
- Y. Goldberg, Neural Network Models for Natural Language Processing, Morgan & Claypool Publishers, 2017.
  Mostly abortors 14, 17

• Mostly chapters 14–17.

- Jurafsky and Martin's, *Speech and Language Processing* is being revised (3<sup>rd</sup> edition) to include DL methods.
  - o <u>http://web.stanford.edu/~jurafsky/slp3/</u>









# Recommended reading

- F. Chollet, *Deep Learning in Python*, 1<sup>st</sup> edition, Manning Publications, 2017.
  - 1<sup>st</sup> edition freely available (and sufficient for this part of the course): <u>https://www.manning.com/books/deep-learning-with-python</u>
  - See mostly sections 6.1–6.3, section 8.1.
  - 2<sup>nd</sup> edition (2022) now available, requires payment. Highly recommended.
- See also the recommended reading and resources of the previous part (NLP with MLPs) of this course.



# Βιβλιογραφία – συνέχεια

 Αν έχετε από το μάθημα της ΤΝ το βιβλίο των Russel & Norvig «Τεχνητή Νοημοσύνη – Μια σύγχρονη προσέγγιση», 4<sup>η</sup> έκδοση, Κλειδάριθμος, 2021, μπορείτε να συμβουλευτείτε τα κεφάλαια 21 και 24.



- Κυρίως τις ενότητες 21.6, 21.8.2, 24.1, 24.2, 24.3.
- Άλλες ενότητες αυτών των κεφαλαίων θα καλυφθούν σε επόμενες διαλέξεις.