Neural Language Models



Tanmoy Chakraborty Associate Professor, IIT Delhi https://tanmoychak.com/

Slides are adopted from the Stanford course 'NLP with DL' by C. Manning

- ▼ Chapter 05. Neural Language Models
 - 5.1 Convolutional Neural Networks
 - ► 5.2 Recurrent Neural Networks
 - ► 5.3 Sequence-to-Sequence Models
 - ► 5.4 Attention Mechanisms
 - 5.5 Limitations of Neural Language Models
 - 5.6 Summary



Generative AI for Text

Tanmoy Chakraborty



Pre-requisite for this chapter

- Loss function, backpropagation
- CNN
- RNN (LSTM/GRU)

Recall: Language Modeling

• Language Modeling is the task of predicting what word comes next



Recall: Language Modeling

- You can also think of a Language Model as a system that assigns a probability to a piece of text.
- For example, if we have some text $x^{(1)}$, ..., $x^{(T)}$, then the probability of this text (according to the Language Model) is:

$$P(\boldsymbol{x}^{(1)}, \dots, \boldsymbol{x}^{(T)}) = P(\boldsymbol{x}^{(1)}) \times P(\boldsymbol{x}^{(2)} | \boldsymbol{x}^{(1)}) \times \dots \times P(\boldsymbol{x}^{(T)} | \boldsymbol{x}^{(T-1)}, \dots, \boldsymbol{x}^{(1)})$$
$$= \prod_{t=1}^{T} P(\boldsymbol{x}^{(t)} | \boldsymbol{x}^{(t-1)}, \dots, \boldsymbol{x}^{(1)})$$

This is what our LM provides







How to Build a Neural Language Model?

- Recall the Language Modeling task:
 - Input: sequence of words $x^{(1)}, x^{(2)}, \dots, x^{(t)}$
 - Output: probability distribution of the next word $Pig(x^{(t+1)}ig|x^{(t)},\dots,x^{(1)}ig)$
- How about a window-based neural model?



Example: NER Task



A Fixed-window Neural Language Model



A Fixed-window Neural Language Model



A Fixed-window Neural Language Model

Improvements over *n*-gram LM:

- No sparsity problem
- Don't need to store all observed ngrams

Remaining problems:

- Fixed window is too small
- Enlarging window enlarges W
- x⁽¹⁾ and x⁽²⁾ are multiplied by completely different weights in W.
 No symmetry in how the inputs are processed.



Approximately: Y. Bengio, et al. (2000/2003): A Neural Probabilistic Language Model

> We need a neural architecture that can process **any length input**

