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



Grok 3 and X Premium+ Update !!

Released on February 19, 2025

Grok 3

Advanced AI Capabilities: Access to Grok 3 for enhanced AI interactions. Enhanced User Experience: Ad-free browsing, post editing, and longer content sharing. Monetization Opportunities: Tools like Get Paid to Post and Creator Subscriptions.

Grok 3 is xAI's latest AI chatbot, offering advanced capabilities, trained on 200,000 GPUs, and available to X Premium+ subscribers. X premium+: Enhanced subscription plan with features like post editing, ad-free experience, and monetization tools.



- ▼ 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**



Recurrent Neural Networks (RNN)



Core idea: Apply the same weights *W* repeatedly









RNN Language Models

RNN Advantages:

- Can process any length input ٠
- Computation for step t can (in ٠ theory) use information from many steps back
- Model size doesn't increase for longer input ۲ context
- Same weights applied on every ٠ timestep, so there is symmetry in how inputs are processed.

RNN Disadvantages:

- Recurrent computation is slow
- In practice, difficult to access information from many steps back





 $oldsymbol{h}^{(0)}$

 \bigcirc

 \bigcirc





Training an RNN Language Model

Training an RNN Language Model

- Get a big corpus of text which is a sequence of words $x^{(1)}, x^{(2)}, \dots, x^{(T)}$
- Feed into RNN-LM; compute output distribution $\hat{y}^{(t)}$ for every step *t*.
 - i.e., predict probability distribution of every word, given words so far
- Loss function on step t is cross-entropy between predicted probability distribution $\hat{y}^{(t)}$, and the true next word $y^{(t)}$ (one-hot for $x^{(t+1)}$):

$$J^{(t)}(\theta) = CE(\boldsymbol{y}^{(t)}, \hat{\boldsymbol{y}}^{(t)}) = -\sum_{w \in V} \boldsymbol{y}_w^{(t)} \log \hat{\boldsymbol{y}}_w^{(t)} = -\log \hat{\boldsymbol{y}}_{\boldsymbol{x}_{t+1}}^{(t)}$$

• Average this to get overall loss for entire training set:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^{T} -\log \hat{y}_{\boldsymbol{x}_{t+1}}^{(t)}$$















Training a RNN Language Model

• However: Computing loss and gradients across entire corpus $x^{(1)}, x^{(2)}, ..., x^{(T)}$ at once is too expensive (memory-wise)!

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} J^{(t)}(\theta)$$

- In practice, consider $x^{(1)}, x^{(2)}, \dots, x^{(T)}$ as a sentence (or a document)
- <u>Recall: Stochastic Gradient Descent</u> allows us to compute loss and gradients for small chunk of data, and update.
- Compute loss $J(\theta)$ for a sentence (actually, a batch of sentences), compute gradients and update weights. Repeat on a new batch of sentences.







Question: What's the derivative of $J^{(t)}(\theta)$ wint the repeated weight matrix W_h ?

Answer:
$$\frac{\partial J^{(t)}}{\partial W_h} = \sum_{i=1}^t \frac{\partial J^{(t)}}{\partial W_h}\Big|_{(i)}$$

"The gradient w.r.t. a repeated weight is the sum of the gradient w.r.t. each time it appears"

Why?







Multivariable Chain Rule

- Given a multivariable function f(x,y), and two single variable functions x(t) and y(t), here's what the multivariable chain rule says:

$$\underbrace{\frac{d}{dt}f(\pmb{x}(t),\pmb{y}(t))}_{dt} = \frac{\partial f}{\partial \pmb{x}}\frac{d\pmb{x}}{dt} + \frac{\partial f}{\partial \pmb{y}}\frac{d\pmb{y}}{dt}$$

Derivative of composition function





Source:

https://www.khanacademy.org/math/multivariable-calculus/multivariablederivatives/differentiating-vector-valued-functions/a/multivariable-chain-rulesimple-version







Training The Parameters of RNNs: Backpropagation for RNNs



Question: How do we calculate this?

Answer: Backpropagate over timesteps i = t, ..., 0, summing gradients as you go. This algorithm is called **"backpropagation through time"** Apply the multivariable chain rule:

= 1

$$\frac{\partial J^{(t)}}{\partial \boldsymbol{W}_{h}} = \sum_{i=1}^{t} \frac{\partial J^{(t)}}{\partial \boldsymbol{W}_{h}} \Big|_{(i)} \frac{\partial \boldsymbol{W}_{h} \Big|_{(i)}}{\partial \boldsymbol{W}_{h}}$$

$$=\sum_{i=1}^{t} \frac{\partial J^{(t)}}{\partial \boldsymbol{W}_{h}}\Big|_{(i)}$$