



NEXT WORD PREDICTION

¹Keerthana N, ²Harikrishnan S, ³Konsaha Buji M, ⁴Jona J B

¹Student, ²Student, ³Student, ⁴Associate Professor

¹Department of Computer Applications,

¹Coimbatore Institute of Technology, Coimbatore, India

Abstract: Writing long sentences is bit boring, however with text prediction within the keyboard technology has created this easy. Next Word Prediction is in addition referred to as Language Modeling. It's the endeavor of predicting what word comes straightaway. It's one in every of the key assignments of human language technology and has various applications. Long short time memory formula can perceive past text and predict the words which can be useful for the user to border sentences and this method uses letter to letter prediction suggests that it predict a letter when letter to form a word.

Index Terms - NLP, LSTM, RNN, Next Word.

I. INTRODUCTION

Word prediction tools were developed which might facilitate to speak and additionally to assist the individuals with less speed writing. during this paper, a language model based mostly framework for fast electronic communication, which will predict probable next word given a group of current words are briefed. Word prediction technique will the task of guesswork the preceding word that's probably to continue with few initial text fragments. Our goal is to facilitate the task of instant electronic communication by suggesting relevant words to the user.

II. LITERATURE REVIEW

Existing systems work on word prediction model, that suggests subsequent immediate word supported this out their word[2]. These systems work victimization machine learning algorithms that has limitation to form correct syntax. Multi-window convolution (MRN N) formula is enforced, additionally they need created residual-connected lowest gated unit(MGU) that is brief version of LSTM during this cnn try and skip few layers whereas coaching end in less coaching time and that they have sensible accuracy out and away victimization multiple layers of neural networks will cause latency for predicting n numbers of words .Developing technologies has been manufacturing additional correct outcomes than the prevailing system technologies, models developed victimization bidirectional LSTM algorithms area unit capable of handling additional knowledge expeditiously and predicts higher[2].

III. THE ALGORITHM

3.1 LSTM

Vanishing gradient descent may be a downside featured by neural networks once considering back propagation. It's Brobbingnag Ian impact and also the weight update method is wide affected and also the model became useless. thus LSTM that incorporates a hidden state and a memory cell with 3 gates that area unit forgotten, scan and input gate.

The forget gate is principally accustomed get sensible management of what data has to be removed that isn't necessary. Input gate makes positive that newer data is additional to the cell and output makes positive what elements of the cell area unit output to subsequent hidden state. The sigmoid operate utilized in every gate equation makes positive we will bring down the worth to either a zero or one.

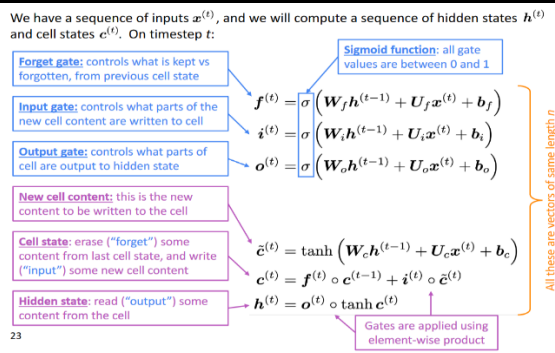


Fig -1 LSTM Architecture

The below Fig-2 is the Illustration of basic LSTM prediction model for our downside. The model encodes the input word sequence describing previous sub events into associate degree embedding and decodes a word sequence describing the doable future sub event out of the embedding.

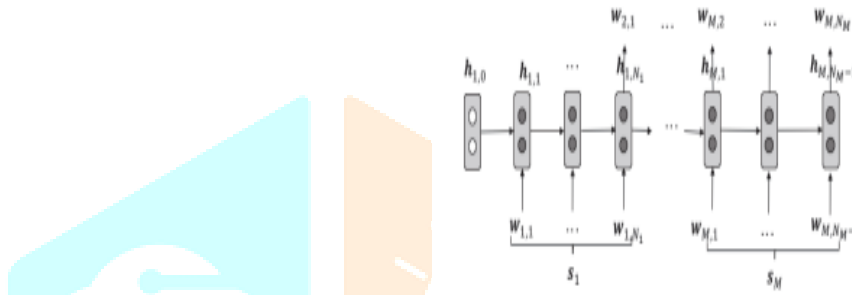


Fig-2 LSTM Prediction Model

3.2 RNN

Continual neural networks area unit generalizations of an instantaneous transmission neural network that has internal memory. The RNN is repetitive in nature as a result of it performs constant operate for every knowledge entry, however at constant time, this output depends on the previous calculation. the choice relies on associate degree analysis of this input and also the output from the previous input. RNNs will use their internal state (memory) to method input sequences once direct communication neural networks cannot. All RNN inputs area unit interconnected.

This state formula is that the following :-

$$h_t = f(h_{t-1}, x_t)$$

Application of the activation operate :-

$$H_t = \tanh(W_{hh}h_{t-1} + W_{hx}x_t)$$

IV. PROPOSED METHODOLOGY

4.1 Data Preprocessing

These area unit easy clean-up procedures that makes it easier to use the information in sequent steps. This method is administered with the assistance of Tensor flow library.

The subsequent area unit few pre-processing steps typically done:-

1. Marking white areas
2. Lower-case conversions
3. Removing numbers
4. Removing punctuation
5. Removing unwanted words
6. Removing non-English words

4.2 Text Analysis

To know the speed of occurrence of terms, Term Document Matrix operate was accustomed to produce term matrixes to achieve the account of term frequencies.

4.3 Tokenization

One in every of the vital social control strategies is named tokenization. It's merely segmenting the continual running text into individual segments of words. One terribly easy approach would be to separate inputs over each house associate degree assign an identifier to every word.

4.4 Pad Sequences

When changing sentences to numerical values, there's still a difficulty of providing equal length inputs to our neural networks. Not each sentence is constant length. pad_sequences operate is employed for artefact the shorter sentences with zeroes, and truncating a number of the longer sequences to be shorter.

Additionally, because it is often specified whether or not to pad and truncate from either the start or ending, relying upon the pre settings and post settings for the padding arguments and truncating arguments. By default, artefact arguments and truncation can happen from the start of the sequence.

V. IMPLEMENTATION AND RESULT

The implementation and results obtained for this project are shortly delineate during this section.

```
In [9]: print(tokenizer.word_index)

{'the': 1, 'and': 2, 'i': 3, 'to': 4, 'a': 5, 'of': 6, 'my': 7, 'in': 8, 'he': 9, 'for': 10, 'you': 11, 'all': 12, 'was': 13, 'she': 14, 'that': 15, 'on': 16, 'with': 17, 'he': 18, 'but': 19, 'as': 20, 'when': 21, 'love': 22, 'is': 23, 'you': 24, 'i': 25, 'will': 26, 'from': 27, 'by': 28, 'they': 29, 'be': 30, 'me': 31, 'so': 32, 'he': 33, 'did': 34, 'no': 35, 'oh': 36, 'ill': 37, 'at': 38, 'one': 39, 'his': 40, 'there': 41, 'were': 42, 'here': 43, 'down': 44, 'now': 45, 'we': 46, 'where': 47, 'young': 48, 'never': 49, 'go': 50, 'come': 51, 'then': 52, 'did': 53, 'not': 54, 'said': 55, 'way': 56, 'their': 57, 'sw': 58, 'can': 59, 'green': 60, 'if': 61, 'take': 62, 'am': 63, 'like': 64, 'right': 65, 'day': 66, 'o': 67, 'out': 68, 'fair': 69, 'this': 70, 'two': 71, 'have': 72, 'can': 73, 'true': 74, 'its': 75, 'how': 76, 'see': 77, 'dear': 78, 'more': 79, 'there's': 80, 'or': 81, 'had': 82, 'would': 83, 'over': 84, 'hear': 85, 'up': 86, 'ie': 87, 'through': 88, 'none': 89, 'again': 90, 'well': 91, 'see': 92, 'and': 93, 'good': 94, 'in': 95, 'ye': 96, 'see': 97, 'left': 98, 'still': 99, 'father': 100, 'long': 101, 'rose': 102, 'could': 103, 'morning': 104, 'wild': 105, 'who': 106, 'eyes': 107, 'came': 108, 'while': 109, 'too': 110, 'back': 111, 'little': 112, 'am': 113, 'took': 114, 'him': 115, 'bow': 116, 'first': 117, 'let': 118, 'man': 119, 'shall': 120, 'know': 121, 'get': 122, 'high': 123, 'game': 124, 'say': 125, 'ever': 126, 'some': 127, 'mary': 128, 'hand': 129, 'till': 130, 'out': 131, 'am': 132, 'time': 133, 'heard': 134, 'dead': 135, 'way': 136, 'bright': 137, 'mountain': 138, 'early': 139, 'rosin': 140, 'gave': 141, 'thee': 142, 'only': 143, 'far': 144, 'maid': 145, 'must': 146, 'find': 147, 'girl': 148, 'sure': 149, 'round': 150, 'dublin': 151, 'once': 152, 'world': 153, 'delight': 154, 'last': 155, 'johnny': 156, 'seen': 157, 'has': 158, 'fine': 159, 'road': 160, 'mother': 161, 'tis': 162, 'what': 163, 'way': 164, 'moon': 165, 'soul': 166, 'see': 167, 'id': 168, 'just': 169, 'that's': 170, 'days': 171, 'darling': 172, 'went': 173, 'white': 174, 'die': 175, 'than': 176, 'hair': 177, 'gus': 178, 'neat': 179, 'today': 180, 'do': 181, 'girls': 182, 'shes': 183, 'thyme': 184, 'thy': 185, 'sin': 186, 'pretty': 187, 'new': 188, 'poor': 189, 'into': 190, 'life': 191, 'irish': 192, 'give': 193, 'boy': 194, 'youre': 195}
```

Fig – 3 Tokenization

In Fig-3 the full dataset is separated into individual phase of words. When segmenting a word index is formed with specified distinctive variety of words in irish-lyrics-eof.txt by distribution fact to every of them.

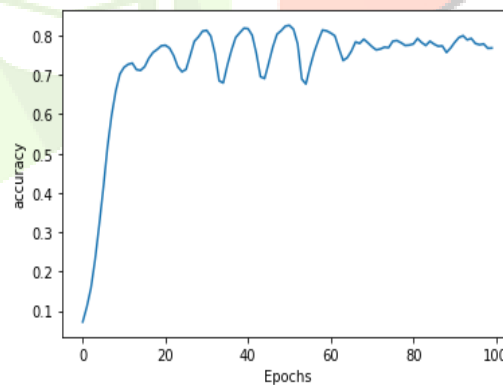


Fig – 4 Accuracy

In Fig-4 when training the model, the accuracy is found to be 0.80 in a hundred epochs.

```

In [13]: seed_text = "I see a beautiful river"
         next_words = int(input("Word Count : "))

         for _ in range(next_words):
             token_list = tokenizer.texts_to_sequences([seed_text])[0]
             token_list = pad_sequences([token_list], maxlen=max_sequence_len-1, padding='pre')
             predicted = np.argmax(model.predict(token_list), axis=-1)
             output_word = ""
             for word, index in tokenizer.word_index.items():
                 if index == predicted:
                     output_word = word
                     break
             seed_text += " " + output_word
         print(seed_text)

Word Count : 2
I see a beautiful river my love

```

Fig – 5 Output

In Fig-5 associate degree input sentence “I see a stunning river” is given. variety of words to be expected when this words is got as input from the user. Because the word count is given as two “my love” is that the expected words.

VI. CONCLUSION

The subsequent word prediction model that was developed is fairly correct on the provided dataset. NLP requires applying various types of pattern discovery approaches aimed at eliminating noisy data. The loss was considerably reduced in concerning a hundred epochs. Files or dataset that are large to process need still some optimizations. However, bound pre-processing steps and bound changes within the model are often created to boost the prediction of the model.

REFERENCES

- [1] Bengio, Y., Simard, P., Frasconi, P., 1994. Learning long dependencies with gradient descent is troublesome. *IEEE transactions on neural networks* five, 157–166.
- [2] R. Kneser and H. Ney, "Improved backing-off for n-gram language modeling", *Conference on Acoustics speech and Signal Process*, pp. 181-184, 1995.
- [3] Mohd. Majid and Piyush Kumar, *Language Modelling: Next word Prediction*, 2019.
- [4] Bengio, Y., Simard, P., Frasconi, P., 1994., Learning NLP with gradient descent is troublesome. *IEEE transactions on neural networks* five, 157–166.
- [5] Serban, I. V.; Sordoni, A.; Bengio, Y.; Courville, A.; and Pineau, J. 2016. Building end-to-end dialogue mistreatments generative stratified neural network models. In *Proceedings of the 30th Conference on Artificial Intelligence. AAAI*.
- [6] J. Yang, H. Wang and K. Guo, "Natural language Word Prediction Model supported Multi-Window Convolution and Residual Network," in *IEEE Access*, vol. 8, pp. 188036-188043, 2020, doi: 10.1109/ACCESS.2020.3031200.
- [7] M. K. Sharma, S. Sarcar, P. K. Saha and D.Samanta, "Visual clue: Associate approach to predict and highlight next character," 2012 fourth International Conference on Intelligent Human pc Interaction (IHCI), Kharagpur, 2012, pp. 1-7, doi:10.1109 /IHCI.2012.6481820.
- [8] Sukhbaatar, S., Weston, J., Fergus, R., et al., 2015. End-to-end memory networks, in: *Advances in neural information processing systems*, pp. 2440–2448.
- [9] Zhou, C., Sun, C., Liu, Z., Lau, F., 2015. A c-lstm neural network for text classification. *arXiv preprint arXiv:1511.08630*.
- [10] Joel Stremmel, Arjun Singh. (2020). Pretraining Federated Text Models for Next Word Prediction using GPT2.
- [11] S. M. Sarwar and Abdullah-Al-Mamun, "Next word prediction for phonetic typing by grouping language models," 2016 2nd International Conference on Information Management (ICIM), London, 2016, pp. 73-76, doi: 10.1109/ INFOMAN.7477536.
- [12] J. Yang, H. Wang and K. Guo, "Natural Language Word Prediction Model Based on Multi-Window Convolution and Residual Network," in *IEEE Access*, vol. 8, pp.188036-188043, 2020, doi: 10.1109 / ACCESS.20202.3031200.