



# INTERNATIONAL JOURNAL OF CREATIVE RESEARCH THOUGHTS (IJCRT)

An International Open Access, Peer-reviewed, Refereed Journal

## Sentence Completion Using NLP Techniques

<sup>1</sup>Tanushree C Hatmode, <sup>2</sup>Prachi Y Duragkar, <sup>3</sup>Radha H Khorgade, <sup>4</sup>Ram B Jaiswal, <sup>5</sup>Shubhashish S Charan, <sup>6</sup>Rajesh Nasare, <sup>7</sup>Hemant Turkar

<sup>1</sup>Student, <sup>2</sup>Student, <sup>3</sup>Student, <sup>4</sup>Student, <sup>5</sup>Student, <sup>7</sup>Assistant Professor, <sup>8</sup>Assistant Professor  
<sup>1</sup>Computer Science engineering,

<sup>1</sup>Rajiv Gandhi College Of Engineering And Research, Nagpur, India

**Abstract:** In this assignment, we plan to use and examine diverse strategies for automatic sentence finishing touch which consist of n-gram modeling, Latent Semantic Indexing and Recurrent Neural Networks. Results from the beyond studies reveal that the LSA version outperforms the traditional n-gram version fashions at the Microsoft Research Sentence Completion Challenge. The LSA version outperforms different non-neural community fashions. Hence, we are able to examine the consequences of the LSA version with RNN version and n-gram version to determine which ones play higher thinking about the size of the statistics. The strategies are first teach on a big corpus of un-annotated textual content, to then attempt to are expecting the lacking phrases within side the check set which includes hundreds of sentences wherein one phrase is lacking and 4 options for the lacking phrase.

**Index Terms -** n-nlp , gram, rnn , lsa

### I. INTRODUCTION

The Sentence Completion performs a pivotal position in English language and communication. It triggers cognitive abilities to interpret and examine the written sentences/phrases. Sentence Completion assists in green propagation of thoughts and clear interplay with every different.

In current years, standardized examinations have proved a fertile supply of assessment statistics for language processing obligations. They are precious for plenty reasons: they constitute sides of language knowledge diagnosed as essential with the aid of using academic experts; they may be prepared in diverse codes designed to assess precise capabilities; they may be yardsticks with the aid of using which society measures academic development and that they have an effect on a big variety of humans. The hassle of Sentence finishing touch is to degree grammatical and semantic correctness and to pick the maximum suitable sentence/phrase. These assessments the capacity of algorithms to differentiate feel from nonsense primarily based totally on a whole lot of sentence-degree phenomena. For every sentence the assignment is to decide which of the alternatives for that phrase the appropriate one is. To enhance accuracy for maximum variety of sentences finishing touch troubles. It must be correct sufficient to offer solutions higher than N-gram fashions. We are going to apply unique NLP strategies to do the same.

In this paper, we are able to check out strategies for answering sentence finishing touch questions. Also we are able to be locating the accuracy of predicting suitable phrase for given SAT fashion sentences. Finally an standard answer that is without problems implementable for completely operating utility functions is evolved.

#### 1.1 Objective

We plan to observe and enforce computational strategies (Algorithms) to condemn finishing touch and strive unique mixtures to boom the accuracy and the use of the maximum correct set of rules to construct an utility.

- BACK OFF N-GRAM LANGUAGE MODEL
- RECURRENT NEURAL NET LANGUAGE MODEL (RNN)
- SENTENCE COMPLETION THROUGH LATENT SEMANTIC ANALYSIS (LSA).

## II. RELATED WORK

The beyond paintings that is maximum much like ours is derived from Microsoft Research Sentence Completion Challenge (G. Zweig, Christopher J. C. Burges, 2011) wherein Microsoft has proposed a fixed of diverse English sentences. [6] Each sentence has been related to the impostor alternatives, wherein every phrase within side the unique sentence is changed with the aid of using an impostor phrase with comparable prevalence statistics. The assignment for every sentence is to decide out of alternatives that is the appropriate choice for the given sentence. This assignment is much like the SAT language check. The query changed into generated in steps. First, the candidate sentence which includes an rare phrase changed into decided on and the opportunity phrase changed into decided robotically with the n-gram language version with the aid of using sampling. The n-gram version used intermediate records as lexicon, which resulted within side the phrases which can be suitable locally, however there may be no different purpose to count on them to make it feel globally. In the second one step, apparent wrong alternatives are removed due to the fact they contained a few grammatical errors. Data used changed into from the 5 of Conan Doyle Sherlock Holmes novels: The Sign of the Four (1890), The Adventures of Sherlock Holmes (1892), The Hound of the Baskervilles (1892), The Memoirs of Sherlock Holmes (1894), and The Valley of Fear (1915). N-gram version changed into educated on 540 texts consisting particularly of 19th-Century Novels.

### 2.1 Procedure :

- Tokenization of sentence takes area because of this that sentence is split into tokens(phrases), this tokens are represented with the aid of using vectors.
- Training the use of every token(phrase) with the aid of using unique hidden layers and shape a version which has the relativity of phrases used for schooling. (first layer of blue balls is hidden layer at T0, 2d at T1 and so on)
- A version is fashioned that is used for checking out, thinking about the instance Ram used to hate this with alternatives a) sequential b) decorrundum c) Moved

2.1.1 Recurrent Neural Network : RNNs are very apt for collection class troubles and the purpose they're so properly at that is that they're capable of hold essential statistics from the preceding inputs and use that statistics to regulate the modern output. If the sequences are pretty lengthy, the gradients (values calculated to music the community) computed at some stage in their schooling (back propagation) both vanish (multiplication of many  $0 < \text{values} < 1$ ) or explode (multiplication of many big values) inflicting it to teach very slowly. Long Short Term Memory is a RNN structure which addresses the hassle of schooling over lengthy sequences and preserving memory. LSTMs clear up the gradient hassle with the aid of using introducing some extra gates that manipulate get admission to to the mobileular state. The device version this is to be applied for sentence finishing touch is defined below:

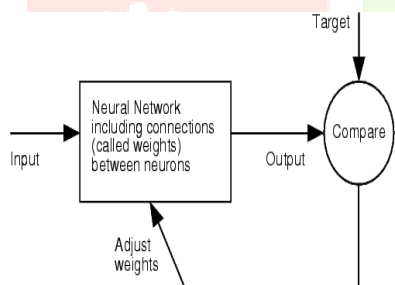


Fig. II.1: Basic Architecture for the use of RNN for sentence finishing touch

2.1.2 Latent Semantic Analysis (LSA) : The approach used to extract and constitute the contextual that means of phrases with the aid of using the use of statistical computations that is carried out to a big corpus of textual content. LSA is an automated statistical method used to extract and infer family members of contextual utilization of phrases in passages. It isn't always any traditional herbal language processing; it makes use of no understanding bases, semantic networks, humanly built dictionaries, syntactic parsers, it takes enter as a uncooked textual content parsed into uniquely person strings which separates into significant passages which include sentences. It takes into consideration the distributional speculation which states that phrases which can be near in that means will arise in comparable portions of textual content.

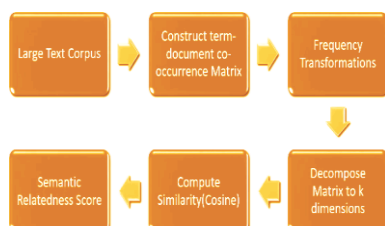


Figure II.2: Steps in LSA

### III. IMPLEMENTATION

- REMOVING TAB AREAS .
- DIVIDING THE SENTENCES THE USE OF LINE STRIP
- REMOVING VARIETY/DIGITS.
- MAKING THE WHOLE LOT INTO LOWERCASE.
- ADDING EACH LINE FROM THE TEXTUAL CONTENT .
- JOIN THE WHOLE LOT WITHIN SIDE THE LIST.
- REPLACE CONSECUTIVE AREAS WITH UNMARRIED AREAS .
- STORE THE BRAND NEW TEXTUAL CONTENT IN ABC

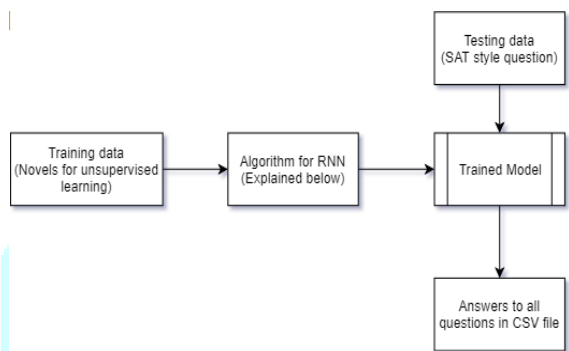


Fig III.1 : Proposed model for implementation

A) Dataset: The dataset used for Sentence Completion is the Microsoft Research Sentence Completion Challenge dataset. This dataset evolved with the aid of using Microsoft Research changed into publicly made to be had in 2012 (G. Zweig, C. Burges, 2012) so that you can enhance the studies within side the discipline of semantic and syntactic analysis. The dataset is primarily based totally on 5 of Conan Doyle Sherlock Holmes books that have been part of Project Gutenberg. We have used Holmes Training Data to teach our version. It has 1040 sentence finishing touch obligations primarily based totally at the SAT query format. Each assignment has 5 alternatives out of which simplest one is accurate however all 5 match probably nicely into the sentence. The schooling set carries 522 books from Project Gutenberg open corpus, every of them having good enough headers. We have a checking out dataset which incorporates 1040 questions together with their 5 unique alternatives. Princeton evaluate dataset with eleven exercise check for SAT is used. Testing dataset is used to look how nicely the device is educated with the intention to solution unique unseen questions .

B) Preprocessing: To follow the unique fashions to the Sentence Completion, preprocessing of dataset is required. Pre-processing in completed in 3 steps, which contain elimination of tab areas, changing the whole lot to lowercase and elimination of punctuation and prevent phrases.

C) Interface: A Desktop utility is created for sentence finishing touch for the consumer base - students, schools, faculties and exam center. There are particularly elements within side the interface. The first one being Upload the CSV file wherein the consumer uploads the checking out questions with 5 alternatives every and the second one. Ask a query wherein consumer can get a solution for the unmarried query from the 5 to be had alternatives. Both those obligations may be accomplished the use of RNN , LSA and both (RNN and LSA). Hence accuracy of those algorithms may be in comparison and evaluated. All the fashions were applied in Python3 the use of Jupyter Notebooks and Collaborator with the aid of using Google Research. We used Python libraries like NLTK, Gensim, SciPy, Numpy and Tensor flow. For RNN, we first ran the code that carries all of the features required for schooling the files. Then, we definitely applied the RNN version. Training changed into executed to extract the checkpoints. Then we executed checking out with the aid of using imparting enter checking out statistics to are expecting the proper solutions, that is then saved within side the separate file. For LSA, we first applied a bag of phrases version to the corpus observed with the aid of using the phrase vectors. Then, we computed the cosine similarity among the candidate solution and the query statement. Finally, we computed common cosine similarity for all of the capabilities and lower back the choice with the best common cosine similarity.

#### IV. RESULT

All the fashions were applied in Python3 the use of Jupiter Notebooks and Collaborator with the aid of using Google Research. We used Python libraries like NLTK, Genesis, SciPy, Numpy and Tensor flow. For RNN, we first ran the code that carries all of the features required for schooling the document case. Then, we definitely applied the RNN version. Training changed into executed to extract the checkpoints. Then we executed checking out with the aid of using imparting enter checking out statistics to are expecting the proper solutions, that is then saved within side the separate file. For LSA, we first applied a bag of phrases version to the corpus observed with the aid of using the phrase vectors. Then, we computed the cosine similarity among the candidate solution and the query statement. Finally, we computed common cosin similarity for all of the capabilities and lower back the choice with the best common cosine similarity.

- The accuracy of RNN version is 42.59.
- The accuracy of LSA Algorithm is 53.75
- Implementation of LSA The accuracy of Latent Semantic Analysis set of rules is extra than that of Recurrent Neural Network set of rules.

#### V. CONCLUSION

The assignment examines diverse strategies in order to are expecting the proper solutions for the given judgment- finishing touch questions. These questions are attractive due to the fact they probe the capacity to differentiate semantically coherent decision from incoherent ones, and but contain no extra context than the unmarried sentence. We have solved the hassle with the aid of using imposing the strategies of namely, recurrent neural networks and latent semantic analysis. The accuracy accomplished with the aid of using lsa changed into appreciably extra than RNN.

#### VI. ACKNOWLEDGMENT

I would like to thank to Prof. Rajesh Nasare, Prof Hemant Turkar ,Assistant Professor of computer science engineering Department from Rajiv Gandhi college of Engineering and Research, Nagpur, India for their guidance

#### VII. REFERENCES

- [1] Joseph Gubbins, Andreas Vlachos. Dependency Language fashions for sentence finishing touch, 2013.
- [2] Geoffrey Zweig, John C. Platt, Christopher Meek, Christopher J. C. Burges, Ainur Yessenalina, Qiang Liu. Computational Approaches to Sentence Completion, 2012.
- [3] Thomas K Landauer, Peter W. Foltz, Darrell Laham. An Introduction to Latent Semantic Analysis, 1998.
- [4] Aubrie M. Woods. Carnegie Mellon University. Exploiting Linguistic Features for Sentence Completion, 2016.
- [5] Tomas Mikolov, Kai Chen, Greg Corrado, Jeffrey Dean. Efficient Estimation of Word Representations in VectorSpace, 2013
- [6] Geoffrey Zweig and Christopher J.C. Burges.