



AI BASED PAPER AUTOMATIC EXAMINATION PAPER EVALUATION SYSTEM

Dr. A.LAXMIKANTH, Associate Professor,

M. RAMYA, R. JAHNAVI, N. SAI KIRAN, T. RAKESH

SRI INDU COLLEGE OF ENGINEERING AND TECHNOLOGY

Sheriguda (V), Ibrahimpatnam (M), RangareddyDist – 501 510

ABSTRACT

Subjective paper evaluation is a tricky and tiresome task to do by manual labor. Insufficient understanding and acceptance of data are crucial challenges while analyzing subjective papers using Artificial Intelligence (AI). Several attempts have been made to score students' answers using computer science. However, most of the work uses traditional counts or specific words to achieve this task. Furthermore, there is a lack of curated data sets as well. This paper proposes a novel approach that utilizes various machine learning, natural language processing techniques, and tools such as Wordnet, Word2vec, word mover's distance (WMD), cosine similarity, multinomial naive bayes (MNB), and term frequency-inverse document frequency (TF-IDF) to evaluate descriptive answers automatically. Solution statements and keywords are used to evaluate answers, and a machine learning model is trained to predict the grades of answers. Results show that WMD performs better than cosine similarity overall. With enough training, the machine learning model could be used as a standalone as well. Experimentation produces an accuracy of 88% without the MNB model. The error rate is further reduced by 1.3% using MNB

INDEX : manual labor, si, wordnet, word2vec, tf-idf

1. INTRODUCTION

1.1 Introduction:

Subjective questions and answers can assess the performance and ability of a student in an open-ended manner. The answers, naturally, are not bound to any constraint, and students are free to write them according to their mindset and understanding of the concept. With that said, several other vital differences separate subjective answers from their objective counterpart. For one, they are much longer than the objective questions. Secondly, they take more time to write. Moreover, they carry much more context and take a lot of concentration and objectivity from the teacher evaluating them. Evaluation of such questions using computers is a tricky task, mainly because natural language is ambiguous. Several preprocessing steps must be performed, such as cleaning the data and tokenization before working on it. Then the textual data can be compared using various techniques such as document similarity, latent semantic structures, concept graphs, ontologies. The final score can be evaluated based on Similarity, keywords presence, structure, language. Several attempts have been made in the past to solve this problem, but there is still room for improvements, some of which is discussed in this paper. Subjective exams are considered more complex and scarier by both students and teachers due to their one fundamental feature,



context. A subjective answer demands the checker check every word of the answer for scoring actively, and the checker's mental health, fatigue, and objectivity play a massive role in the overall result. Therefore, it is much more time and resource-efficient to let a system handle this tedious and somewhat critical task of evaluating subjective answers. Evaluating objective answers with machines is very easy and feasible. A program can be fed with questions and one-word answers that can quickly map students' responses. Nevertheless, subjective answers are much more challenging to tackle. They are varied in length and contain a vast amount of vocabulary. Furthermore, people tend to use synonyms and convenient abbreviations, which makes the process that much trickier.

2. LITERATURE SURVEY

TITLE: "Measurement of text similarity: A survey," Information

ABSTRACT: Text similarity measurement is the basis of natural language processing tasks, which play an important role in information retrieval, automatic question answering, machine translation, dialogue systems, and document matching. This paper systematically combs the research status of similarity measurement, analyzes the advantages and disadvantages of current methods, develops a more comprehensive classification description system of text similarity measurement algorithms, and summarizes the future development direction. With the aim of providing reference for related research and application, the text similarity measurement method is described by two aspects: text distance and text representation. The text distance can be divided into length distance, distribution distance, and semantic distance; text representation is divided into string-based, corpus-based, single-semantic text, multi-semantic text, and graph-structure-based representation. Finally, the development of text similarity is also summarized in the discussion section.

TITLE: "A survey on the techniques, applications, and performance of short text semantics similarity,"

ABSTRACT: Short text similarity plays an important role in natural language processing (NLP). It has been applied in many fields. Due to the lack of sufficient context in the short text, it is difficult to measure the similarity. The use of semantics similarity to calculate textual similarity has attracted the attention of academia and industry and achieved better results. In this survey, we have conducted a comprehensive and systematic analysis of semantic similarity. We first propose three categories of semantic similarity: corpus-based, knowledge-based, and deep learning (DL)-based. We analyze the pros and cons of representative and novel algorithms in each category. Our analysis also includes the applications of these similarity measurement methods in other areas of NLP. We then evaluate state-of-the-art DL methods on four common datasets, which proved that DL-based can better solve the challenges of the short text similarity, such as sparsity and complexity. Especially, bidirectional encoder representations from transformer model can fully employ scarce information of short texts and semantic information and obtain higher accuracy and F1 value. We finally put forward some future directions.



TITLE: “Subjective answer evaluation using machine learning,”

ABSTRACT: This project proposes a novel approach that utilizes various machine learning, natural language processing techniques, to evaluate descriptive answers automatically. Solution statements and keywords are used to evaluate answers, and a machine learning model is trained to predict the grades of answers. With enough training, the machine learning model could be used as a standalone as well. Experimentation produces an accuracy of 97% with the Proposed model. Interestingly, artificial intelligence is utilized extensively as an efficient tool for predicting such a problem. The proposed work utilizes the deep learning technique along with some preprocessing steps to improve the prediction of Answer Evaluation.

TITLE: “Automated assessment system for subjective questions based on LSI,”

ABSTRACT: Subjective question is capable of examining the adopting ability of knowledge of the student, but the assessment for it suffers from a number of questions such as trickiness, synonymy and polysemy. This reduces the advantage of subjective question for online exercise. In this paper we explore an approach to automated assessment system for subjective question based on latent semantic indexing. Chinese automatic segmentation techniques and subject ontology are used for transferring the reference answers to a term-document matrix, which is then projected to a k-dimensional LSI space by the statistical technique Singular Value Decomposition to solve the problem of synonymy and polysemy. A reference unit vector is introduced to alleviate the problem of trickiness. The system then concludes the quality of the solution according to the similarity between the projected vectors. The experimental results prove the feasibility of our theoretical architecture and flow for automated assessment of subjective question.

TITLE: “From word embeddings to document distances,”

ABSTRACT: We present the Word Mover’s Distance (WMD), a novel distance function between text documents. Our work is based on recent results in word embeddings that learn semantically meaningful representations for words from local co-occurrences in sentences. The WMD distance measures the dissimilarity between two text documents as the minimum amount of distance that the embedded words of one document need to “travel” to reach the embedded words of another document. We show that this distance metric can be cast as an instance of the Earth Mover’s Distance, a well studied transportation problem for which several highly efficient solvers have been developed. Our metric has no hyperparameters and is straight-forward to implement. Further, we demonstrate on eight real world document classification datasets, in comparison with seven state-of-the-art baselines, that the WMD metric leads to unprecedented low k-nearest neighbor document classification error rates.

TITLE: “Similarity analysis of law documents based on Word2vec,”

ABSTRACT: With the increasing demand for computer-assisted wisdom in justice, deep learning has gradually become an effective means of helping intelligent justice. The similarity analysis of law documents is the basis of intelligent justice, while law documents based on several types of cases are quite different in terms of format and length, which causes trouble in analyzing similarities. For that



we propose a more specific approach to dealing with law documents, combining Word2vec with legal documents corpus. To measure the efficiency of the proposed method, we designed two sets of controls. The experimental results show that the Word2vec model can improve the accuracy by 0.20 compared with the bag of words (BOW) model, and the equipped law documents corpus can increase by 0.05-0.10 based on the Word2vec model. Thus, the combination of Word2vec and the law documents corpus is more compatible with the simple and efficient application of similarity analysis of law documents.

3. PROBLEM STATEMENT

It appears there might be some confusion. The information you provided in your previous message seems to be more focused on the proposed system and its objectives rather than describing the existing system. The existing system typically refers to the state of affairs or methodologies in place before the implementation of the proposed system. If you have information about the existing system, you could provide details on how the assessment of descriptive answers is currently handled, whether it's manual evaluation by teachers or any existing tools or methods in use.

LIMITATIONS OF SYSTEM:

Subjectivity and Bias: Issue: Manual evaluation can be subjective, leading to variations in grading among different evaluators. **Impact:** Inconsistencies in grading may result in unfair assessments and disparities in students' grades. **Time-Consuming:** Issue: Manual grading of descriptive answers is a time-consuming process, especially in scenarios with a large number of students or complex questions. **Impact:** Teachers may face challenges in providing timely feedback to students, and the overall assessment process may be delayed. **Scalability Challenges:** Issue: As the number of students and assessments increases, scalability becomes a significant challenge for manual evaluation. **Impact:** Educational institutions may struggle to efficiently manage and scale the assessment process, particularly during peak times.

4. PROPOSED SYSTEM & IT'S ADVANTAGES:

The proposed system, "A Descriptive Answer Evaluation System Using Cosine Similarity Technique," offers a transformative approach to address the limitations of traditional manual evaluation methods for descriptive answers. The primary objective is to leverage computer-assisted assessment tools, particularly in the context of the evolving challenges posed by the COVID-19 pandemic. The system aims to alleviate the subjectivity and bias inherent in manual grading by introducing an automated evaluation process based on the cosine similarity technique. This method allows for a more objective assessment of descriptive answers, irrespective of their length, enabling a fairer and more consistent grading system. One of the key features of the proposed system is its ability to significantly reduce the time required for assessment. By automating the evaluation process, teachers can allocate more time to providing detailed and timely feedback to students. The scalability of the system addresses the challenges associated with handling a large number of assessments efficiently. This shift towards a computer-assisted solution not only streamlines the evaluation process but also minimizes the resource-intensive nature of manual grading, potentially leading to



costsavings for educational institutions. Moreover, the proposed system enhances the feedback loop for students by providing pictorial representations of the results using the cosine similarity technique. This visual representation not only facilitates a quick understanding of the assessment outcome but also serves as a valuable learning aid. The web-based application aspect further modernizes the assessment approach, making it more adaptable to the digital learning environment. In essence, the proposed system strives to revolutionize the assessment of descriptive answers, making it more objective, efficient, and accessible in the contemporary educational landscape.

4.1 ADVANTAGES:

Objective Assessment: The system employs the cosine similarity technique, a quantitative measure that provides an objective evaluation of descriptive answers. This helps eliminate subjective biases often associated with manual grading, ensuring fairness and consistency in assessments.

Time Efficiency: Automated evaluation significantly reduces the time required for grading descriptive answers. This efficiency benefits both teachers and students by expediting the feedback process, allowing for quicker identification of areas of improvement and enhancing the overall learning experience.

Scalability: The proposed system is designed to handle a large volume of assessments efficiently. As the number of students and evaluations increases, the automated approach ensures scalability, addressing the challenges posed by manual grading in terms of time and resource constraints.

Enhanced Feedback: The system provides visual representations of assessment results using the cosine similarity technique. These pictorial representations offer a clear and concise overview of performance, aiding students in understanding their strengths and weaknesses. This enhanced feedback supports a more targeted approach to learning and improvement.

5. IMPLEMENTATION

1 KEYWORDS:

Keywords are question-specific things that are essential for answering that question. These keywords play a significant role in penalizing or promoting the score evaluated by the similarity measurement module and must only contain the essential words in lower case.

2 SOLUTION:

The solution is a subjective answer that is being used to map students' responses. This solution must contain all the keywords and contexts discussed in the answers in separate lines/paragraphs. The teacher/evaluator typically prepares the solution to the question.

3 ANSWER

The answer is a subjective response from the student that is to be evaluated. It usually contains some or all of the keywords and spans 1 to a few sentences depending on the type of question and the



student's writing style. It almost always contains synonym words compared to the solution and, therefore, requires much more semantic care when processing.

4 DATA COLLECTION

To train and test the proposed model, there is a need for a massive amount of corpus containing subjective question answers, but there is no publicly available labeled subjective question answers corpus to the best of our knowledge. In this work, we create subjective answers labeled corpus. For generating corpus, the important thing is to target those websites and blogs where subjective questions and answers exist. We crawl various websites and collect a subjective question answers corpus, and the crawl data belong to various domains such as computer science and general knowledge.

5 DATA ANNOTATION

After getting crawled data, there is a further need to annotate data because that crawled data is unlabeled. To annotate data, a group of different volunteers is selected, which belong to the domain of our subjective question answers corpus. We hire 30 different annotators from different colleges and universities and reside in Pakistan's different cities. Most of them are students and teachers. The average age of annotators is in the 21-25 range, whereas some annotators are in the age range of 27-51. We task annotators to best score the subjective question answers according to the answers given by students.

6 PREPROCESSING MODULE

After taking inputs from the user, both the solution and the answer go through some preprocessing steps, which involve tokenization, stemming, lemmatization, stop words removal, case folding, finding, and attaching synonyms to the text. Note that stop words are not removed when passing the data to word2vec because word2vec contains a vast vocabulary and can utilize those stop words to make better semantic sense of the text. However, stop words are removed before passing to a machine learning model such as Multinomial Naive Bayes because they hinder the machine's ability to learn the patterns.

7 SIMILARITY MEASUREMENT MODULE

This module consists of WDM and Cosine Similarity functions which take two sentences or word vectors and return their Similarity. WDM tells us the dissimilarity while Cosine Similarity measures Similarity. Our approach uses both of these similarity measures one at a time and compares the results at the end. Various similarity (or dissimilarity) thresholds . 1) THRESHOLDS ANALYSIS Various thresholds used in this paper have been experimentally deduced to produce the optimal result. WDM thresholds of WDM_LOWER and WDM_UPPER represent the dissimilarity between two sentences, where more dissimilarity represents high similarity. 0.7 threshold for WDM_LOWER was experimentally observed to represent semantically very similar sentences, and 1.6 thresholds for WDM_UPPER were observed to represent semantically less similar sentences. Anything beyond 1.6 is assumed to be too dissimilar to consider viable for comparison. Similarly, Cosine similarity



thresholds COS_LOWER and COS_UPPER represent the similarity between two sentences. It should be noted that cosine similarity does not take the context of two sentences into account when measuring similarity as opposed to WDM, hence the usage of both of these similarity (or dissimilarity) measuring approaches.

8 RESULT PREDICTING MODULE

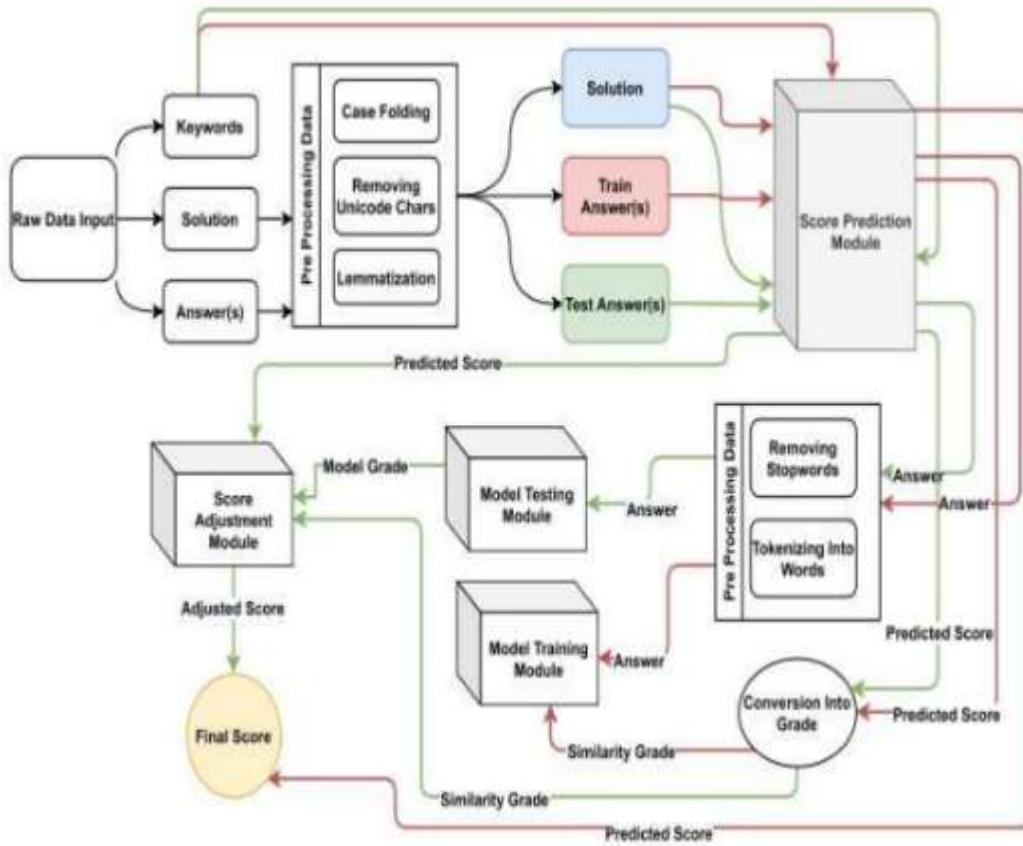
Result Predicting Module is the core of this work. It shows the working of this module. It operates on the following Algorithm 1: We now have the overall score calculated by our module using either WDM or Cosine Similarity while considering the maximum matched solution/answer sentence pairs. This result can be compared to an actual score or fed into a machine learning model to be trained.

9 MACHINE LEARNING MODEL MODULE

This model consists of Machine learning models trained on the data obtained from the result prediction module. Its working is as follows:

- Input data from Result Prediction Module.
- Preprocess the solution and answer, removing stop words, and use Countvectorizer to represent them in either Bag of Words or TF-IDF form.
- Convert the overall score obtained from Result Prediction Module into some category. Four categories A, B, C, and D, are used in the paper, representing 1st, 2nd, 3rd, and 4th quarter of a 100. For example, A represents marks from 0 to 25, and B represents 26 to 50.
- The number of categories is kept to a minimum because of the unavailability of the actual dataset. Practically, these categories can be extended to cover smaller score ranges.
- A machine learning model such as Multinomial Naive Bayes, which performs well for multi-class classification, is chosen.
- The preprocessed answer is used as testing data with the machine learning model to predict its class/category, and that category is checked with the result obtained from Result Prediction Module. This gives us confidence in the predicted result from the model.

6. SYSTEM ARCHITECTURE



7.EXPECTED RESULTS

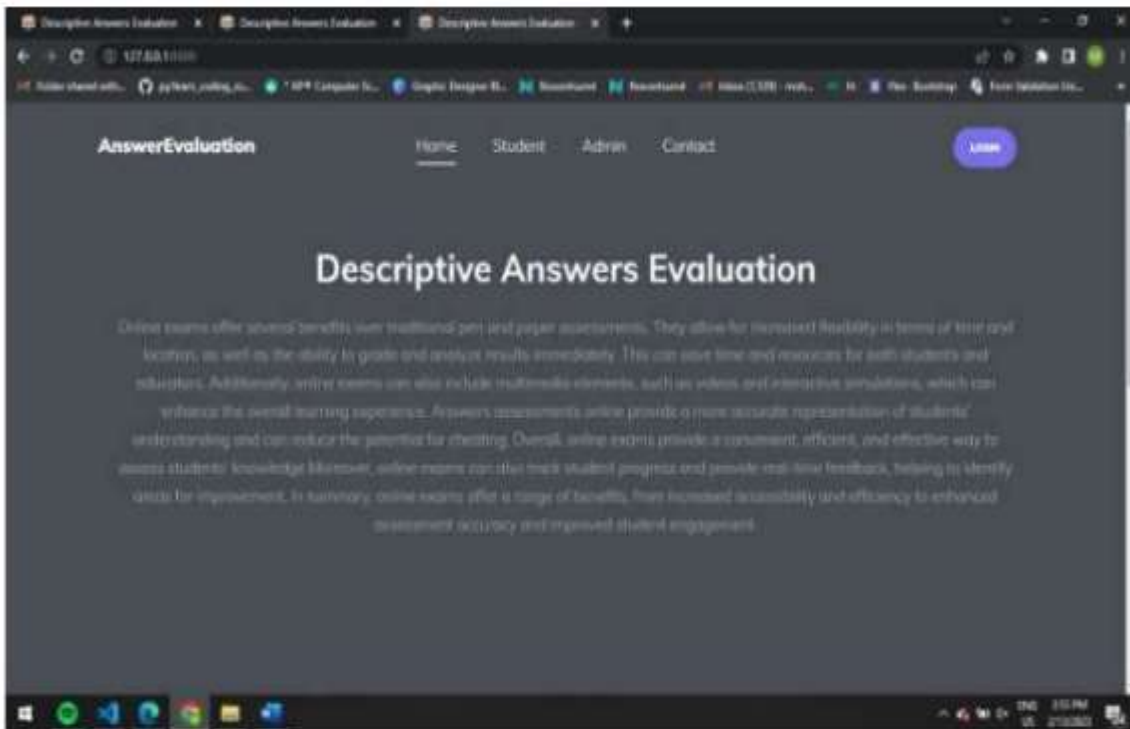


Fig.1

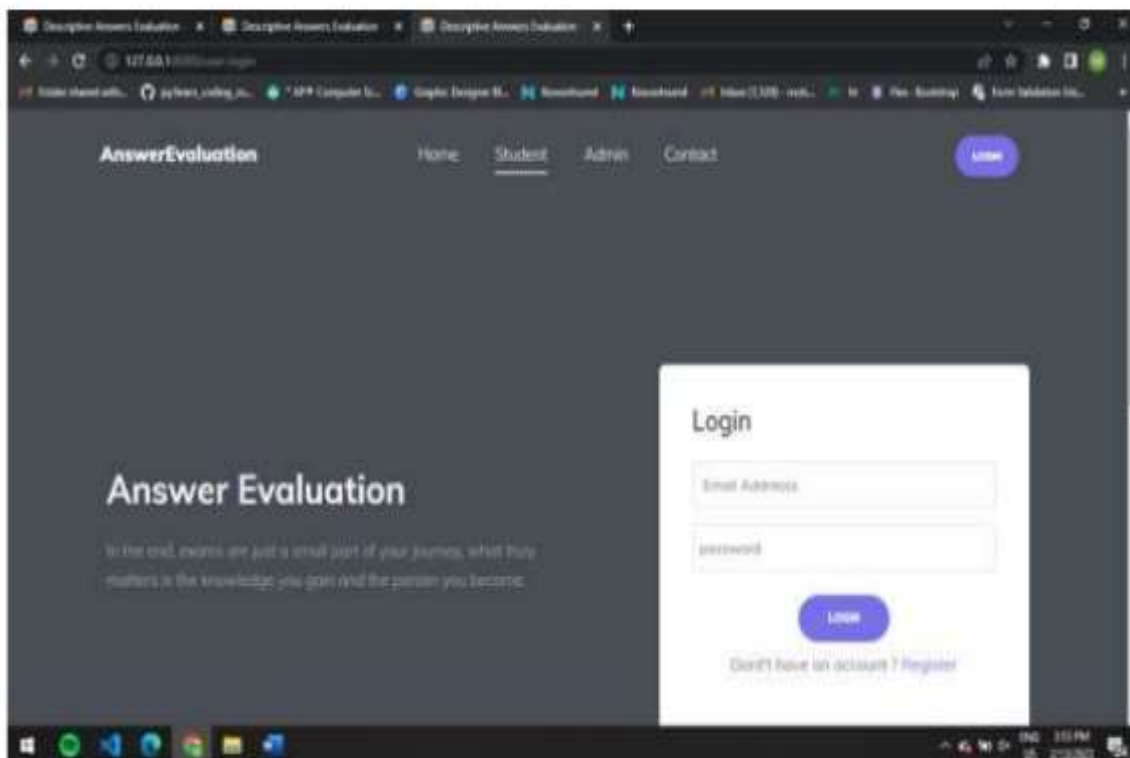


Fig.2

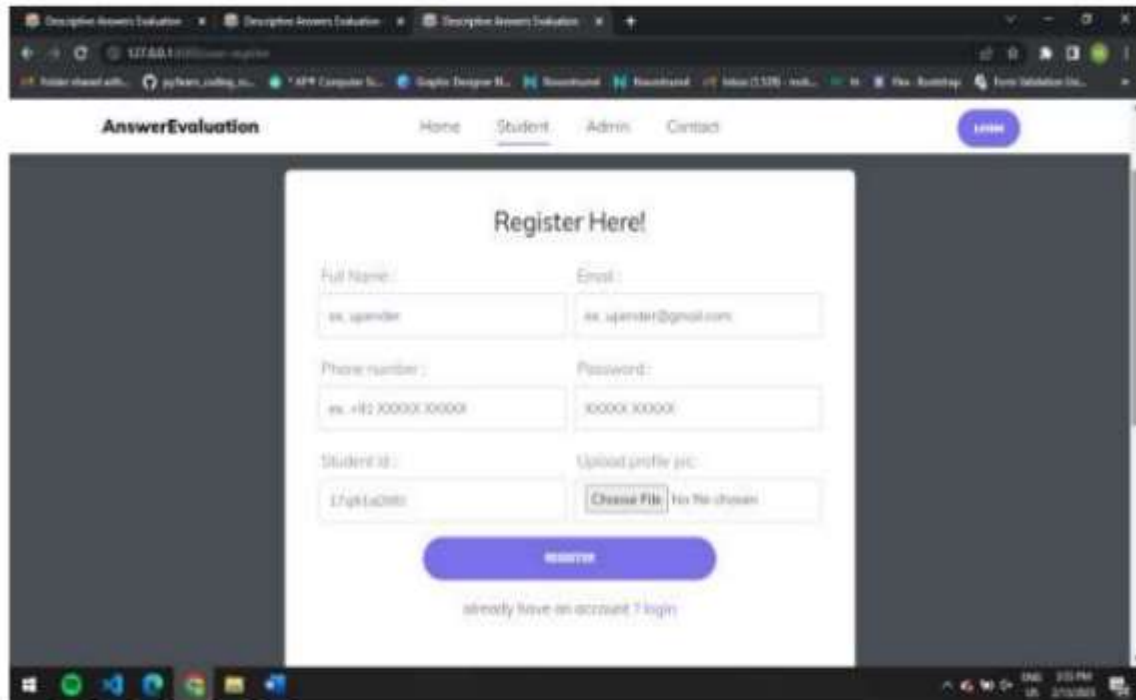


Fig.3

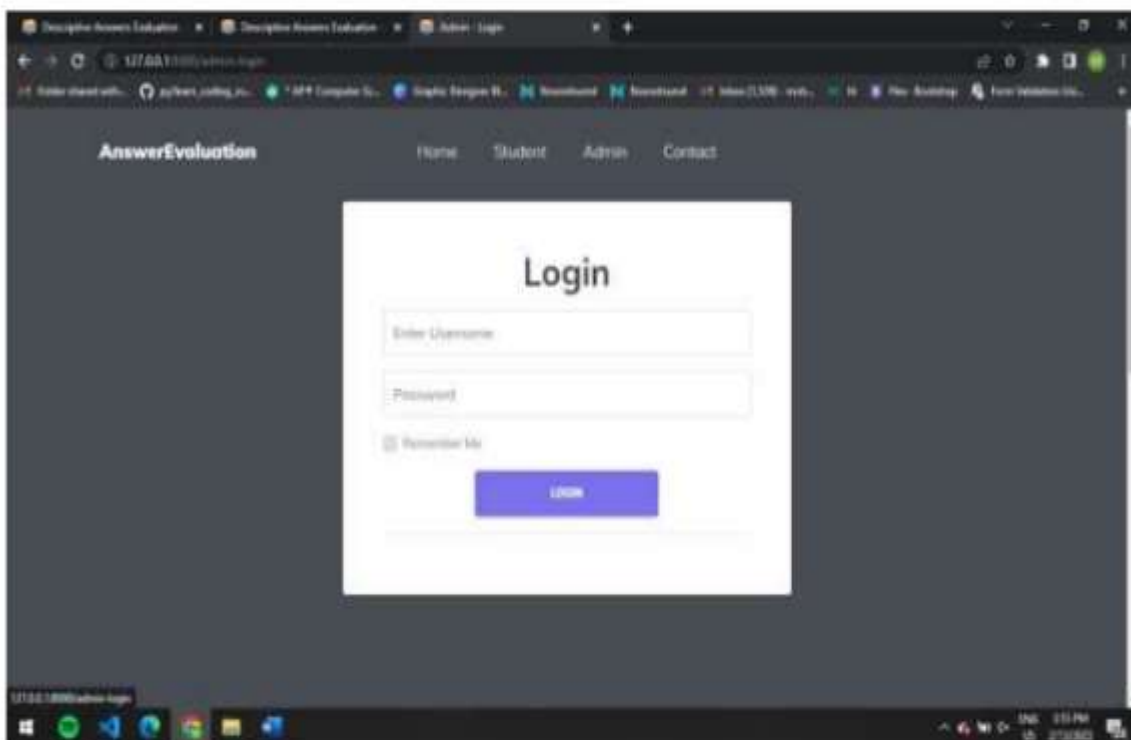


Fig.4

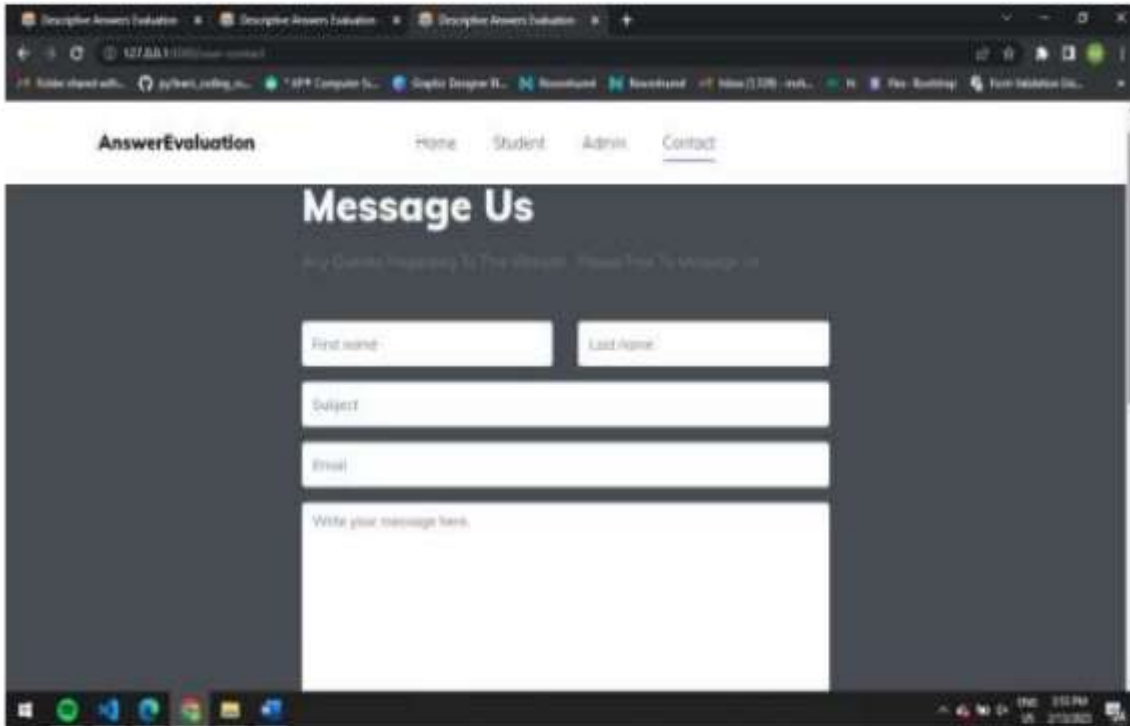


Fig.5

User:

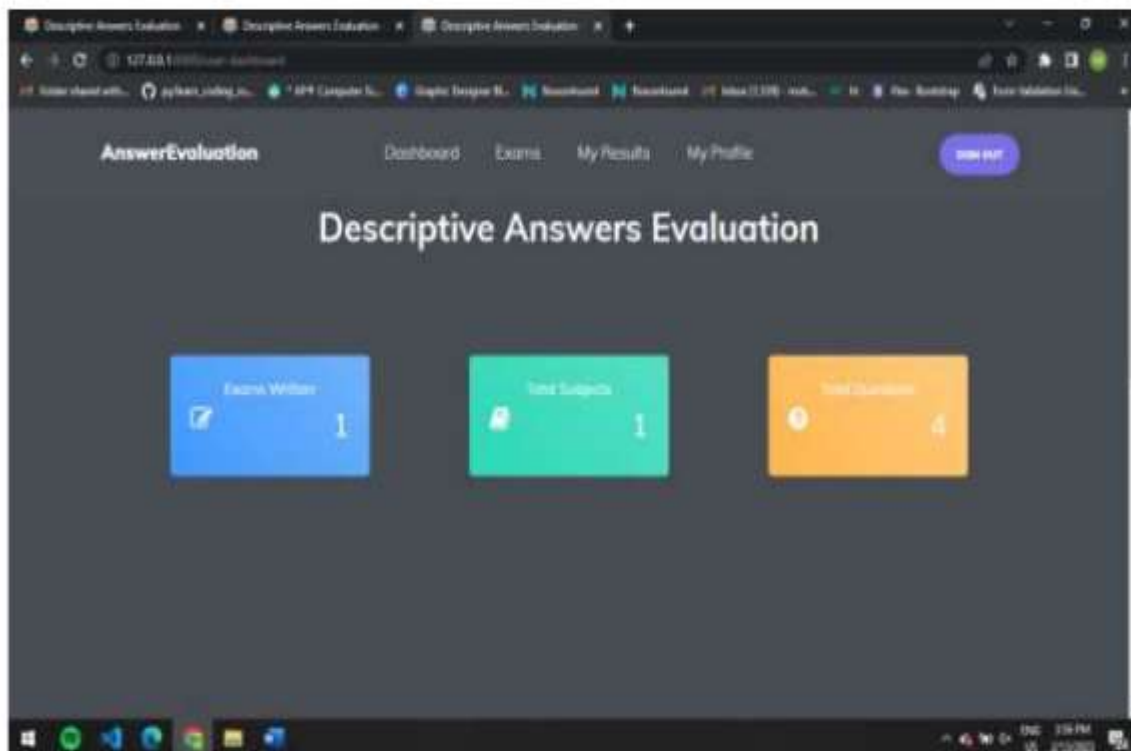


Fig.6

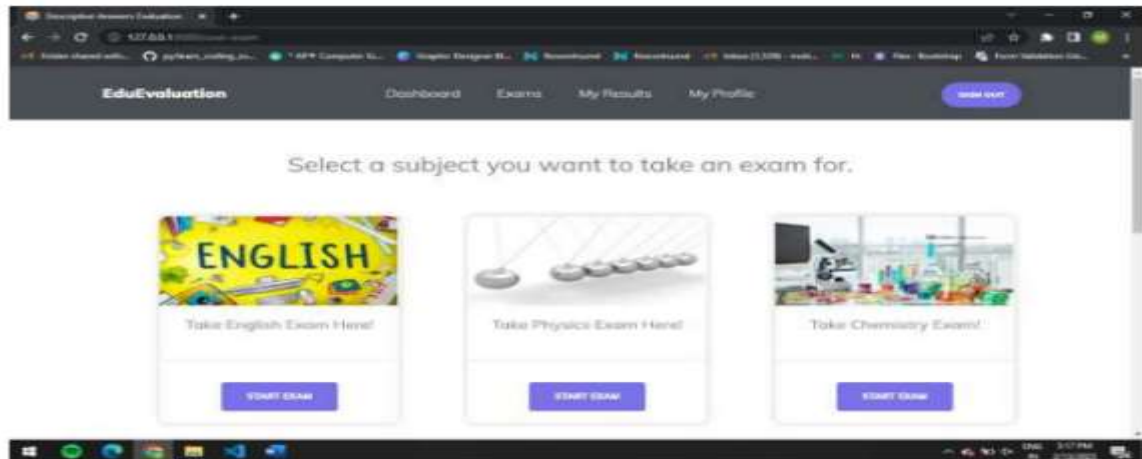


Fig.7

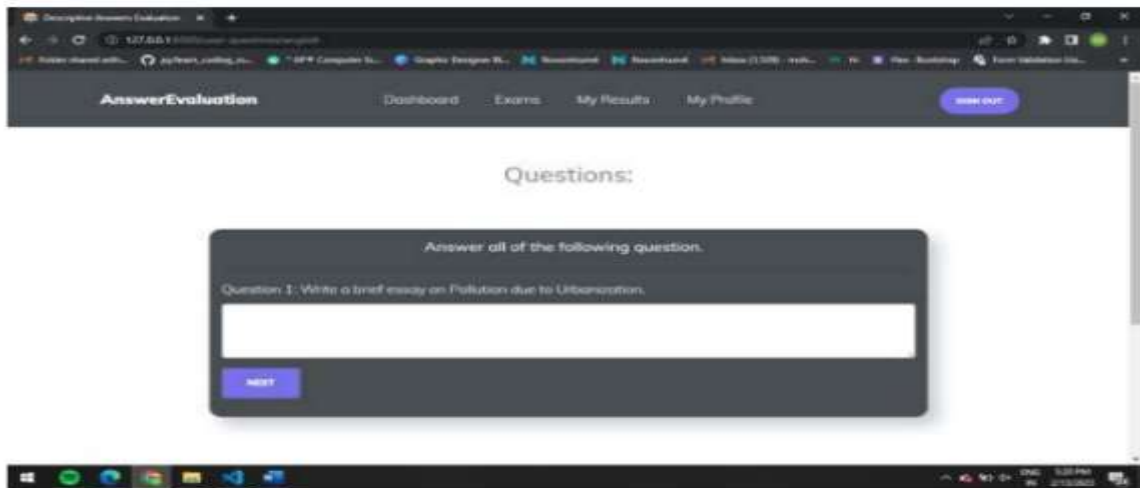


Fig.8



Fig.9

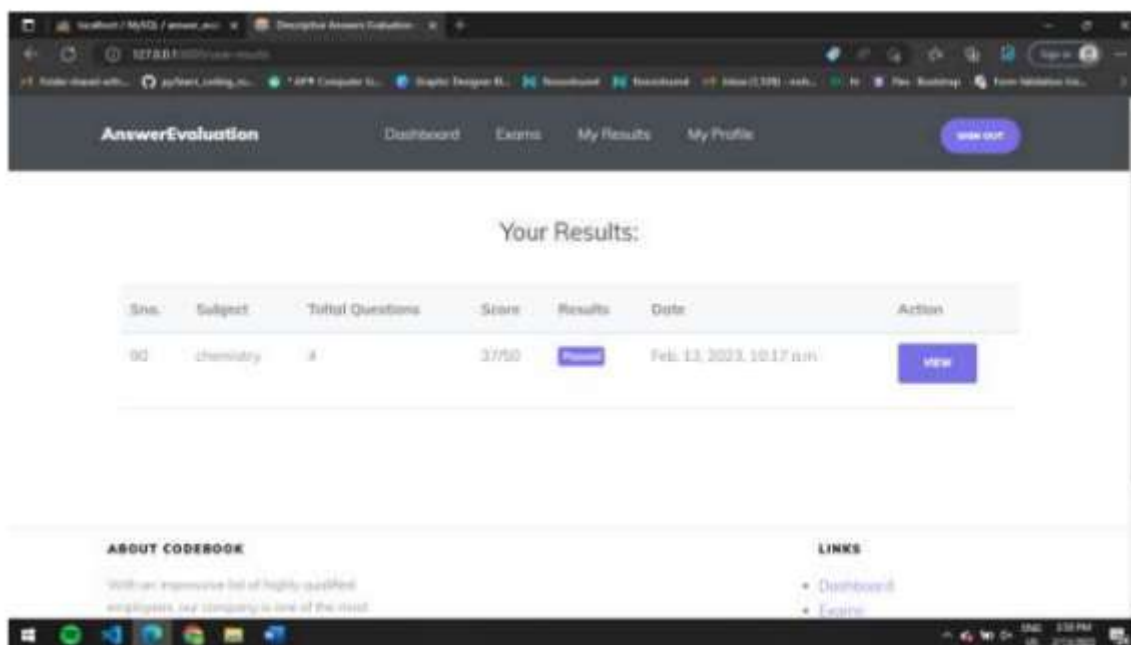


Fig.10

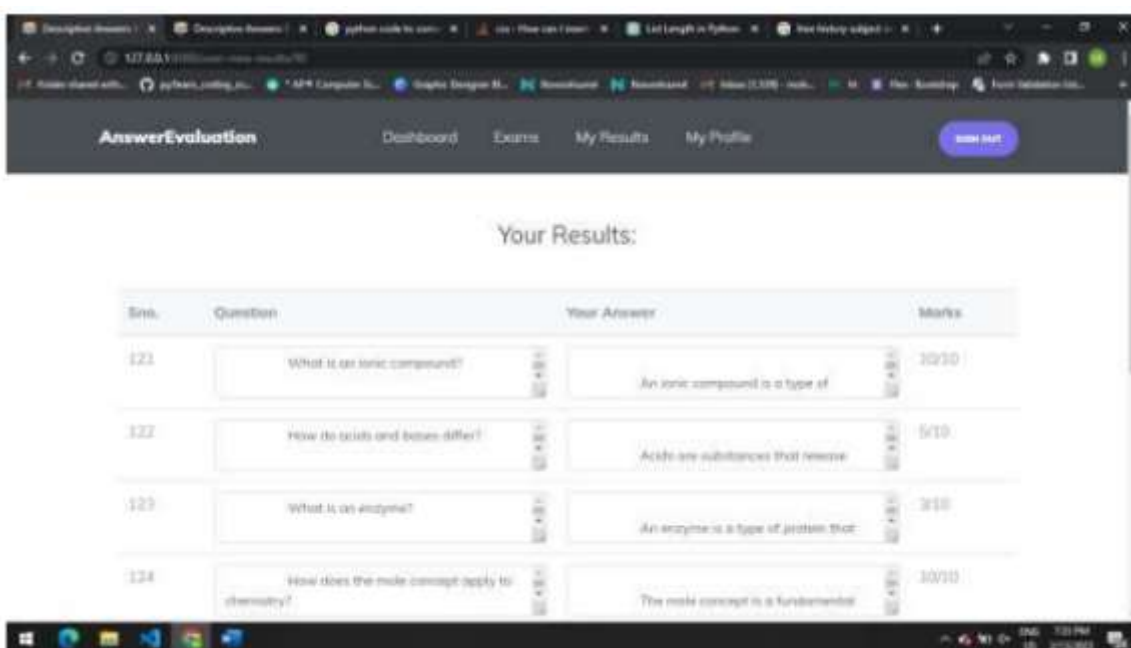


Fig.11

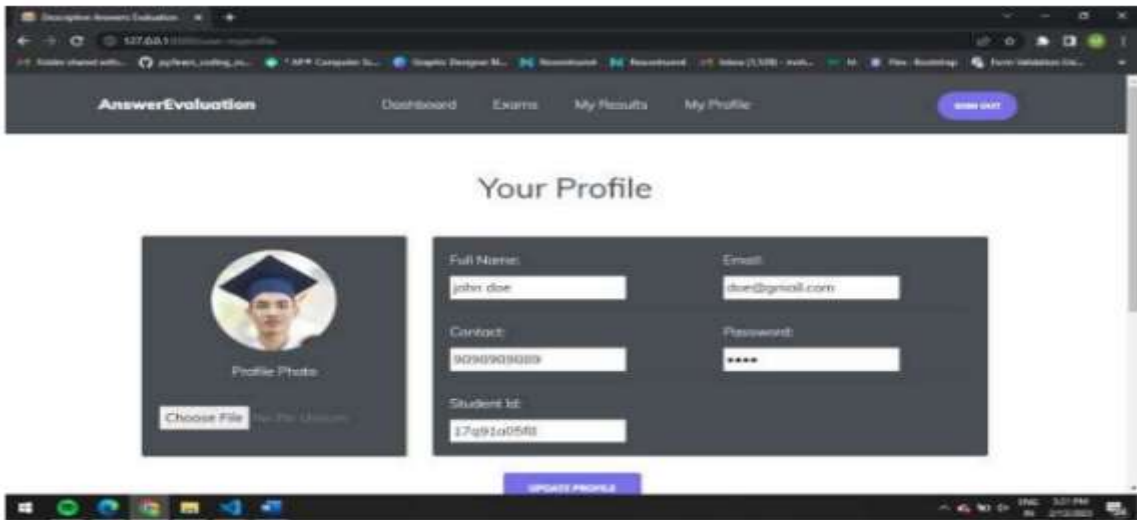


Fig.12

Admin:

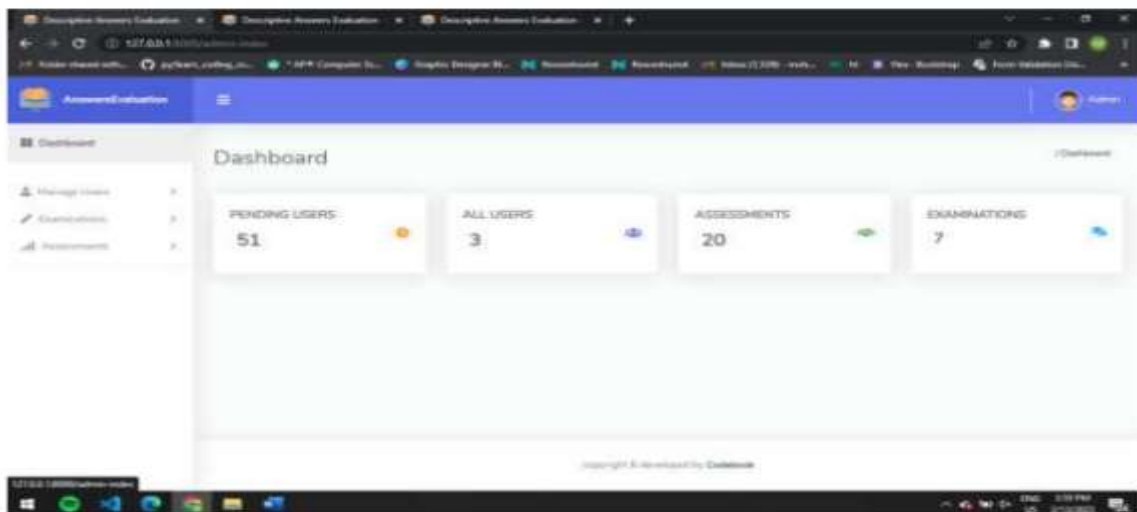


Fig.13

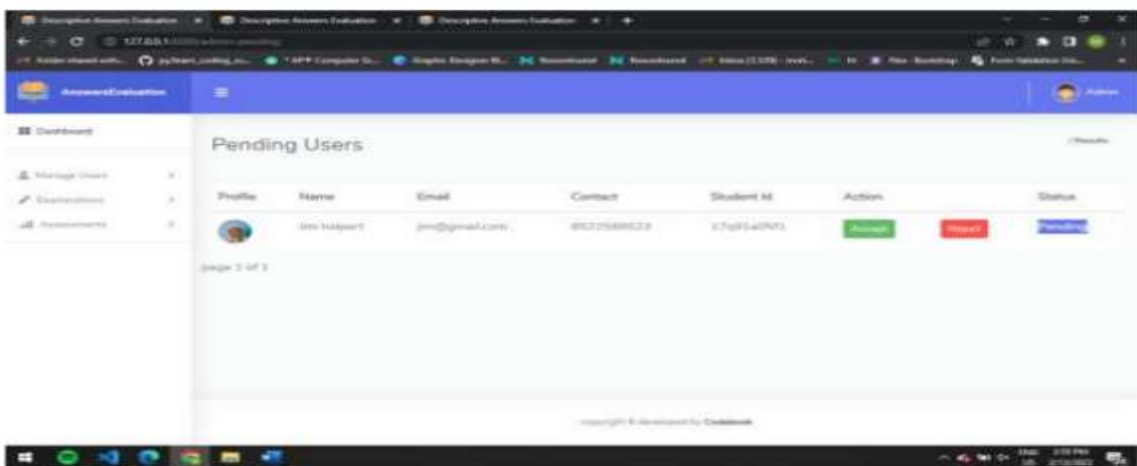


Fig.14

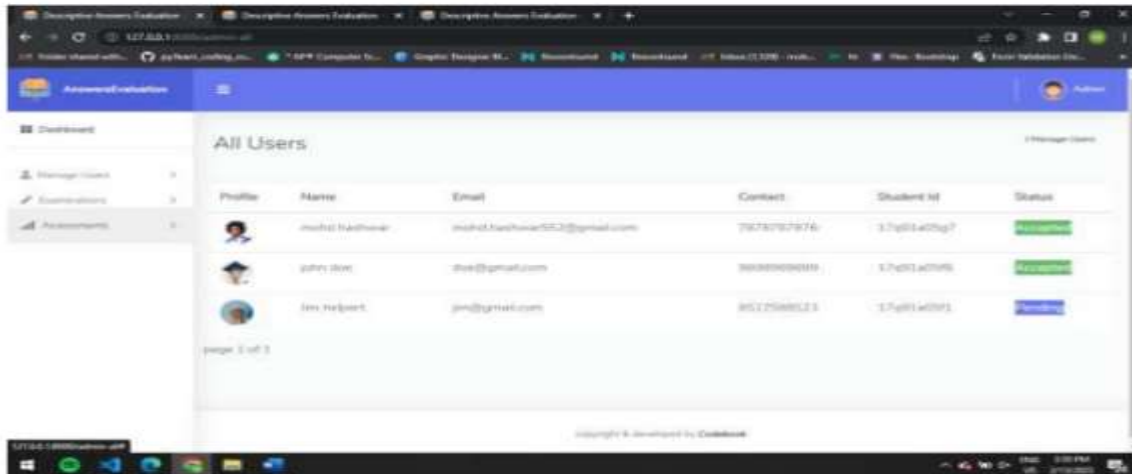


Fig.15

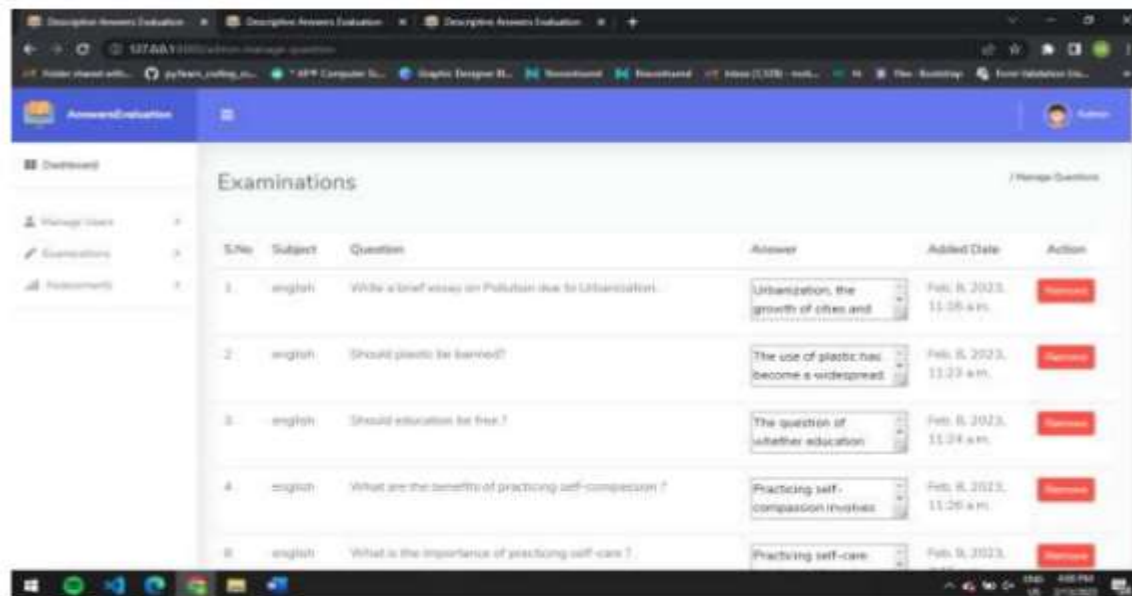
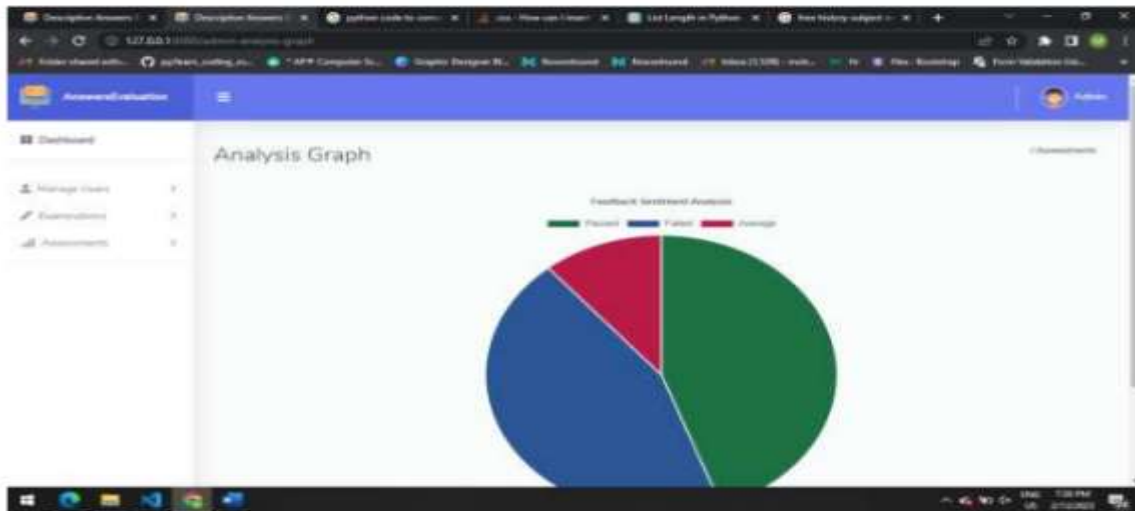


Fig.18

**Fig.20**

8. CONCLUSIONS

This paper proposed a novel approach to subjective answers evaluation based on machine learning and natural language processing techniques. Two score prediction algorithms are proposed, which produce up to 88% accurate scores. Various similarity and dissimilarity thresholds are studied, and various other measures such as the keyword's presence and percentage mapping of sentences are utilized to overcome the abnormal cases of semantically loose answers. The experimentation results show that on average word2vec approach performs better than traditional word embedding techniques as it keeps the semantics intact. Furthermore, Word Mover's Distance performs better than Cosine Similarity in most cases and helps train the machine learning model faster. With enough training, the model can stand on its own and predict scores without the need for any semantics checking. In terms of future improvements, the word2vec model can be trained especially for subjective answers evaluation of a particular domain, and with large data sets, the number of classes or grades in the model can be significantly increased. Subjective answers evaluation remains an interesting problem to tackle, and in the future, we hope to find more efficient ways to solve this problem.

6. REFERENCES

- [1] J. Wang and Y. Dong, "Measurement of text similarity: A survey," *Information*, vol. 11, no. 9, p. 421, Aug. 2020.
- [2] M. Han, X. Zhang, X. Yuan, J. Jiang, W. Yun, and C. Gao, "A survey on the techniques, applications, and performance of short text semantic similarity," *Concurrency Comput., Pract. Exper.*, vol. 33, no. 5, Mar. 2021.
- [3] M. S. M. Patil and M. S. Patil, "Evaluating Student descriptive answers using natural



language processing,” *Int. J. Eng. Res. Technol.*, vol. 3, no. 3, pp. 1716–1718, 2014.

[4] P. Patil, S. Patil, V. Miniyar, and A. Bandal, “Subjective answer evaluation using machine learning,” *Int. J. Pure Appl. Math.*, vol. 118, no. 24, pp. 1–13, 2018.

[5] J. Muangprathub, S. Kajornkasirat, and A. Wanichsombat, “Document plagiarism detection using a new concept similarity in formal concept analysis,” *J. Appl. Math.*, vol. 2021, pp. 1–10, Mar. 2021.

[6] V. Suma and Shavige Malleshwara Hills, "Data Mining based Prediction of Demand in Indian Market for Refurbished Electronics," *Journal of Soft Computing Paradigm (JSCP)* 2, pp. 101-110, 2020.

[7] M. Kusner, Y. Sun, N. Kolkin, and K. Weinberger, “From word embeddings to document distances,” in *Proc. Int. Conf. Mach. Learn.*, 2015, pp. 957–966.

[8] C. Xia, T. He, W. Li, Z. Qin, and Z. Zou, “Similarity analysis of law documents based on Word2vec,” in *Proc. IEEE 19th Int. Conf. Softw. Qual., Rel. Secur. Companion (QRS-C)*, Jul. 2019, pp. 354–357.

[9] H. Mittal and M. S. Devi, “Subjective evaluation: A comparison of several statistical techniques,” *Appl. Artif. Intell.*, vol. 32, no. 1, pp. 85–95, Jan. 2018.

[10] L. A. Cutrone and M. Chang, “Automarking: Automatic assessment of open questions,” in *Proc. 10th IEEE Int. Conf. Adv. Learn. Technol.*, Sousse, Tunisia, Jul. 2010, pp. 143–147.

[11] G. Srivastava, P. K. R. Maddikunta, and T. R. Gadekallu, “A two-stage text feature selection algorithm for improving text classification,” *Tech. Rep.*, 2021.

[12] H. Mangassarian and H. Artail, “A general framework for subjective information extraction from unstructured English text,” *Data Knowl. Eng.*, vol. 62, no. 2, pp. 352–367, Aug. 2007.

[13] B. Oral, E. Emekligil, S. Arslan, and G. Eryigit, “Information extraction from text intensive and visually rich banking documents,” *Inf. Process. Manage.*, vol. 57, no. 6, Nov. 2020, Art. no. 102361.

[14] H. Khan, M. U. Asghar, M. Z. Asghar, G. Srivastava, P. K. R. Maddikunta, and T. R. Gadekallu, “Fake review classification using supervised machine learning,” in *Proc. Pattern*



Recognit. Int. Workshops Challenges (ICPR). New York, NY, USA: Springer, 2021, pp. 269–288.

[15] W. K. Michael and A. G. Thomas, "A Framework for the Evaluation of Statistical Prediction Models," CHEST, vol. 158, no. 1, pp. S29-S38, 2020.

[16] N. Madnani and A. Cahill, "Automated scoring: Beyond natural language processing," in Proc. 27th Int. Conf. Comput.Linguistics (COLING), E. M. Bender, L. Derczynski, and P. Isabelle, Eds. Santa Fe, NM, USA: Association for Computational Linguistics, Aug. 2018, pp. 1099–1109.

[17] S. Deepa, A. Alli, Sheetac and S. Gokila, "Machine learning regression model for material synthesis prices prediction in agriculture," in materialstoday, 2021.

[18] G. Grefenstette, "Tokenization," in Syntactic Wordclass Tagging. Springer, 1999, pp. 117–133.

[19] K. Sirts and K. Peekman, "Evaluating sentence segmentation and word Tokenization systems on Estonian web texts," in Proc. 9th Int. Conf. Baltic (HLT) (Frontiers in Artificial Intelligence and Applications) vol. 328, U. Andrius, V. Jurgita, K. Jolantai, and K. Danguole, Eds. Kaunas, Lithuania: IOS Press, Sep. 2020, pp. 174–181.

[20] S. Matthew and Lewis, "Identifying airline price discrimination and the effect of competition," International Journal of Industrial Organization, vol. 78, 2021.