

# Early identification of Family Medicine residents at risk of failure using Natural Language Processing and Explainable Artificial Intelligence

## Authors:

<sup>1\*</sup>Abhisht Joshi, <sup>1,2</sup>Pouria Mortezaagha, <sup>1</sup>Diana Inkpen, <sup>1</sup>Edward Seale, <sup>1,3</sup>Douglas Archibald, <sup>1</sup>Kendall Noel, <sup>1,2</sup>Arya Rahgozar

<sup>1</sup>*University of Ottawa, Ontario, Canada.*

<sup>2</sup>*Ottawa Hospital Research Institute (OHRI), Ontario, Canada.*

<sup>3\*</sup> *Bruyère Research Institute, Ottawa, Ontario, Canada.*

Email: [abhishtjoshi16@gmail.com](mailto:abhishtjoshi16@gmail.com) , [ajoshi@uottawa.ca](mailto:ajoshi@uottawa.ca)

## Abstract

**Background:** During residency, each resident is observed and receives feedback based on their performance. Residency training is demanding, with a few residents struggling in their academic performance. A competency-based residency training program's success depends on its ability to identify residents with difficulty during their first year of post-graduate education and to provide them with timely intervention and support.

**Objective:** In large training programs such as Family Medicine, identifying residents at risk of failing their certification exams is difficult. We develop a AI system using state-of-the-art technologies in Machine Learning (ML), Deep Learning (DL), Natural Language Processing (NLP) and Explainable AI (XAI) to detect at-risk residents automatically.

**Methods:** We implemented ML, DL and NLP models for the prediction and its performance analysis. The target variable chosen for the prediction was the determination of whether the resident would fail or pass their certification exam. XAI was used to enhance the understanding of the model's inner workings.

**Results:** In total, there were 1382 data points of residents. The champion model, Support Vector Machine (SVM), achieved an accuracy of 89.05% and an F1 score of 74.54 for the multiclass classification when multimodal (text and tabular) data was used. This model outperformed the models that only used qualitative or quantitative data exclusively.

**Conclusion:** Combining qualitative and quantitative data represents a novel approach and has provided better classification results. This research demonstrates the feasibility of an automated AI system for the early identification of residents at risk of academic struggle.

## Introduction

Residency programs play a crucial role in molding future specialist and generalist physicians by providing them extensive clinical experience. When residents encounter problems late in their residency, it is often more difficult to rectify them. Ideally, these problems should be identified early<sup>a</sup> in residency, which would allow them to be resolved promptly, avoiding ramifications, such as failure on certification examinations. Early identification of at-risk residents is a vital responsibility of a residency's program of evaluation. Artificial intelligence (AI) is suggested to be used for early detection of at-risk residents.

According to a survey of internal medicine program directors, struggling residents display characteristics of inadequate medical knowledge, poor clinical judgment, and inefficient use of their time<sup>1</sup>. Additionally, struggling residents may feel overburdened, unsure of the objectives of training, unclear about their performance evaluation and thus incapable of prioritizing areas of improvement<sup>2</sup>. The demanding nature of residency programs, the long work hours, and the stress of transitioning to residency, all play a role in these challenges<sup>3</sup>.

The use of Machine Learning (ML) and Natural Language Processing (NLP) in medical education to improve students' and healthcare professionals' teaching and learning processes has gained popularity in recent years. Several promising strategies have emerged because of studies exploring the potential applications of ML<sup>8,10,12-19</sup> and NLP<sup>4-8,9,11</sup> in medical education. This emphasis on NLP has practical applications. Automated essay scoring (AES) for medical knowledge examinations and constructed-response assignments, automates grading by linking language to human scores using word processing and NLP<sup>4,5</sup>. NLP was applied to examine faculty and medical residents' feelings about entrustable professional activity (EPA) evaluations revealing that general surgery participants expressed fewer positive emotions compared to those in emergency medicine<sup>6</sup>. Explainable AI (XAI) has shown to improve decision making and educational results, by improving a prediction model's transparency and comprehension<sup>20-24</sup>. The explainability of the otherwise complex reasons behind ML predictions play an important role in establishing end users' trust in AI systems' credibility, by mapping the important patterns, phrases and terms associated with the highest information value for the predictions<sup>25</sup>.

The identification of at-risk residents and the subsequent educational interventions can be greatly improved with the use of AI. We propose a multimodal approach of combining Qualitative (Text) and Quantitative (Numerical) data utilizing Post Graduate Year 1 (PGY 1) Family Medicine rotation data, along with XAI for model explanation. This integration of data in the form of text and narration, allows AI systems to consider real-life experiences<sup>26</sup>. Integration of humanistic opinions can enhance AI's data-centric methods. Our study attempts to close the research gap, by investigating the efficiency of combining qualitative and quantitative data, via three independent sets of experiments. By examining the synergy between the two types of data, we hope to show that their combination produces more valid results. The study makes use of advanced models such as XLNET<sup>29</sup> for the prediction of the residents at risk, along with the implementation of XAI to help improve the end user's understanding.

The objectives of our study are twofold. First, we will explore the feasibility of the automating the identification and prediction of at risk residents in educational trouble using ML and NLP. Second, we will use XAI and language models (LM), to identify the important latent

---

<sup>a</sup> By Early we mean, as soon as the resident evaluation data becomes available.

characteristics, discriminating patterns and generate insights.

## Methods

### *Research Setting and Participants*

The dataset utilized for the study is the Family Medicine In-Training Evaluation Report (ITER) of first-year residents (PGY-1) and their certification exam data. The ITER data is sourced from the Department of Family Medicine at the University of Ottawa. The exams data is based on the CFPC<sup>b</sup> certification exam, which has two components Simulated office oral (SOO) and Short-answer management problems (SAMP) and they make up the Canadian Family Medicine Certification exam. The exams are offered in the Spring and the Fall of each year. The exam data used in this analysis covers the following periods: Spring: 2018-2022, Fall: 2018 - 2022. There was no data for Spring 2020, as the exam was cancelled due to the COVID-19 pandemic.

Given this exam data, we focused on analyzing the performance outcomes represented by the target variable "Pass/Fail," which is a multiclass classification with four different classes, Fail in SOO, Fail in SAMP, Fail in both and Pass (See Figure 1). The dataset is imbalanced<sup>c</sup>, meaning some classes have a much greater count than others. Exam results are recorded as z-scores for the SOO and the SAMP where a z-score of less than -2.0 marks a failure. Of the 1382 SOO scores, 61 fell below this level. Of 1382 SAMP scores, 21 fell below -2.0. Only 8 out of 1382 failed both parts. Most candidates passed: 1292 out of 1382 either passed both components or one if only one was available. This distribution shows a significant skew towards the "PASS" category, highlighting the imbalance in the dataset.

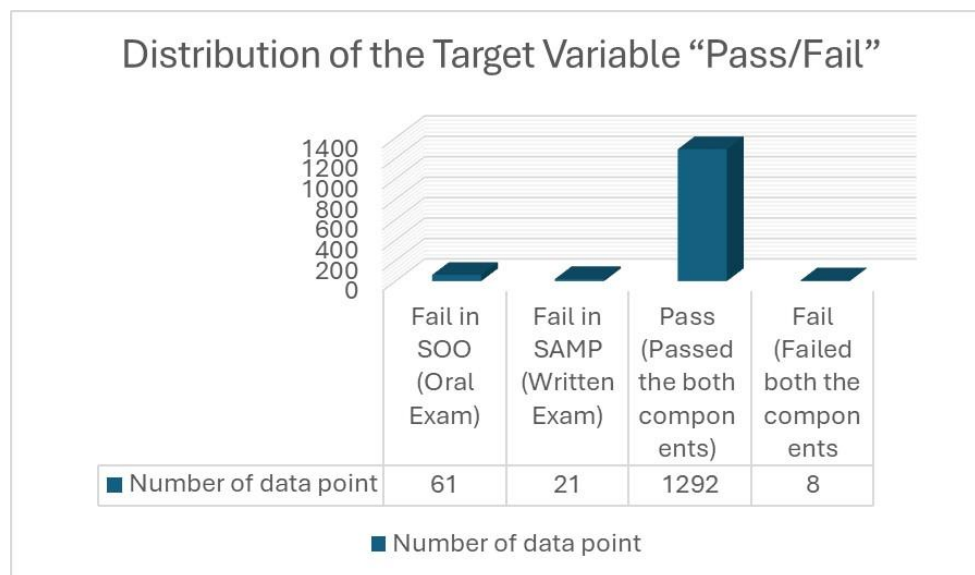


Figure 1. Distribution of Target Variable "Pass/Fail"

<sup>b</sup> CFPC: College of Family Physicians of Canada

<sup>c</sup> We used SMOTE to deal with the imbalance, please refer to methods section for more details.

## Interventions

This section presents the techniques employed to develop an AI-based medical residency intervention system. *Figure 2* illustrates the comprehensive approach implemented in our methodology.

The dataset comprises of multiple columns representing both qualitative and quantitative evaluations of the resident. This data was constructed from the assessment conducted by faculty members based on the resident's performance during their family medicine rotations in residency during their first year. We aim to utilize a multimodal approach, which involves combining quantitative (tabular) and qualitative (text) data, which includes feedback, and comments offered by the faculty to each resident. Furthermore, the examined dataset also includes the residents' final Pass or Fail status on their first certification exam attempt.

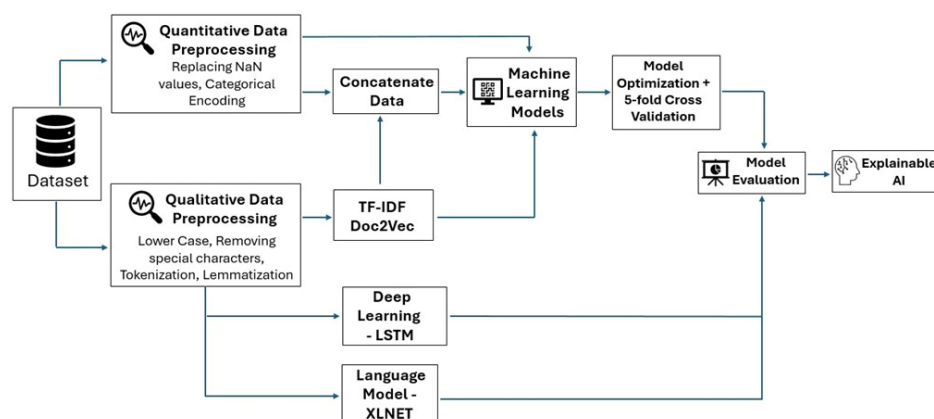


Figure 2. Research Methodology for Medical Residency Intervention

## AI Modelling

Data were cleaned in preparation for the experiment. The cleaning process involved filling empty entries and performing categorical encoding to the quantitative part of the data. Subsequently, a multi-step cleaning process was implemented for the qualitative data, encompassing converting the text to lowercase, removing special characters, tokenizing, and lemmatizing the text. We then vectorized the data using TF-IDF<sup>27</sup> to convert the text into numerical representations and produced dense vector representations encapsulating the semantic meaning of the text using Doc2Vec<sup>28</sup> embeddings for additional model processing. A Language Model XLNET<sup>29</sup> and a Deep Learning model Long Short-Term Memory (LSTM)<sup>30</sup> which is a type of Recurrent Neural Network designed to remember long-term dependencies were applied on qualitative data.

## Outcomes Measured

Our study involved a series of three AI experiments utilizing “Pass/Fail” as the target variable for all the three experiments. For comparison purposes, we experimented with three kinds of data within the dataset: qualitative, quantitative and multimodal.

### *Series of Experiments I: Qualitative Data*

In these experiments, we employed four approaches focusing solely on qualitative data via a) ML models on text vectorization (TF-IDF), b) ML Models on text embeddings (Doc2Vec), c) fine tuning the pre-trained XLNET on the data and d) LSTM, a Deep Learning model.

### *Series of Experiments II: Quantitative Data*

We utilized only quantitative data. SMOTE<sup>36</sup> was applied to balance the data. We used SMOTE to synthetically increase the short segments within the dataset without introducing any bias. The experimental methodology followed these steps: feature scaling and normalization, Principal Component Analysis (PCA), feature selection (Mutual Information gain), hyperparameter tuning (Grid search) and model evaluation (Cross-validation) and prediction.

### *Series of Experiments III: Multimodal (Combination of Qualitative and Quantitative Data)*

This final series of experiments aimed to enhance information richness and evaluate whether multimodality benefits the AI model's performance. Before concatenating the data, we applied TF-IDF and Doc2Vec to text data, followed by SMOTE, as done previously. Finally, we plugged the transformed data into the algorithm.

## ***Analysis of Outcomes***

### *Evaluation Metrics and Cross Validation*

The F1 (macro) score and accuracy were selected as evaluation metrics to measure model performance (*See Supplementary Material 1*). A 5-fold cross-validation was performed on the stratified data of which four parts were used for training the AI model and one part for testing. Through the cross-validation we reported the average of the indices with a variance of up to 3%<sup>33</sup>.

## ***Explainable AI (XAI)***

Explaining an AI model means rendering its output understandable to a human being<sup>35</sup>. XAI plays a crucial role in understanding what is going on behind a complex ML model. To gain insights into our AI model, SHAP<sup>31</sup> and BERTopic<sup>32</sup> were employed for model explanation and to understand why certain predictions were made. SHAP helped us identify important input variables, akin to the industrial term "fine classing," which is comparable to the PCA (Principal Component Analysis) used in our process. On the other hand, BERTopic is a topic modelling technique that leverages language model and assists in "explaining" output predictions by identifying and grouping topics within the textual data. Our utilization of BERTopic as a tool for "Global Explainability" is reinforced using a Language Model (XLNET). The key features were then shared with an expert.

### ***Shapley Additive Explanations (SHAP)***

SHAP<sup>31</sup> was used for model interpretation to identify crucial features and provide their order of importance. Together this information would allow the faculty and supervisors to understand which aspects of the residency to focus on when a resident is at risk.

### ***BERTopic***

BERTopic<sup>32</sup>, is a cutting-edge topic modelling technique that leverages BERT embeddings and class-based TF-IDF to create dense clusters, enabling the generation of easily interpretable topics while retaining important words in the topic descriptions. This approach is particularly valuable in the context of XAI, as it emphasizes transparency and interpretability in the modelling process.

### ***Feasibility and Acceptability***

These models were feasible in time and resources, making them an effective screening tool. Sharing key features with experts ensured feasibility, relevance, and acceptability.

### ***REB statement***

The project's scope was reviewed by the Office of Research Ethics and Integrity at the University of Ottawa and was determined deemed exempt from further review.

### **Results**

We reported the score averages of performance of the ML models for three sets of experiments in Table 1.

#### ***Series of Experiments I: Qualitative Data***

XLNET outperformed the other models with an accuracy of 72.45% and an F1 score of 55.48. LSTM performed second best with an accuracy of 70.33% and an F1 score of 53.73.

#### ***Series of Experiments II: Quantitative Data***

The SVM model, followed by the CatBoost model, achieved the highest performance with an accuracy of 81.71% and 80.93%, and an F1 score of 63.43 and 63.01, respectively.

#### ***Series of Experiments III: Multimodal Data***

The SVM (TF-IDF) model outperformed all other models with an accuracy of 89.05% and F1 score of 74.54, followed by SVM (Doc2Vec) with an accuracy of 82.10 and an F1 score of 72.40.

The results clearly indicate that the SVM (TF-IDF) model from Experiment III was the champion. It used TF-IDF on text data and a multimodal approach, outperforming other models

in Experiments I, II, and III. Interestingly, despite Doc2Vec being context-dependent, the TF-IDF model outperformed it, likely due to TF-IDF's better performance with smaller datasets<sup>33</sup>.

**Table 1. Model Performances of each Experiment**

Model	Experiment No.	Data Type	Accuracy	Precision	Recall	F1 score
MLR <sup>d</sup> (Doc2Vec)	1	Qualitative	67.16	42.25	42.65	42.44
BNB <sup>e</sup> (Doc2Vec)	1	Qualitative	62.93	38.10	34.91	36.43
GNB <sup>f</sup> (Doc2Vec)	1	Qualitative	67.55	39.45	46.27	42.58
SVM(Doc2Vec)	1	Qualitative	70.33	58.61	43.75	50.10
MLR (TF-IDF)	1	Qualitative	66.50	41.20	47.25	44.02
BNB (TF-IDF)	1	Qualitative	61.68	35.01	32.84	33.89
GNB (TF-IDF)	1	Qualitative	68.81	38.91	45.65	42.01
SVM (TF-IDF)	1	Qualitative	70.93	58.35	45.44	51.09
LSTM	1	Qualitative	70.33	60.22	48.51	53.73
XLNET	1	Qualitative	72.45	63.10	49.51	55.48
SVM	2	Quantitative	81.71	64.98	61.97	63.43
RF <sup>g</sup>	2	Quantitative	76.43	62.60	59.10	60.79
MLP <sup>h</sup>	2	Quantitative	77.92	64.04	59.79	61.84

<sup>d</sup> MLP: Multinomial Logistic Regression

<sup>e</sup> BNB: Bernoulli Naïve Bayes

<sup>f</sup> GNB: Gaussian Naïve Bayes

<sup>g</sup> RF: Random Forest

<sup>h</sup> MLP: Multi-layer Perceptron

CatBoost	2	Quantitative	80.93	66.10	60.94	63.01
MLR (Doc2Vec)	3	Combined	75.42	69.48	67.85	68.65
BNB(Doc2Vec)	3	Combined	70.79	52.50	49.85	51.14
GNB (Doc2Vec)	3	Combined	77.60	71.80	69.85	70.81
SVM(Doc2Vec)	3	Combined	82.10	74.80	70.15	72.40
MLR (TF-IDF)	3	Combined	80.26	73.45	70.10	71.73
BNB (TF-IDF)	3	Combined	70.24	54.22	48.45	51.17
GNB (TF-IDF)	3	Combined	82.15	72.32	69.69	70.98
SVM (TF-IDF)	3	Combined	89.05	76.11	73.04	74.54

### ***Explainable AI***

#### ***SHAP***

The summary plot for the most important features is shown in Figure 3 and is explained below: The summary plot shows the feature in order of importance, where the bar length for each class represents its impact on that class. Target variables were labelled as follows, Class 0 for students who fail SAMP, Class 1 for those who fail SOO, Class 2 for those who fail both, and Class 3 for students who pass both SAMP and SOO. The most crucial feature was, "Were the rotation objectives discussed with the resident?" This feature has two possible answers: Yes or No. It was observed that the residents who had discussed their objectives with their supervisor had a failure rate of 7.2%, but those who did not have a failure rate of just 4.1%. It is hypothesized that during residency, there is a higher chance that the rotation's objectives will be discussed if the resident is having difficulty. This contrasts with a resident who is known or perceived to be doing well and who has not undertaken these measures with the faculty.



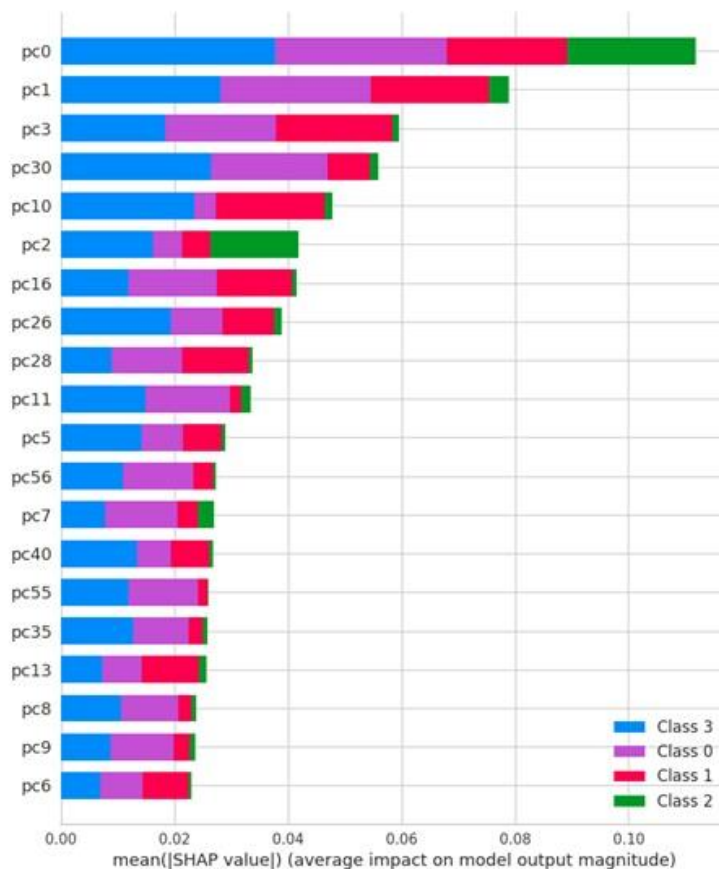


Figure 3. SHAP Summary plot

### ***BERTopic***

BERTopic has the potential to produce several topics, for the sake of clarity we have selected the top seven topics. Figure 4 presents the top terms for each topic using a bar chart that shows the importance of each word for the topic. The horizontal axis is the c-TF-IDF scores of the top five most representative words for each topic. For instance, Topic 1 “Trajectory” is related to the progress of the medical resident during their residency and their supervisor comments on whether the resident is off-trajectory or on-trajectory. It was observed that residents who were on-trajectory had a lower failure rate compared to the ones who were off-trajectory. Furthermore, Topic 7, for includes “teamwork”, and is strongly correlated with favorable outcomes. Residents who received comments like "team player with excellent communication skills," "excellent team player!" and "always a team player" passed their certification exams and were designated “PASS”. This shows that teamwork is a crucial indicator of positive performance outcomes. (See *Supplementary Material 2* for detailed analysis).



Figure 4. Topic generated through BERTopic

## Discussion

The SVM (TF-IDF) model from experiment 3, which utilized a multimodal approach combining qualitative and quantitative data, emerged as the model with the highest performance across experiments 1, 2, and 3, establishing itself as the champion model. This multimodal approach enhanced the prediction performance, comprehensive information and Cross-Domain learning of the model, leveraging the richness of information. Furthermore, the implementation of XAI including SHAP and BERTopic offered insightful analysis of the decision-making process of the model, hence stressing the explainability and transparency of the AI system.

The study's findings correspond with earlier studies using AI to forecast resident performance. An ML based model using previous data to predict resident performance reached 72% accuracy, showcasing the potential of ML in competency-based education<sup>18</sup>. The critical role of narrative feedback in decision-making is demonstrated by the ability of NLP and ML models to predict underperforming residents with 87% accuracy using narrative feedback from workplace-based assessments (WBAs)<sup>8</sup>. Additionally, classification and regression tree methods identified specific keywords in evaluation reports signaling residents at risk of failure, with precision rates of 23.3 and 23.4%<sup>17</sup>. The enhanced accuracy of our multi-modal technique implies that combining diverse data sources provides a more complete analysis than past studies depending on just one data type. Moreover, our use of XAI provides a better understanding of model decision-making, thereby differentiating our work.

Despite promising results, these techniques have limits. First, a relatively small and skewed sample of PGY 1 residents that might compromise generalizability of the model. Our study had a limited dataset of 1,382 data points. Although we attempted to mitigate potential bias by synthetically generating additional data using SMOTE, relying on authentic data would be more ideal for model accuracy. A second limitation could be restricted vocabulary coverage of the AI

models, highlighting the necessity for ongoing updates to the model's training data to include medical terminology and abbreviations.

Our future research aims to enhance the department's practical application by translating domain-specific knowledge into a usable tool through techniques like domain-specific language modeling and fine-tune models pre-trained on medical corpora such as Bio-BERT<sup>37</sup> and Bio-GPT<sup>38</sup>. Additionally, we plan to implement data augmentation techniques like GraphRAG<sup>39</sup> to further improve the model's performance.

## **Conclusion**

The Multi-modal SVM (TF-IDF) model showed the best performance among all the three experiments and was useful in predicting at-risk residents in the medical residency training program. This research emphasizes the advantages of amalgamating qualitative and quantitative data as well as using Explainable AI methods to give insightful analysis for early identification and intervention of residents in difficulty.

## References

1. Yao DC, Wright SM. The challenge of problem residents. *J Gen Intern Med*. 2001;16(7):486-492. doi:[10.1046/j.1525-1497.2001.016007486.x](https://doi.org/10.1046/j.1525-1497.2001.016007486.x)
2. Perez AR, Boscardin CK, Pardo M. Residents' Challenges in Transitioning to Residency and Recommended Strategies for Improvement. *J Educ Perioper Med*. 2022;24(1):E679. doi:10.46374/volxxiv\_issue1\_boscardin
3. ALQahtani DA, Mahzari MM, ALQahtani AA, Rotgans JI. Time Pressure Experienced by Internal Medicine Residents in an Educational Hospital in Saudi Arabia: A Qualitative Study. *Health Professions Education*. 2020;6(3):354-367. doi:10.1016/j.hpe.2020.05.005
4. Gierl MJ, Latifi S, Lai H, Boulais AP, De Champlain A. Automated essay scoring and the future of educational assessment in medical education. *Med Educ*. 2014;48(10):950-962. doi:10.1111/medu.12517
5. Burstein MDS Jill, ed. *Handbook of Automated Essay Evaluation: Current Applications and New Directions*. Routledge; 2013. doi:10.4324/9780203122761
6. Johnson D, Chopra S, Bilgic E. Exploring the Use of Natural Language Processing to Understand Emotions of Trainees and Faculty Regarding Entrustable Professional Activity Assessments. *J Grad Med Educ*. 2024;16(3):323-327. doi:[10.4300/JGME-D-23-00526.1](https://doi.org/10.4300/JGME-D-23-00526.1)
7. Wang M, Sun Z, Jia M, et al. Intelligent virtual case learning system based on real medical records and natural language processing. *BMC Medical Informatics and Decision Making*. 2022;22(1):60. doi:10.1186/s12911-022-01797-7
8. Yilmaz Y, Nunez AJ, Ariaeinejad A, Lee M, Sherbino J, Chan TM. Harnessing Natural Language Processing to Support Decisions Around Workplace-Based Assessment: Machine Learning Study of Competency-Based Medical Education. *JMIR Medical Education*. 2022;8(2):e30537. doi:[10.2196/30537](https://doi.org/10.2196/30537)
9. Park A, Hartzler AL, Huh J, McDonald DW, Pratt W. Automatically Detecting Failures in Natural Language Processing Tools for Online Community Text. *Journal of Medical Internet Research*. 2015;17(8):e4612. doi:10.2196/jmir.4612
10. Joshi A, Saggarr P, Jain R, Sharma M, Gupta D, Khanna A. CatBoost — An Ensemble Machine Learning Model for Prediction and Classification of Student Academic Performance. *Adv Data Sci Adapt Data Anal*. 2021;13(03n04):2141002. doi:10.1142/S2424922X21410023
11. Solano QP, Hayward L, Chopra Z, et al. Natural Language Processing and Assessment of Resident Feedback Quality. *Journal of Surgical Education*. 2021;78(6):e72-e77. doi:10.1016/j.jsurg.2021.05.012
12. Shailaja K, Seetharamulu B, Jabbar MA. Machine Learning in Healthcare: A Review. In: *2018 Second International Conference on Electronics, Communication and Aerospace Technology (ICECA)*. ; 2018:910-914. doi:10.1109/ICECA.2018.8474918
13. Pucchio A, Eisenhauer EA, Moraes FY. Medical students need artificial intelligence and machine learning training. *Nat Biotechnol*. 2021;39(3):388-389. doi:10.1038/s41587-021-00846-2
14. Saleem TJ, Chishti MA. Exploring the Applications of Machine Learning in Healthcare. *International Journal of Sensors Wireless Communications and Control*. 2020;10(4):458-472. doi:10.2174/2210327910666191220103417

15. Chen PHC, Liu Y, Peng L. How to develop machine learning models for healthcare. *Nat Mater*. 2019;18(5):410-414. doi:10.1038/s41563-019-0345-0
16. Javaid M, Haleem A, Pratap Singh R, Suman R, Rab S. Significance of machine learning in healthcare: Features, pillars and applications. *International Journal of Intelligent Networks*. 2022;3:58-73. doi:10.1016/j.ijin.2022.05.002
17. Tremblay G, Carmichael PH, Maziade J, Grégoire M. Detection of Residents With Progress Issues Using a Keyword-Specific Algorithm. *J Grad Med Educ*. 2019;11(6):656-662. doi:[10.4300/JGME-D-19-00386.1](https://doi.org/10.4300/JGME-D-19-00386.1)
18. Ariaeinejad A, Patel R, Chan TM, Samavi R. P031: Using machine learning algorithms for predicting future performance of emergency medicine residents. *Canadian Journal of Emergency Medicine*. 2017;19(S1):S88-S88. doi:10.1017/cem.2017.233
19. Kibble J, Plochocki J. Comparing Machine Learning Models and Human Raters When Ranking Medical Student Performance Evaluations. *J Grad Med Educ*. 2023;15(4):488-493. doi:10.4300/JGME-D-22-00678.1
20. Burkart N, Huber MF. A Survey on the Explainability of Supervised Machine Learning. *jair*. 2021;70:245-317. doi:10.1613/jair.1.12228
21. Khosravi H, Shum SB, Chen G, et al. Explainable Artificial Intelligence in education. *Computers and Education: Artificial Intelligence*. 2022;3:100074. doi:10.1016/j.caeai.2022.100074
22. Amann J, Blasimme A, Vayena E, Frey D, Madai VI. Explainability for artificial intelligence in healthcare: a multidisciplinary perspective. *BMC Med Inform Decis Mak*. 2020;20:310. doi:10.1186/s12911-020-01332-6
23. Albahri AS, Duhaim AM, Fadhel MA, et al. A systematic review of trustworthy and explainable artificial intelligence in healthcare: Assessment of quality, bias risk, and data fusion. *Information Fusion*. 2023;96:156-191. doi:10.1016/j.inffus.2023.03.008
24. Burkart N, Huber MF. A Survey on the Explainability of Supervised Machine Learning. *jair*. 2021;70:245-317. doi:10.1613/jair.1.12228
25. Jacovi A, Marasović A, Miller T, Goldberg Y. Formalizing Trust in Artificial Intelligence: Prerequisites, Causes and Goals of Human Trust in AI. Published online January 20, 2021. doi:10.48550/arXiv.2010.07487
26. Ostherr K. Artificial Intelligence and Medical Humanities. *J Med Humanit*. 2022;43(2):211-232. doi:10.1007/s10912-020-09636-4
27. Havrlant L, Kreinovich V. A simple probabilistic explanation of term frequency-inverse document frequency (tf-idf) heuristic (and variations motivated by this explanation). *International Journal of General Systems*. 2017;46(1):27-36. doi:10.1080/03081079.2017.1291635
28. Distributed Representations of Sentences and Documents | BibSonomy. Accessed February 13, 2024. <https://www.bibsonomy.org/bibtex/8dfe5d45e8d4b9c4812fd5590fbadef>
29. Yang Z, Dai Z, Yang Y, Carbonell J, Salakhutdinov R, Le QV. XLNet: Generalized Autoregressive Pretraining for Language Understanding. In: ; 2019. Accessed February 13, 2024. <https://www.semanticscholar.org/paper/XLNet%3A-Generalized-Autoregressive-Pretraining-for-Yang-Dai/e0c6abdbdecf04ffac65c440da77fb9d66bb474c>
30. Hochreiter S, Schmidhuber J. Long Short-Term Memory. *Neural Computation*. 1997;9(8):1735-1780. doi:10.1162/neco.1997.9.8.1735
31. Lundberg S, Lee SI. A Unified Approach to Interpreting Model Predictions. Published online November 24, 2017. doi:10.48550/arXiv.1705.07874

32. Grootendorst M. BERTopic: Neural topic modeling with a class-based TF-IDF procedure. Published online March 11, 2022. doi:10.48550/arXiv.2203.05794
33. Yadav S, Shukla S. Analysis of k-Fold Cross-Validation over Hold-Out Validation on Colossal Datasets for Quality Classification. In: 2016 IEEE 6th International Conference on Advanced Computing (IACC). ; 2016:78-83. doi:10.1109/IACC.2016.25
34. Cahyani DE, Patasik I. Performance comparison of TF-IDF and Word2Vec models for emotion text classification. Bulletin EEI. 2021;10(5):2780-2788. doi:10.11591/eei.v10i5.3157
35. Burkart N, Huber MF. A Survey on the Explainability of Supervised Machine Learning. jair. 2021;70:245-317. doi:10.1613/jair.1.12228
36. Chawla NV, Bowyer KW, Hall LO, Kegelmeyer WP. SMOTE: Synthetic Minority Over-sampling Technique. Journal of Artificial Intelligence Research. 2002;16:321-357. doi:10.1613/jair.953
37. Lee J, Yoon W, Kim S, et al. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. arXiv.org. doi:[10.1093/bioinformatics/btz682](https://doi.org/10.1093/bioinformatics/btz682)
38. Luo R, Sun L, Xia Y, et al. BioGPT: Generative Pre-trained Transformer for Biomedical Text Generation and Mining. arXiv.org. doi:[10.1093/bib/bbac409](https://doi.org/10.1093/bib/bbac409)
39. Edge D, Trinh H, Cheng N, et al. From Local to Global: A Graph RAG Approach to Query-Focused Summarization. arXiv.org. April 24, 2024. Accessed September 5, 2024. <https://arxiv.org/abs/2404.16130v1>