## Language Model Classifier Aligns Better with Physician Word Sensitivity than XGBoost on Readmission Prediction

Ming Cao<sup>\*†‡</sup> Grace Yang<sup>\*†‡</sup> Lavender Y. Jiang<sup>†‡</sup> Xujin C. Liu<sup>‡§</sup> Alexander T. M. Cheung<sup>‡¶</sup> Hannah Weiss<sup>‡</sup> David Kurland<sup>‡</sup> Kyunghyun Cho<sup>†||</sup> Eric K. Oermann<sup>ঠ</sup> MC7787@NYU.EDU GY654@NYU.EDU LYJ2002@NYU.EDU CHRIS.LIU@NYU.EDU ALEXANDER.CHEUNG@NYULANGONE.ORG HANNAH.WEISS@NYULANGONE.ORG DAVID.KURLAND@NYULANGONE.ORG KC119@NYU.EDU

ERIC.OERMANN@NYULANGONE.ORG

## Abstract

Traditional evaluation metrics for classification in natural language processing such as accuracy and area under the curve fail to differentiate between models with different predictive behaviors despite their similar performance metrics. We introduce sensitivity score, a metric that scrutinizes models' behaviors at the vocabulary level to provide insights into disparities in their decision-making logic. We assess the sensitivity score on a set of representative words in the test set using two classifiers trained for hospital readmission classification with similar performance statistics. Our experiments compare the decision-making logic of clinicians and classifiers based on rank correlations of sensitivity scores. The results indicate that the language model's sensitivity score aligns better with the

professionals than the xgboost classifier on tf-idf embeddings, which suggests that xgboost uses some spurious features. Overall, this metric offers a novel perspective on assessing models' robustness by quantifying their discrepancy with professional opinions. Our code is available on GitHub.<sup>1</sup>

**Keywords:** hospital readmission prediction, sensitivity analysis, model interpretability, MIMIC-III

## 1. Introduction

Predicting 30-day all-cause hospital readmission is a classical problem in medical informatics as 30-day readmissions are associated with longer hospital stays, higher mortality rates, and significant operating expenses. Several natural language processing (NLP) tasks have used models based on BioClinicalBERT, bag-of-words, and BI-LSTM; however, they all achieve similar test performances despite difference in modelling tech-

 $<sup>^{\</sup>ast}$  These authors contributed equally

 $<sup>^\</sup>dagger$  NYU Center for Data Science

 $<sup>^\</sup>ddagger$  NYU Langone Health

 $<sup>\</sup>$  NYU Tandon School of Engineering

 $<sup>\</sup>P$ NYU Grossman School of Medicine

<sup>|</sup> Prescient Design

https://github.com/nyuolab/Model\_ Sensitivity

niques (Jamei et al.; van Walraven et al.; Xiao et al., b; Alsentzer et al.; Huang et al.). It is challenging to select the best model for deployment in such cases with similar performance statistics, particularly in the case of complex machine-learning models whose decision-making logic may differ substantially from professionals. This lack of interpretability and insight into model predictions is a substantial barrier to medical deployment for the fear of potentially causing harm or incurring additional costs (Jackson et al.; Xiao et al., a).

To keep the model accountable, we propose a general-purpose evaluation framework to quantify the classifier's sensitivity to individual words based upon text perturbations. Text perturbations have been used for model bias detection (Prabhakaran et al.), but our goal is to recast it for model sensitivity quantification as a function of model vocabulary. Similar attempts to evaluate models have been proposed, but with their respective limitations. We discuss the existing methods and the motivation to create our metric in Appendix D. For a particular model, we compute the sensitivity score for each target word as the  $L_1$  distance between the classifier outputs before and after perturbations of the word in the test corpora. We test the language model classifier versus an extreme gradient boosted tree (xgboost) classifier (Chen and Guestrin) with tf-idf (Term Frequency-Inverse Document Frequency) embeddings as features (Sparck Jones).

The sensitivity score quantifies the significance of tokens as perceived by the model, which enables us to inspect the classifier's decision-making logic in the context of professional opinions. Ideally, if the classifier is sensitive to clinically significant words deemed by professionals, and indifferent to words considered trivial, then the classifier conforms to professional opinions and is less likely to have learned spurious correlations. For example, a reliable readmission classifier could be sensitive to words like "dementia" and indifferent to words like "parent". Sanity checks like this example enables us to rule out bias and lend credence to the model.

The major contributions of this abstract is: we offer a novel perspective of evaluating the model's accountability. Models could make correct predictions simply by exploiting spurious patterns from the data. Such models are susceptible to distribution shifts, and therefore less reliable in a changing clinical environment once deployed. Traditional metrics typically do not examine which features models use to make predictions, but rather focus on evaluating performance statistics. Thus, it is necessary to zoom into the model's behavior at the vocabulary level. A divergence between clinicians and the model in word significance rankings could indicate liability to distribution shifts or potential insights overlooked by humans.

## 2. Methods

In our study, we test the sensitivity of two text-based readmission predictors: a finetuned BioClinicalBERT and an xgboost on tf-idf embeddings. These models' overall performances are evaluated by their AU-ROC (The area under the receiver operating characteristic) and AUPRC (The area under the precision recall curve) scores on the test dataset. Their accountabilities are evaluated by the correlation of their sensitivity rankings on a set of 49 hand-selected words with those of physicians.

#### 2.1. Dataset

Both of the models are trained on a readmission dataset derived from MIMIC-III, a database of electronic health records from ICU patients at the Beth Israel Hospital (Johnson et al.). Each note has a binary label for readmission. A positive label suggests readmission within 30 days following the patient's discharge. The labeled dataset has 6% positive labels, with 52,725 examples and a 70% train, 15% validation and 15% test split. For more details see Appendix A.

For assessing model robustness and generalizability, a second readmission dataset is obtained from the NYU Langone Health System as part of an IRB approved study of inpatient readmissions. The NYU readmissions dataset consists of 45,120 examples from all clinical departments with 10.67% of positive labels.

#### 2.2. Readmission Classifier

Language-model based classifier: Bio-ClinicalBERT is a transformer encoder model pretrained on PubMed abstracts and MIMIC-III. It uses the original vocabulary of BERT and a Wordpiece tokenizer (Song et al.).

We finetune the pretrained model on our readmission dataset for 10 epochs and searched for a learning rate that gives the lowest validation loss. We use weighted cross entropy due to label imbalance. See appendix B.1 for more details.

For inference, we use a threshold of 0.35 to convert the redicted probabilities to binary labels (label 1 is assigned if the model's predicted probability is above 0.35 and vice versa) such that we reach 70% recall on the validation set.

**Baseline Model (tf-idf+xgb):** As a baseline model, we build a xgboost classifier using tf-idf embeddings as features. We select xgboost to compare the language model based classifier with a traditional machine learning classifier. See Appendix B.2 for more details.

Tf-idf has a less informative embedding, but is faster to train. The tf-idf based xgboost model does not use self-attention and positional embedding to incorporate semantic context and word orders into the representation. This makes its embeddings simply reflect how important each word is to a note compared to the entire corpus. On the other hand, tf-idf has a linear computational complexity for training with respect to the input sequence length, whereas attentionbased model such as BioClinicalBERT has quadratic complexity.

## 2.3. Metric: Token Sensitivity Score

We present the token **sensitivity score**, a metric to gauge the difference between classifier outputs before and after text perturbations. We say a classifier is sensitive to a token if the probability output changes notably after we perturb that token.

In this work, we explore swapping words as the perturbation. In our experiment, we used 3 types of perturbations: the uniform perturbations, the 1-gram perturbations, and the context perturbations. The uniform perturbation replaces a token of interest with another token uniformly sampled from the vocabulary of BioClinicalBert. The 1-gram perturbation replaces a token of interest with one of the five most frequent tokens. The context perturbation replaces a token of interest with one of the five most likely words according to BioClinicalBert's masked language modeling.

To illustrate the intuition of sensitivity score, consider determining whether or not a patient needs to be hospitalized. If "the patient has stroke", then we think it's likely that the patient is hospitalized. After context perturbation, if the text is swapped to "the patient has flu", then the patient is probably fine. Since the swap changes our opinion significantly, we say "stroke" is a sensitive word for determining hospitalization. Now we perturb the statement by swapping "has" with "got". In this case, the semantic meaning does not change much, and we say "has" is a relatively nonsensitive word for ample is: determining hospitalization.

In order to compare the model's behaviors with human doctors, we further propose the token **sensitivity score rank**. For each token, we compute the token's rank in terms of sensitivity scores among all the target words. We can then use the correlation between a model's rank and human's rank to assess how much a model's decisions align with human doctors' beliefs. For more details on the mathematical formulation, see the remaining subsections.

#### 2.3.1. Notations

To formalize the token sensitivity score, we need to first introduce some notations. A classifier f is a function that takes in a note x comprising of a set of tokens sampled from a vocabulary and outputs a probability p. For example, a constant classifier that always predict 30-day readmission using medical discharge summaries is  $f_{\text{const}}(x) = 1$ .

A note x is a sequence of n tokens/words  $(w_1, w_2, \ldots, w_n)$ . For example, a note could be ("his", "mom", "visited").

We can perturb a note with a **perturba**tion function q, which is parameterized by a token of interest u and a perturbation filter h. The perturbation function g replaces the first occurrence of the token of interest uwith a perturbed token given by the filter h. We only replace the first occurrence (denoted as one-swap) to reduce the impact of token frequency, which positively correlates with changes in predicted probabilities as shown in Appendix E.1.

For example, if we are interested in how sensitive a classifier f is to seeing "dad" rather than "mom", then our word of interest is u = "mom", our perturbation filter is h(``mom'') = ``dad'', and a perturbation ex-

$$g_{u,h}((\text{"his"},\text{"mom"},\text{"visited"})))$$
$$=(\text{"his"},\text{"dad"},\text{"visited"}).$$

More generally,

$$g_{u,h}(x) = g((w_1, \dots, w_n))$$
  
=  $(w_1, \dots, w_{k-1}, h(w_k), w_{k+1}, \dots, w_n),$ 

where k is the first location where the word of interest u appears. Note that for a fixed perturbation function  $g_{u,h}$ , the input x must contain u. Otherwise, the perturbation function does not change the note.

## 2.3.2. QUANTIFYING SENSITIVITY W.R.T. PERTURBATION FUNCTION

With a perturbation function  $g_{u,h}$ , we can quantify the sensitivity as the difference in predicted readmission probabilities before and after perturbing the token of interest uwith the filter h, as measured by  $L_1$  distance:

$$d_f(g_{u,h})(x) = |f(x) - f(g_{u,h}(x))|$$

For example, given note x = ("his", "mom", "visited"), the example classifier  $f_{\text{const}}$ , and the perturbation function  $g_{u,h}$  defined in section **D**, we have

$$d_{f_{\text{const}}}(g_{u,h})(x) = |1 - 1| = 0,$$

since the constant classifier always predicts positive labels regardless of the input text.

## 2.3.3. Averaging across Different Perturbations and Notes

We are interested in a general-case estimate of the sensitivity of a classifier with respect to perturbing a token. This motivates us to look at the difference in predicted probability across different perturbations and in different notes. For example, the change ("his mom visited"  $\rightarrow$  "his dad visited") might lead to a smaller difference than ("his mom is pregnant" $\rightarrow$  "his dad is pregnant") because the semantic change is less substantial.

To approximate the sensitivity of classifier f with respect to various perturbations, we consider a set of perturbation filters  $H = \{h_1, \ldots, h_m\}$ . The choice of H is up to the users. We use uniform filters to introduce randomness in perturbations. To limit the degree of perturbations out of the training distribution, we add filters based on distributions induced from the training set, i.e., the 1-gram distribution of a well-trained masked language model.

We define the **note-level sensitivity** score with respect to a token u in a note xas the average difference in predicted probabilities after applying a filter in H.

$$\overline{d_f}(g_{u,H})(x) = \frac{1}{m} \sum_{i=1}^m d_f(g_{u,h_i})(x)$$

To approximate the sensitivity across various notes, we consider a set of notes  $X = \{x_1, \ldots, x_l\}$  and define the **overall sensitivity score** with respect to a token u as the average of note-level sensitivity score:

$$\overline{d_f}(g_{u,H})(X) = \frac{1}{l} \sum_{j=1}^l \overline{d_f}(g_{u,H})(x_j)$$

## 2.3.4. Comparing Sensitivity of Different Models

Different models  $f_1$ ,  $f_2$  have respective ranges of output probabilities, making it unfair to directly compare the changes in predicted probabilities. For example,  $f_1$  could mostly predict  $p \in [0.1, 0.3]$ , whereas  $f_2$  predict  $p \in [0.4, 1]$ . In this case, comparing the overall sensitivity scores of  $f_1$  and  $f_2$  is biased towards the conclusion that  $f_2$  is more sensitive, because its range of outputs is higher.

To address this issue, we compare relative sensitivity as opposed to the absolute sensitivity. That is, we want to say whether  $f_1$  is *more* sensitive to u compared to  $f_2$ . We measure relative sensitivity with the sensitivity rank. Given a set of token of interest U, we calculate f's sensitivity scores of each token within it. A token u has a sensitivity rank of  $r_{u,f}$  if it has the  $r_{u,f}$ -th highest sensitivity score among U.

Now we can compare relative sensitivity of different classifiers using the sensitivity ranks. For example, we say  $f_1$  is more sensitive to u than  $f_2$  if  $r_{u,f_1} < r_{u,f_2}$ .

#### 2.3.5. EXAMPLE WITH MIMIC-III

We choose 49 tokens of interest to evaltwo readmission classifiers, uate our For each  $f_{\text{BioClinicalBERT}}$  and  $f_{\text{tfidf+xgb}}$ . token u of interest, we calculate the overall sensitivity score  $d_f(g_{u,H})(X)$ , with H as one of the fifteen perturbation filters, and X as the subset of our dataset that contains the token u. Our perturbation filters H is partitioned into 3 sets of 5 replacement filters: the uniform perturbations, the 1-gram perturbations, and the context perturbations. For more details check Appendix E.2.

## 3. Experiment Results

Tf-idf+xgb has slightly better AUC than BioClinicalBERT. As shown in Table 1, The standard deviation of BioClinicalBERT is 2.77 times that of tf-idf+xgb in AUROC, and 3.25 times that of tf-idf+xgb in AUPRC. While the marginal advantage of tf-idf+xgb suggests that word order and semantic context are not crucial, our next result shows that tf-idf+xgb's predictive behaviors align worse with the physicians.

Language-model based readmission classifier correlates better with clinicians' sensitivity. To investigate which model is more reliable, we select 49 target words based on professional inputs and collect the sensitivity rankings of 3 readmission predictor: the finetuned BioClinicalBERT, tfidf+xgb, and 3 human clinicians. We in-

Model	AUROC	AUPRC
bioclinical tf-idf+xgb	0.6995±0.0036 <b>0.7150</b> ±0.0013	$\begin{array}{c} 0.1224{\pm}0.0065\\ \textbf{0.1340}{\pm}0.0020\end{array}$

Table 1: Comparison of readmission AUC between BioClinicalBERT and tfidf+xgb. Test statistics are from 5 trials with distinct random seed.

vite three clinicians to rank 49 target words with a score from 1 to 5 to reflect the significance of each word in readmission prediction. A higher rank (smaller number) indicates a more important word for decision making. We obtain the overall clinician rankings by averaging the ranks across all three clinicians. We then assess the model's similarity to professional judgements through the Spearman rank correlation (Spearman) between models' rankings and the overall clinician ranking. Table 2 shows that BioClinical-BERT has a higher rank correlation, despite a slightly lower AUC in Table 1. Check Appendix E.3 and Appendix F for more details.

classifier	${\rm rank}_{-}{\rm correlation}$
BioClinicalBERT	<b>0.5754</b>
tf-idf+xgb	0.1259

Table 2: Spearman rank correlation between the two classifiers' rankings and the physicians' ranking.

In addition to the quantitative difference in Spearman rank correlations, to understand the difference between BioClinical-BERT and tfidf+xgb, we perform a qualitative analysis on where these models disagree based on Table 8 and Table 9. Words like "tumor", "pancreatic", and "dementia" are considered important by clinicians and BioClinicalBERT but neglected by the tfidf+xgb model. Meanwhile, words like "increase", "prescribed", and "blood" are considered important by the tfidf model but trivial by the other two. The inconsistency between the rankings of clinicians and the tfidf+xgb model shows its potential reliance on spurious features, rendering it less robust.

To empirically verify that the tfidf+xgboost model is less robust, we make zero-shot inferences on a held-out readmission dataset from NYU Langone Health. The result in Table 3 partially verifies the legitimacy of our sensitivity metric's ability in evaluating robustness.

Model	AUROC	AUROC_drop
bioclinical tf-idf+xgb	0.000	$0.0665 \\ 0.11$

Table 3: tfidf+xgboost has a greater drop in AUROC than BioClinicalBERT when inferencing on new data

## 4. Discussions

**Limitation:** Our metric has a large spatial and computational complexity. Generally, the complexity is worse than NM, where N is the dataset size and M is the number of perturbations. We can reduce complexity with the Monte Carlo method: subsampling the dataset and the perturbation filters. Methods of improving computational inefficiency is a direction of future research.

**Implications:** BioClinicalBERT might use the semantic context to extract more holistic information from a patient, as opposed to relying on statistical correlations that do not apply to specific subgroups. For example, an "increase" of blood pressure might be dangerous for patients with heart problems, but is a sign of recovery for patients with low blood pressure.

## Acknowledgments

We would like to acknowledge Michael Costantino, Ph.D. and Kevin Yie, M.S., from the NYU Langone High Performance Computing (HPC) team. Without their tireless assistance in building and maintaining our GPU cluster, none of this research would have been possible. We would also like to thank Ben Guzman from the NYU Langone Predictive Analytics Unit and Vincent J. Major from NYU Grossman School of Medicine for their help with learning the SQL data structures used as part of this work. This work is supported in part by NSF under grants 1922658. Additional funding comes from NYU Grossman School of Medicine.

## References

- Emily Alsentzer, John R. Murphy, Willie Boag, Wei-Hung Weng, Di Jin, Tristan Naumann, and Matthew B. A. McDermott. Publicly available clinical BERT embeddings. URL http://arxiv.org/ abs/1904.03323.
- Tianqi Chen and Carlos Guestrin. XGBoost: A scalable tree boosting system. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pages 785–794. doi: 10.1145/2939672.2939785. URL http:// arxiv.org/abs/1603.02754.
- Kexin Huang, Jaan Altosaar, and Rajesh Ranganath. ClinicalBERT: Modeling clinical notes and predicting hospital readmission. URL http://arxiv.org/abs/1904. 05342. Number: arXiv:1904.05342.
- Stuart Jackson, Maha Yaqub, and Cheng Xi Li. The agile deployment of machine learning models in healthcare. 1:7. ISSN 2624-909X. doi: 10.3389/fdata.2018.00007. URL https:

# //www.frontiersin.org/article/10. 3389/fdata.2018.00007/full.

- Mehdi Jamei, Aleksandr Nisnevich, Everett Wetchler, Sylvia Sudat, and Eric Liu. Predicting all-cause risk of 30-day hospital readmission using artificial neural networks. 12(7):e0181173. ISSN 1932-6203. doi: 10.1371/journal.pone. 0181173. URL https://dx.plos.org/10.1371/journal.pone.0181173.
- Alistair E.W. Johnson, Tom J. Pollard, Lu Shen, Li-wei H. Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. MIMIC-III, a freely accessible critical care database. 3(1): 160035. ISSN 2052-4463. doi: 10.1038/ sdata.2016.35. URL http://www.nature. com/articles/sdata201635.
- Richard Liaw, Eric Liang, Robert Nishihara, Philipp Moritz, Joseph E Gonzalez, and Ion Stoica. Tune: A research platform for distributed model selection and training. arXiv preprint arXiv:1807.05118, 2018.
- Scott M Lundberg and Su-In Lee. A Unified Approach to Interpreting Model Predictions. page 10.
- Vinodkumar Prabhakaran, Ben Hutchinson, and Margaret Mitchell. Perturbation sensitivity analysis to detect unintended model biases. URL http://arxiv.org/ abs/1910.04210.
- Xinying Song, Alex Salcianu, Yang Song, Dave Dopson, and Denny Zhou. Fast WordPiece tokenization. URL http:// arxiv.org/abs/2012.15524.
- Karen Sparck Jones. A STATISTICAL INTERPRETATION OF TERM SPECI-FICITY AND ITS APPLICATION IN RETRIEVAL. 28(1):11–21. ISSN 0022-0418. doi: 10.1108/eb026526. URL https:

//www.emerald.com/insight/content/ doi/10.1108/eb026526/full/html.

- C. Spearman. The proof and measurement of association between two things. 15 (1):72. ISSN 00029556. doi: 10.2307/ 1412159. URL https://www.jstor.org/ stable/1412159?origin=crossref.
- Carl van Walraven, Jenna Wong, and Alan J Forster. LACE+ index: extension of a validated index to predict early death or urgent readmission after hospital discharge using administrative data. page 11.
- Cao Xiao, Edward Choi, and Jimeng Sun. Opportunities and challenges in developing deep learning models using electronic health records data: a systematic review. 25(10):1419–1428, a. ISSN 1067-5027, 1527-974X. doi: 10.1093/jamia/ocy068. URL https://academic.oup.com/ jamia/article/25/10/1419/5035024.
- Cao Xiao, Tengfei Ma, Adji B. Dieng, David M. Blei, and Fei Wang. Readmission prediction via deep contextual embedding of clinical concepts. 13(4):e0195024,
  b. ISSN 1932-6203. doi: 10.1371/journal. pone.0195024. URL https://dx.plos. org/10.1371/journal.pone.0195024.

## Appendix A. Data

## A.1. Data Source

We use three tables from the MIMIC-III in our study: patients, admissions, and noteevents. The three tables have "subject\_id", "hadm\_id", and "row\_id" as their respective primary keys. We use the discharge note as our input since it contains the most information as a summary of the entire visit.

## A.2. Label Generation

For each distinct discharge note associated with an encounter, we generate a binary readmission label based on the patient's medical record. A positive label suggests readmission within 30 days following the patient's discharge.

We generate the readmission labels as follows: for each patient, we order their encounters by admission time. For each of their encounters, we calculate the readmission interval between the discharge time for the current visit and the subsequent admission time. If an encounter does not have subsequent admissions or the readmission interval is longer than 30 days, we assign a negative label. Otherwise, we assign a positive label.

To properly handle boundary cases, we only consider encounters that are discharged at least a month before the latest admission time. This prevents false negative labels with unobserved readmissions outside the dataset.

## Appendix B. Training Details

#### **B.1.** Finetuning BioClinicalBERT

In finetuning, we search learning rate from {2e-5, 2e-6, 2e-7, 2e-8, 2e-9} using ray-tune (Liaw et al., 2018). For each learning rate, we finetune the model for 10 epochs using a Nvidia-3090 GPU for around 140 minutes. We select the model with best validation loss, which uses a learning rate of 2e-5 and 5 epochs.

To address the label imbalance issue, we use a weighted cross entropy loss function to increase the penalty for misclassifying the minority class. Specifically, we weigh the positive example with the ratio of negative examples in the entire dataset; Similarly, we weigh the negative example with the ratio of positive examples.

#### B.2. Training tf-idf+xgb

The model first convert the input texts to a word-count vector, then calculate the tfidf embedding. Next, xgboost uses this embedding matrix for binary classification. We repeat the training five times with distinct random seeds (24, 42, 61, 67, 70) for a total of three minutes.

## Appendix C. Figures



Figure 1: Comparison of BioClinicalBERT word sensitivities using the oneswap and multiple-swap scheme on readmission prediction. The oneswap scheme counteracts the effect of frequency on sensitivity score.

#### Appendix D. Metric Formalization

## **D.1.** Previous Works

Existing metrics for model interpretation mainly focus on unraveling individual predictions or mimicking local behaviors of the model. However, these methods have two limitations: first, they focus on attributing decision making to local features in individual examples; second, they cannot quantify explainability as measured by alignment with some ground-truth references. For example, SHAP (Shapley Additive exPlanations) (Lundberg and Lee) may be able to tell us that "alcohol" is important for predicting that an individual depressed patient



Figure 2: Comparison of tf-idf word sensitivities using the one-swap and multiple-swap scheme on readmission prediction. The one-swap scheme counteracts the effect of frequency on sensitivity score.

would be readmitted. But SHAP cannot tell us whether "alcohol" is a key feature for predicting readmission over an entire readmission dataset (not necessarily the case, because the use of "rubbing alcohol" indicates good hygiene). Further, if SHAP tells us that "alcohol" is important for an individual prediction, we cannot interpret this finding without checking its alignment with some experts. One way is to consult some physicians: "Is it good that my model thinks alcohol is an important feature?" Our proposed sensitivity score fills this gap by: first, quantifying the sensitivity of each token over an entire dataset (e.g., overall, "alcohol" is not a sensitive feature); second, quantifying the alignment of model's sensitivity with a reference (e.g., overall, language model's sensitivity to words are more similar to trustworthy physicians than tf-idf models).

## Appendix E. Derivation details

## E.1. Reasons for Substituting the First Occurrence

This section explains the rationale behind word substitution only once regardless their total occurrences.

Previously, we substitute all occurrences of any target word (denoted as multiple-swap) and found that frequent words in general have smaller sensitivity rankings, i.e., higher sensitivity scores. The sensitivity ranking and word frequency are correlated with a Pearson coefficient of -0.62 for tf-idf as shown in Figure 2 and Appendix F. The Pearson coefficient is -0.60 for BioClinicalBERT, as shown in Figure 1 and Appendix F. Thus, the metric is biased if we compare sensitivities of tokens with diverging frequencies.

To tackle this problem, we substitute the *first* occurrence of the target word (called "one-swap" for brevity) rather than multiple-swap. Figure 2 and 1 shows that with this method, the sensitivity ranking are distributed more evenly across word frequency. Nevertheless, it should be acknowledged that under rare circumstances, swapping the first occurrence could introduce a small bias when its particular context lend higher importance to this occurrence than others.

To get some insights about the correlation between word frequency and sensitivity, we consider a toy example of a linear classifier using the TD-IDF embedding. Suppose we have an input

$$x = (\text{"hi"}, \text{"there"})$$

and a perturbed input with respect to the token of interest u= "there":

$$g(x) = (\text{``hi''}, \text{``hi''})$$

The tf-idf+linear classifier is defined as

$$f_{\text{tf-idf+linear}}(x) = \operatorname{softmax}(W\phi(x)).$$

Here  $\phi$  is the tf-idf embedding function with two vocab: "hi" and "there":

$$\phi(x) = \begin{bmatrix} \operatorname{tf}(\text{``hi''}, x) \log(\frac{|X|}{\operatorname{df}(\text{``hi''}, x)}) \\ \operatorname{tf}(\text{``there''}, x) \log(\frac{|X|}{\operatorname{df}(\text{``there''}, x)}) \end{bmatrix},$$

where term frequency tf(w,x) is the number of w in a note x and document frequency df is the number of notes  $x \in X$  that contain token w.

The key insight is that a higher term frequency of the token of interest tf(u, x) would lead to a larger change in  $L_1$  norm of the perturbed tf-idf embedding. For example, both "hi" and "there" appears once in x, so

$$\phi(x) = \begin{bmatrix} \log(\frac{|X|}{\mathrm{df}(\text{``hi'',x)}}) \\ \log(\frac{|X|}{\mathrm{df}(\text{``there'',x)}}) \end{bmatrix}.$$

After perturbation, "there" completely disappear while there are 2 "hi"s, so

$$\phi(g(x)) = \begin{bmatrix} 2\log(\frac{|X|}{\mathrm{df}(\text{``hi'',x)}}) \\ 0 \end{bmatrix}$$

The difference in tf-idf embedding is

$$\begin{aligned} \|\phi(x) - \phi(g(x))\|_1 \\ = \log(\frac{|X|}{\mathrm{df}(\mathrm{``hi'',x)}}) + \log(\frac{|X|}{\mathrm{df}(\mathrm{``there'',x)}}) \end{aligned}$$

More generally, given a word of interest u with term frequency n, the difference in tf-idf embeddings scales with n:

$$\begin{aligned} \|\phi(x) - \phi(g(x))\|_1 \\ &= n \log(\frac{|X|^2}{\mathrm{df}(\text{``hi"}, \mathbf{x})\mathrm{df}(\text{``there"}, \mathbf{x})}) \\ &\implies \|\phi(x) - \phi(g(x))\|_1 \propto n \end{aligned}$$

In our toy example, since f is linear, we know the difference in predicted probability would be larger if the perturbed word of interest has a larger term frequency n:

$$|f(x) - f(g(x))| \propto ||\phi(x) - \phi(g(x))||_1 \propto n$$
(1)

On average, a token with a higher word frequency (meaning that it appears more often in the dataset) have a higher term frequency in each note. By Equation 1, such high-frequency token has a higher sensitivity score after perturbation.

#### E.2. Perturbation Filter

Each filter in uniform perturbations replaces the token of interest u with another token uniformly sampled from the vocabulary of BioClinical Bert.

$$U = \{h : h(u) = w', w' \in W_{\text{uniform}}\}$$

Each filter in 1-gram perturbations replaces the token of interest u with one of the five most frequent tokens from the subset Xthat contains u.

$$G = \{h : h(u) = w', w' \in W_{1-\text{gram}}\}$$

Each filter in context perturbations replaces the token of interest u with one of the five most likely predictions according to nonfinetuned BioClinicalBERT's MLM probability. We specifically used the *non-finetuned* BioClinicalBERT to avoid using the same model for both prediction and assessment.

$$O = \{h : h(u) = w', w' \in W_{\text{context}}\}$$

#### E.3. Spearman Rank Correlation

$$r_s = 1 - 6 \cdot \frac{\sum D^2}{n(n^2 - 1)} \tag{2}$$

We use the Spearman rank correlation (Spearman) to quantify the divergence between the model ranking and the manual ranking. In Equation 2, D is the difference between ranks, and n is the number of pairs of data.

## Appendix F. Tables

The sensitivity ranking and word frequency are negatively correlated for both the tf-idf model and BioClinicalBERT when we swap multiple occurrences.

Despite similar AUROC scores, BioClinicalBERT is more sensitive to disease names compared to general words, while tf-idf+xgb displays a more irregular distribution of sensitivity.

Table 4 and Table 5 display the Sensitivity Score of 10 words with different frequencies with respect to the language model and the tfidf+xgboost model using multiple swaps of occurrences. The one-swap counterparts are shown in Table 6 and Table 7. To assess the models' alignment with professional opinions, Table 8 lists the token sets and words' relative significance ranking.

word	test_fre	$\overline{d_f}(g_{u,H})(X)$	rank
cancer	1912	0.033884	1
mg	58910	0.027513	2
colon	395	0.020076	3
expired	1234	0.018921	4
deceased	399	0.017083	5
heparin	1898	0.014974	6
died	1871	0.012532	7
father	1890	0.007779	8
mother	1897	0.006713	9
mouthwash	78	0.005541	10
regimen	395	0.004930	11
congenital	78	0.002269	12
thinner	78	0.002439	13

word	test_fre	$\overline{d_f}(g_{u,H})(X)$	rank
hypoglycemia	168	0.025779	1
fall	789	0.021095	2
ulcer	358	0.011134	3
prematurity	164	0.010929	4
arthritis	158	0.009826	5
father	1890	0.007510	6
mother	1897	0.006404	7
patient	7900	0.005129	8
blood	6565	0.004980	9
labor	16	0.003916	10
vaccination	78	0.001465	11

Table 4: The language model's Sensitivity – Score of 10 words within different T frequency range (multiple-swap)

Table 6:	BioClinicalBERT's			Sensit	ivity		
	Score	of	a	list	of	ailments	and
	words	for	co	mpar	isor	n (one-swa	ιp)

word	${\bf test\_fre}$	$\overline{d_f}(g_{u,H})(X)$	rank
expired	1234	0.019003	1
mg	58910	0.007416	2
heparin	1898	0.000583	3
mouthwash	78	0.005541	4
deceased	399	0.017083	5
died	1871	0.012532	6
cancer	1912	0.033884	7
regimen	395	0.004930	8
colon	395	0.020076	9
mother	1897	0.006713	10
father	1890	0.007779	11
thinner	78	0.002439	12
$\operatorname{congenital}$	78	0.002269	13

Table 5: The tf-idf + xgboost model's Sensitivity Score of 10 words within different frequency range (multipleswap)

word	${\bf test\_fre}$	$\overline{d_f}(g_{u,H})(X)$	rank
fall	789	0.003331	1
blood	6565	0.000481	2
hypoglycemia	168	0.000426	3
patient	7900	0.000342	4
ulcer	358	0.000319	5
prematurity	164	0.000249	6
arthritis	158	0.000248	7
mother	1897	0.000248	8
father	1890	0.000239	9
vaccination	78	0.000024	10
labor	16	0.000000	11

Table 7: tf-idf+xgb's Sensitivity Score of a list of ailments and words for comparison (one-swap)

word	manual	language	tf-idf	word	manual	language	tf-idf
chemotherapy	1.0	15	25	urinary	26.0	24	5
hypoglycemia	2.0	3	12	faint	26.0	46	42
tumor	3.0	19	40	refills	26.0	9	19
overdose	3.0	1	2	immunizations	26.0	38	39
dementia	3.0	2	16	blood	26.0	34	9
anticoagulation	6.0	36	36	family	26.0	35	37
delirium	6.0	26	17	diarrhea	32.0	16	22
debridement	6.0	7	10	female	32.0	18	44
$\operatorname{arrhythmia}$	6.0	5	20	prescribed	32.0	30	7
pancreatic	6.0	4	35	medication	32.0	45	29
amputation	6.0	10	23	electrolytes	32.0	40	43
fall	6.0	6	3	allergies	32.0	12	46
cardiovascular	13.0	17	45	aspirin	32.0	13	6
neurosurgery	13.0	31	33	increase	32.0	48	4
diabetes	13.0	39	24	tylenol	40.0	20	13
ablation	16.0	11	15	care	40.0	37	26
expired	16.0	21	1	benign	40.0	29	38
dehydrated	16.0	25	41	mother	40.0	28	31
palpitations	16.0	8	28	cartridge	40.0	44	14
obesity	16.0	41	34	labor	40.0	43	49
wheeze	21.0	14	27	moderate	40.0	49	47
vaccination	21.0	47	48	tablet	40.0	33	21
arthritis	21.0	22	30	mg	40.0	42	11
pain	21.0	27	32	patient	40.0	32	18
dysfunction	21.0	23	8				

Table 8: Sensitivity Score rankings of 49 hand-chosen words for model comparison (one-swap). The language model's sensitivity ranking aligns better with the clinicians' manual rankings. (first half) Table 9: Sensitivity Score rankings of 49 hand-chosen words for model comparison (one-swap). The language model's sensitivity ranking aligns better with the clinicians' manual rankings. (second half)