

1 **Supplemental Table 1. Dictionary Used in NLP-Lexicon Model**  
 2

<b>Terms and phrases indicative of professionalism lapses</b>
\b(make made)\b.*\bfun\b
(joke joking joked)
(insensitive)
(mean-spirited)
(overweight obese)
(however)
([""])- phrase in quotations
(bitch bitches)
((make made makes)\b.*\bcomments)
(resent resented resenting)
(intimidated intimidating afraid humiliating humiliated humiliate)
(disparage disparaging)
(uncomfortable embarrassment hostile angry)
(\bbody\b.*\bhabitus)
(behind patients' backs)
((not)\b.*\bbacknowledged)
((bad badly)\b.*\bI was treated)
(disinterested)
(malignant)
(active indifference)
(bully bullied)
(out of proportion)
(bros)
(yelling yell)
(favoritism)
((play)\b.*\bfavorites)
(felt concerned about)
((uncomfortable unfavorable hostile) (learning work) environment)
(commenting)
(unprofessional disrespectful rude rudest lack respect)
((concerned)\b.*\b(about)\b.*\b(ability))
((not)\b.*\b(safe))
(penis)
(problematic)
(toxic)
(judgmental)
(condescending condescendingly condescended condescend condescends)
(fearful)
(inappropriate inappropriately)
(stupid)
(impatient)
(shame shaming shamed)
(gossip gossiped)

(hatefully)  
(errands|errand)  
(concern|concerns)  
(defame|disdain|curse)  
(shocked)  
(tardy|tardiness|late)  
(reprimand|reprimanded)  
(personal life)  
(genitalia)  
(facebook)  
(preconceived notions)  
(dismissive)  
(punitive)  
(spoke ill)  
(unkind)  
(unwelcome)  
(abrasive)  
(confrontational)  
(demean|demeaning|demeaned|demeans)  
(not appropriate)  
(physical contact)  
(swatted)  
(misrepresentation|dishonest)  
(hit|shove|shoved)  
(slap|slapped)  
(bad-mouth|bad-mouthing)  
(derogatory)  
(racial slurs)  
(berate|berated|berates)  
((raise|raising)\b.\*\bvoice)  
(ignore|ignored|ignoring|ignores)  
(demoralized|demoralize|demoralizes|demoralizing)  
(avoid|avoided|avoiding|avoids)  
(sarcasm|sarcastic)  
(sexist)  
(homophobic)  
(nasty)  
(difficult to work with)  
(undermine|undermines|undermined|undermining)  
(discredits|discredit|discrediting|discredited)  
(biases)  
(frequent flyer|frequent flyers|frequent flier|frequent fliers)  
(passive aggressive)  
(patronizing)  
(hinder|hinder|hindered|hinders)

**Terms and phrases indicative of excellent professionalism/no professionalism lapses**

(role model)

(pleasure)

(amazing)

(awesome)

(fantastic)

(favorite)

(wonderful)

(admire)

(one of the best)

**Negation Terms**

(never)

(protect|protects|protected|protecting)

(not)

(rather than)

(instead of)

1

2

1 **Supplemental Table 2. Application of Positive Predictive Value (PPV) for Identification of**  
 2 **Professionalism Lapses.**

	<b>Professionalism Lapse Identified by Manual Review</b>	<b>No Professionalism Lapse Identified on Manual Review</b>	<b>Total</b>
<b>Professionalism Lapse Identified by NLP</b>	<b>A</b> 160	<b>B</b> 80	<b>(A+B)</b> 240
<b>No Professionalism Lapse Identified on NLP</b>	<b>C</b> 40	<b>D</b> 720	<b>(C+D)</b> 760

3  
 4 This table highlights the PPV for the current study. A represents the True positives, with  
 5  $A/(A+B)$  representing the PPV. In this example, the PPV of the NLP model for accurate  
 6 identification of a professionalism lapse would be 67% (160/240). This suggests that 67% of  
 7 NLP-identified professionalism lapses would be confirmed on manual review. The NPV of the  
 8 NLP model, in this example, would be 95% (720/760).  
 9

1 **Supplemental Digital Appendix 1. Additional Information on Natural Language Processing**  
2 **Methods and Predictive Modeling Approach**

3 We created features for the predictive models using several different approaches. First, we used  
4 term-frequency inverse document frequency weighted n-grams. We then selected the top 100 1-,  
5 2-, and 3-word phrases that appeared in at least 5 documents. Second, we converted words in  
6 each comment to a vector using the Spacy built-in 96-dimension word embedding model. Each  
7 comment was represented by a concatenation of the element-wise min, mean, and max over each  
8 word vector, thus creating a 288-dimension input vector. Third, we used a sentiment score for  
9 each comment using the first 250 tokens as input to Bidirectional Encoder Representations from  
10 Transformers (BERT). Finally, we created stacked ensemble models using the predicted  
11 probabilities of other models as inputs. Each of these featurization approaches was used  
12 separately or in combination as input variables into each model.

13  
14 We trained three types of models including a penalized logistic regression model with L1 and L2  
15 penalties, a random forest model, and a neural network model. Model tuning parameters and  
16 determination of model performance was determined by complete grid search using 5-times-  
17 repeated 10-fold cross validation to identify the model with the highest scaled Brier score (which  
18 provides an estimate of the calibration and discrimination of the model's performance).

19  
20 For additional information on natural language processing approaches, please refer to:

21       Nadkarni PM, Ohno-Machado L, Chapman WW. Natural language processing: an  
22       introduction. J Am Med Inform Assoc. 2011 Sep-Oct;18(5):544-51. doi:  
23       10.1136/amiajnl-2011-000464. PMID: 21846786; PMCID: PMC3168328.

1

2 For additional information on predictive modeling and performance metrics, please refer to:

3 Steyerberg EW, Vickers AJ, Cook NR, Gerds T, Gonen M, Obuchowski N, Pencina MJ,

4 Kattan MW. Assessing the performance of prediction models: a framework for traditional

5 and novel measures. *Epidemiology*. 2010 Jan;21(1):128-38. doi:

6 10.1097/EDE.0b013e3181c30fb2. PMID: 20010215; PMCID: PMC3575184.

7