

# Machine Learning Models for the Segmentation of Clinical Exchange

Mehedi Hasan, BS<sup>1\*</sup>, Alexander Kotov, PhD<sup>1\*</sup>, Ming Dong, PhD<sup>1</sup>, Sylvie Naar, PhD<sup>2</sup>, Gwen L. Alexander, PhD<sup>3</sup>, April Idalski Carcone, PhD<sup>4</sup>

<sup>1</sup>Department of Computer Science, Wayne State University, Detroit, Michigan

<sup>2</sup>Director, Center for Translational Behavioral Research, Department of Behavioral Sciences and Social Medicine, Florida State University, Tallahassee, Florida

<sup>3</sup>Department of Public Health Sciences, Henry Ford Health System, Detroit, Michigan

<sup>4</sup>Department of Family Medicine and Public Health Sciences, School of Medicine, Wayne State University, Detroit, Michigan

**Abstract** Communication science to understand clinical process like Motivational Interviews (MI) is limited by traditional qualitative coding methods. Qualitative coding of communication data is very resource-intensive and time-consuming process, requires auto coding. This study utilized e-coaching data where email is used to deliver motivation-enhancing coaching to encourage healthy eating. A critical step toward automatic annotation of communication coding process is the segmentation of text data. In this study, we transformed segmentation task into a classification task and developed several state-of-the-art machine learning models including Support Vector Machine, Naive Bayes, K-Nearest Neighbor (KNN), Recurrent Neural Networks by utilizing contextual, topic and punctuation mark features. Results indicate that KNN is the best model and achieved 0.986 F1-measure in overall, 0.779 and 0.993 F1-measures for detecting “boundary” and “not boundary” classes, respectively. This study has a great implication to save money and accelerate the pace of identifying effective communication strategies.

## Introduction

Communication science to understand clinical process like Motivational Interviews (MI) is limited by traditional qualitative coding methods. Motivational Interviewing (MI) is an evidence-based communication technique to increase intrinsic motivation and self-efficacy for behavior change<sup>1-3</sup>. Patient “change talk”, statements of intrinsic motivation about their desire, ability, reasons, need for and commitment to behavior change, is an established mediator of health behavior change<sup>4</sup>. Qualitative coding of communication data has been traditionally performed manually with pre-defined codes by trained coders, which is a tedious resource-intensive task, requiring an iterative process of human (subjective) interpretation of the text.

Rapidly developing computational technologies, specifically, machine learning models, offer a unique opportunity to accelerate this process. Machine learning-based classification methods have been successfully applied to a variety of analytical tasks on textual data including classification<sup>5</sup>, sentiment analysis<sup>6</sup>, and digital forensics<sup>7</sup>. A few studies have been done for the automatic annotation of clinical interactions in the setting of MI intervention. Lacson et al.<sup>8</sup> applied AdaBoost<sup>9</sup> classifier for annotating the interactions in hemodialysis phone dialog as Clinical, Technical, Backchannel, and Miscellaneous categories. Our research group has recently applied several machine learning-based models to MI sessions<sup>10,11</sup>. A simple communication code scheme was automated to characterize patient communication and achieved accuracy comparable to human coders<sup>10</sup>.

Experimental data utilized in the above MI studies were prepared by transcribing audio conversation which involves a counselor, a patient, and some cases a caregiver<sup>10,11</sup>. However, this study utilized e-coaching data which have different structure and context compared to traditional clinical data. E-coaches use email to deliver motivation-enhancing coaching to encourage healthy eating, grounded in the principles of motivational interviewing. E-coaching data is comprised of email responses which are unsegmented, unlike more traditional dyadic clinical communication where segmentation occurs naturally due to its conversation nature. Discourse analysis on e-coaching text can help to understand the communication process in e-coaching text by revealing the socio-psychological characteristics of a patient<sup>12-14</sup>. By utilizing e-coaching data and leveraging innovative machine learning models, the ultimate goal of this research study is to fully automate the communication coding process.

---

\* Authors provided equal contribution.

A significant barrier to fully automated the communication coding process is the segmentation of the text data. During traditional qualitative coding, coders determine where to segment text, that is, where one code ends and another begins. Automatic segmentation of e-coaching intervention sessions is a challenging task due to the 2 important reasons. First, the email is an unstructured text that contains informal email exchange as a non-traditional format. Second, a text segment does not necessarily belong to the entire sentence or collection of sentences. One sentence can be segmented into multiple MI behaviors or several sentences may represent a single MI code. Figure 1 illustrates the segmentation of an e-coaching email exchange where the first sentence segmented into 2 different MI behaviors. On the other hand, fourth and fifth segments contain 1 and 3 sentences, respectively.

---

On Mon Nov 10 20:40:02 2014, XXX wrote:

(Hi XXX, I haven't had a chance to look through MD 5 or 6, but I've found a few veggies that I like to pack and take with me. I just have to prep them more. Thanks XXX)

---

(Email Date: 2014-11-11 10:29:18)

Hi XXX,

It's good to hear from you. It sounds like you found a plan that works for you as long as you are able to find time to prep veggies for on-the-go snacks. Sometimes people find inspiration for making a change by considering things that are important to them. There is some evidence that behavior change is often easier when it relates to your own values and goals. This might be helpful in finding reasons to keep up with what you are now doing. You stated that being considerate, respected, and responsible are important to you. How, if at all, would you say that eating better and having more energy would help you be considerate and respected? How about to be more responsible?

I look forward to hearing from you again soon,

YYY

---

**Figure 1: Segmentation of e-coaching text depicts the main challenges of boundary detection.**

In this paper, we address the text segmentation problem by developing several state-of-the-art machine learning models to promote the automatic identification of best communication strategies without human interference. More specifically, we developed Support Vector Machine (SVM), Naive Bayes (NB), K-Nearest Neighbor (KNN), Long Short Term Memory (LSTM), and Gated Recurrent Unit (GRU) by utilizing contextual, topic and punctuation mark features, to find the best model for the segmentation of e-coaching text.

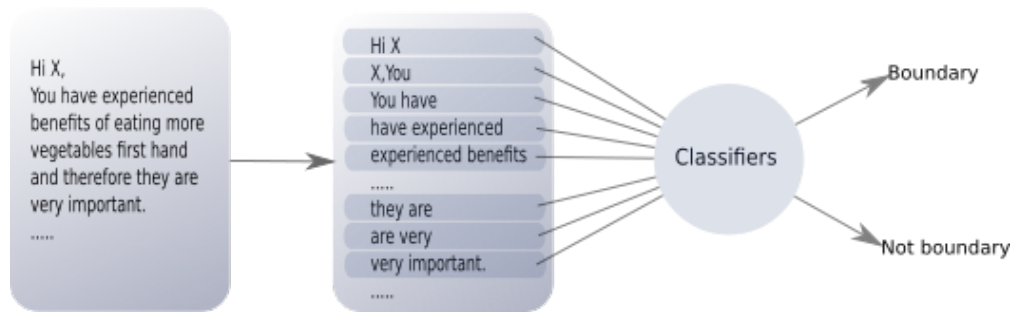
Previous studies mainly focus on segmentation of text into sections and headers<sup>15-18</sup> or sentence boundary detection<sup>19-21</sup> in the medical domain. Apostolova et al.<sup>15</sup> applied SVM by utilizing word-vector cosine similarity metric combined with several heuristics to classify clinical report into semantic sections such as demographics, history, exam procedure, finding, impression, etc. After identification of each line in the document, Tepper et al.<sup>17</sup> trained Maximum Entropy models for the section classification. In 2009, Denny et al.<sup>16</sup> proposed a SecTag algorithm, which combined natural language processing technique, terminology-based rule, and naive Bayesian score for identifying sections and headers that achieved 99% recall with 95.6% precision. On the other hand, SVM exploiting with linear kernel and recurrent convolutional neural networks with prosodic, part of speech features and word embeddings, were trained by Kreuzthaler et al.<sup>20</sup> and Griffis et al.<sup>19</sup>, respectively, for the detection of sentence boundary.

Recently, an online clinical intervention called MENU GenY<sup>22</sup> (Making Effective Nutrition Choices for Generation Y) was proposed to test its efficacy. MENU GenY is a technology-based public health intervention to encourage increased fruit and vegetable intake among young adult age 21-30, utilized personalized e-coaching data. The goal of MENU GenY was to develop a better coding dictionary among GenY to improve eating habits. However, segmentation of clinical conversation in the context of electronically delivered interventions, in particular, segmentation of clinical interaction text into groups of MI behaviors, is still ignored while relying on manual hand-coded approach. Therefore, this study introduces a novel approach and the authors are not aware of any other work this approach has been considered for the segmentation of e-coaching text.

## Methods

### Data collection

The experimental dataset for this work was constructed from 49 e-coaching sessions, which include a total of 3,138 segmented and annotated MI behaviors. Each session represents an MI intervention delivered via email. To filter out noise from the dataset, non-ascii characters are removed and then applied stemming to obtain a general form of word from different word representations, such as “eating”, “eats”, and “eat”. We formulate the text segmentation task into a binary classification, as shown in Figure 2. Clinical exchange is given as the input, it is partitioned as adjacent word pairs by sliding them. Each pair is classified into either “boundary” or “not boundary” class. The original text is segmented at the position, where an adjacent word pair classified into “boundary” class. If all adjacent pairs of a block of text classified into “not boundary”, the whole block of text is then treated as one segment about a single MI behavior. Totally, we obtained 95,421 word pairs, which include 3,138 “boundary” and 92,283 “not boundary” instances.



**Figure 2: Transformation of text segmentation task into text classification task.**

For the experiment, we utilized three type of features including word (textual feature), topic, and punctuation mark. Each word represented in a binary format, where 1 indicates the appearance of the word and 0 for absence. Topics are considered as features since topic models are very effective<sup>11,23,24</sup> to represent text documents. In this paper, we exploit a topic model named Labeled LDA<sup>25</sup> model, where topics correspond to labels. The generative process of Labeled LDA drawing the multinomial topic distributions over vocabulary for each topic. In our experiment, we represent each word in a vector of 2 topics, where the number of topics is experimentally determined by the model performance<sup>11</sup>. Punctuation mark containing one of the symbols {':', ',', '!', '?', ':', ';', '-'} is also employed as feature. This is one of the most important features as they indicate the boundary of a sentence, clause, and phrase.

### Segmentation classifiers

Several state-of-the-art classifiers, including Naive Bayes (NB)<sup>26</sup>, Support Vector Machine (SVM)<sup>27</sup>, K-Nearest Neighbor (KNN)<sup>26</sup>, two variant of Recurrent Neural Networks (RNN)<sup>28</sup>: Long Short Term Memory (LSTM)<sup>29</sup> and Gated Recurrent Unit (GRU)<sup>30</sup>, are employed to estimate the classification performance.

**Naive Bayes:** this model is constructed by using the training data and estimate the prior probability of classes, and each feature has given the class. Then, the posterior probability is computed to predict the class label by applying the Bayes theorem with the assumption that features are conditionally independent. This study utilized a specialized version of Naive Bayes called Multinomial Naive Bayes, which is best suitable for discrete features such as word.

**Support Vector Machine:** we used this model as one of the state-of-the-art classification technique proven to perform well in text categorisation<sup>31</sup> for its ability to cope with very high dimensional input feature space. SVM finds the best hyperplane in the feature space that maximizes the separation between the closest “boundary” and “not boundary” training examples. In this experiment, the polynomial kernel is employed to train the SVM model for the segmentation of e-coaching text.

**K-Nearest Neighbor:** By this model, each training sample represented as a point in the input feature space. For a new

test sample, Euclidean distance is calculated to find the k-nearest neighbors. Finally, the test sample is classified into majority class of the k-nearest neighbors. We experimentally determined that best performance was achieved with  $k = 3$  for the classification of word pairs.

**Recurrent Neural Networks:** RNN is a neural network architecture designed to capture sequential patterns present in temporal sequence such as text data. When we predict the “boundary” point, adjacent word pair will help to understand the pattern of the sequence. Long Short Term Memory networks usually referred as LSTMs<sup>29</sup>, are a special type of RNN capable of handling variable size input sequence, contains internal memory. GRU<sup>30</sup> is a variant of LSTM mathematically represented by the following formula:

$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z) \quad (1)$$

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r) \quad (2)$$

$$\tilde{h}_t = \tanh(W_h x_t + r_t \odot U_h h_{t-1} + b_h) \quad (3)$$

$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t \quad (4)$$

In Eq. 1-4,  $\sigma$  corresponds to sigmoid function and  $\odot$  designates an element-wise product. The update gate  $z_t$  and reset gate  $r_t$  at time step  $t$  are computed by the Eq. (1) and (2), where  $W_z$ ,  $W_r$ ,  $W_h$ ,  $U_z$ ,  $U_r$ ,  $U_h$  are the weight matrices and  $b_z$ ,  $b_h$  and  $b_r$  are bias vectors. The activation  $h_t$  of the GRU at time  $t$  is a linear combination of previous activation  $h_{t-1}$  and the candidate activation  $\tilde{h}_t$ , which is represented by Eq. (4) and (3). We build our RNN model with one hidden layer, output layer, and input layer. One segmented text was taken as input sequence which get one hot encoding of word vector for the model input. A binary label was assigned as output sequence where “boundary” and “not boundary” words were associated with label 1 and 0, respectively. We reset our model state after feeding each input sequence. Since one-hot vector is given in the input layer, results are reported with textual and punctuation mark features only. We experimentally determined that the best performance is achieved when the number of hidden units = 25, batch size = 1, and optimizer = adam.

### Evaluation metrics

In this experiment, standard metrics: precision, recall, and F1-measure, are applied to evaluate the performance of binary classifiers<sup>32</sup>. However, accuracy is not reported as a performance metric because accuracy is highly sensitive to the prior class probabilities and does not fully describe the actual difficulty of the decision problem for an unbalanced dataset. We conduct the experiment with 5 folds cross-validation and weighted macro-averaging of these metrics over the folds. All models have trained on 80% of the word pairs and remaining 20% of the data is used as a test set for reporting the performance of the model. We also estimate the area under the receiving operating characteristics (ROC) curve<sup>33</sup> (AUC) metric due to its effectiveness in measuring the quality of binary classifiers for imbalanced datasets<sup>34</sup>.

### Results

Experimental results are evaluated with “boundary” and “not boundary” classes as well as their weighted average, which are shown in Table 1, 2, and 3, respectively.

**Table 1:** Performance of NB, SVM, KNN, and RNN methods for detecting segmentation boundary in e-coaching text. The highest value for each performance metric is highlighted in bold.

Method	contextual features only			contextual + punctuation marks (+ topics except RNN)		
	Precision	Recall	F1-measure	Precision	Recall	F1-measure
NB	0.594	0.662	0.626	0.590	0.666	0.626
SVM	0.742	<b>0.679</b>	0.709	0.774	0.696	0.733
KNN	<b>0.808</b>	0.663	<b>0.728</b>	<b>0.820</b>	<b>0.742</b>	<b>0.779</b>
LSTM	–	–	–	0.800	0.646	0.714
GRU	–	–	–	0.800	0.715	0.741

As follows from Table 1, KNN performs best among all machine learning models in terms of precision and F1-measure, achieved 0.808 precision with 0.728 F1-measure when contextual features are used, and 0.820 precision with 0.779 F1-measure when a combination of contextual, topic, and punctuation mark features are used. However, NB demonstrates the lowest performance among all models in terms of all performance metrics. In this study, GRU appears as the second highest model, obtains 10.68% higher recall with 3.78% higher F1-measure than LSTM. On the other hand, SVM exhibits highest recall 0.679 when only textual features are used. When textual features are used in combination with topic and punctuation mark features, recall increases by 0.6%, 2.5%, and 11.92%; and F1-measure increases by 0%, 3.39%, and 7% for NB, SVM, and KNN models, respectively. Nevertheless, precision increases by 4.31% and 1.49% for SVM and KNN methods while decreases by 0.7% in NB.

**Table 2:** Performance of NB, SVM, KNN, and RNN methods for the identification of “not boundary” class. The highest value for each performance metric is highlighted in bold.

Method	contextual features only			contextual + punctuation marks (+ topics except RNN)		
	Precision	Recall	F1-measure	Precision	Recall	F1-measure
NB	0.988	0.985	0.987	0.989	0.984	0.986
SVM	<b>0.989</b>	0.992	0.991	0.990	0.993	0.991
KNN	<b>0.989</b>	<b>0.995</b>	<b>0.992</b>	0.991	<b>0.994</b>	0.993
LSTM	–	–	–	0.993	<b>0.994</b>	<b>0.994</b>
GRU	–	–	–	<b>0.994</b>	<b>0.994</b>	<b>0.994</b>

Table 2 summarizes the performance of NB, SVM, KNN, and RNN models for detecting “not boundary” class in e-coaching text. We observed that performance is remarkably high in “not boundary” class compared to boundary detection which is expected because 96.71% instances are from “not boundary” class. KNN achieves 22.40%, 50.07%, and 36.26% higher precision, recall, and F1-measure for contextual features; and 20.85%, 33.96%, and 27.47% higher precision, recall, and F1-measure for combined features compared to boundary class. In contrast to boundary detection, RNN demonstrates the highest performance among all models. LSTM obtains 0.993 precision with 0.994 recall and F1-measure while LSTM exhibits 0.994 for all performance metrics. Impact of additional features is also consistent with “not boundary” classification. Results show that F1-measure increases by 0%, and 0.1% for SVM and KNN models although decreases by 0.1% for NB.

**Table 3:** Weighted average performance of NB, SVM, KNN, and RNN methods for the segmentation of e-coaching text in detecting both “boundary” and “not boundary” classes. The highest value for each performance metric is highlighted in bold.

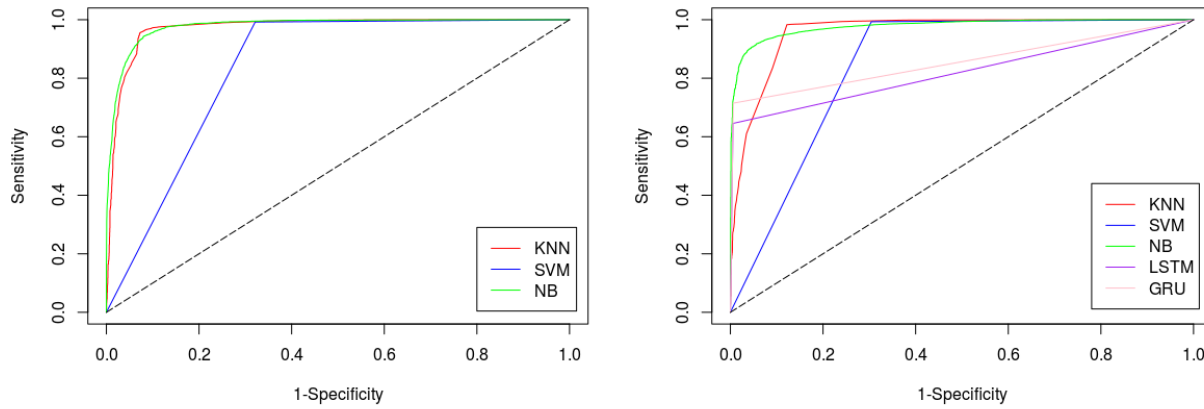
Method	contextual features only			contextual + punctuation marks (+ topics except RNN)		
	Precision	Recall	F1-measure	Precision	Recall	F1-measure
NB	0.975	0.974	0.975	0.976	0.974	0.975
SVM	0.981	0.982	0.981	0.983	0.983	0.983
KNN	<b>0.983</b>	<b>0.984</b>	0.983	<b>0.986</b>	<b>0.986</b>	<b>0.986</b>
LSTM	–	–	–	<b>0.986</b>	0.983	0.984
GRU	–	–	–	<b>0.986</b>	0.985	<b>0.986</b>

Table 3 outlines the weighted average results of the experiment on the models for the segmentation of e-coaching text by classifying them into “boundary” and “not boundary” classes. Overall, KNN obtains the best performance with all metrics and NB denotes the lowest performance among all methods. GRU demonstrates the same result as KNN for precision and F1-measure except for recall. SVM shows moderate performance, obtains 0.981, 0.982, and 0.983 for precision, recall, and F1-measure, respectively when textual features are used. Influence of the additional features is also consistent as above, precision increases by 0.1%, 0.2%, and 0.3%; recall increases by 0%, 0.1%, and 0.2%; and F1-measure increases by 0%, 0.2%, and 0.3% for NB, SVM, and KNN methods, respectively, when combined features are used.

## Discussion

This study is the first efforts to evaluate the automatic segmentation of e-coaching text. Experimental results indicate that KNN is the best model among all machine learning methods considered for this study. KNN achieved 0.986 F1-measure in overall, 0.779 and 0.993 F1-measures for detecting “boundary” and “not boundary”, respectively. The robust performance of KNN provides the evidence that machine learning models are capable to learn information from clinical exchange. Although the domain of this study was intentionally quite small, we believe that our study is not limited to the e-coaching domain, and it can be successfully applied to other domain as well.

The additional topic and punctuation mark feature made a significant improvement in performance of all machine learning methods. Nearly all cases, the model performs better when contextual features are used in combination with topic and punctuation mark features. This results also mean that segmentation performance might be improved by adding more relevant features including human insight into the problem.



**Figure 3: Receiver operating characteristic curves showing the performance of binary classifiers for the segmentation of e-coaching text when textual features (left) and combination of textual and other features (right) are used.**

In this paper, results are reported by each class to avoid confusion about the overall model performance. In addition, standard metrics: precision, recall, and F1-measure were used to eliminate doubt about the model performance because accuracy is misleading for imbalance dataset. AUC values are also outlined due to its effectiveness in measuring the quality of binary classifiers for imbalanced datasets<sup>34</sup>, which are demonstrated by the ROC curves in Figure 3. NB shows the highest AUC values, achieved 0.978 for both cases while provides lowest classification results. On the other hand, KNN and SVM exhibit 0.972 and 0.835 AUCs when only textual features are used; and 0.959 and 0.844 AUCs when a combination of textual, topic and punctuation mark features are used. Finally, LSTM demonstrates lowest AUC values among all machine learning models, achieved 0.82 AUC while GRU achieves 4.15% higher AUC than LSTM. The conclusion drawn from the ROC curves also confirmed the robustness and superiority of KNN model for the segmentation of clinical exchange.

We observed the second highest performance of RNN, in particular, LSTM and GRU for the text segmentation. We believe that RNN will perform better if a large set of data is utilized. In this study, we employed 3,138 examples of boundary case which failed to achieve best results. We also observed that GRU performs better than LSTM which was observed in other previous study<sup>35</sup>.

Although punctuation mark plays an important role in segmentation boundary detection, and large numbers of errors were encountered by the false positive of boundary identification. For example, a text block “A1 A2 A3. B1 B2 B3. C1 C2 C3 C4.” can be incorrectly segmented at position A3, B3, and C4 where a punctuation mark was encountered. Similarly, additional information is the common reason for the classified original segment into multiple segments. For

instance, the above text block can also be incorrectly segmented at position B3 and C4 because third sentence (C) only supports the MI code confirmed by first two sentences (A and B).

Our proposed approach is novel for the segmentation of e-coaching text because previous studies mainly focus on the segmentation of text into sections, headers, and sentences in other medical domain. However, this study segmented clinical exchange into groups of MI behaviors which will significantly reduce the amount of resource and time required to segment clinical exchange manually. Furthermore, these segmentation models could be integrated with auto coding classifiers to build a software package of automatic coding procedure of clinical exchange.

The limitation of this study is that e-coaching text is collected from a single medical institute; formatting, style, and email segment can be different in other settings. Therefore, there is a need to replicate the experiments with different data sets. As our future work, we plan to evaluate our approach on other datasets involves in discourse analysis. We also plan to use more relevant features to improve model performance. For example, part-of-speech tagging<sup>36</sup> and distance from the beginning and end of the current sentence might significantly enhance the classification performance.

## Conclusion

Segmentation of e-coaching text is an integral part of developing an automated e-coaching intervention. Although several studies have done in clinical interventions, they are limited by the qualitative coding of clinical interactions. In addition, previous studies in the medical domain mainly segmented clinical text into sections and sentences, none of them are used for the segmentation of text into groups of MI behaviors in the setting of discourse analysis with email under the principle of motivational interviews. In this paper, we compared the performance of machine learning models for the task of segmentation of e-coaching text. We found out that k-nearest neighbor provides the best performance for the segmentation of text in terms of all performance metrics. Manual segmentation of e-coaching data is very resource-intensive and time-consuming task, which can significantly decrease the time and effort required to develop an effective behavioral intervention. Our proposed methods can help to identify individual text segments, which can be annotated directly with a classification model. This approach will also help for developing fully automated e-coaching and accelerate the pace of identifying effective communication strategies.

## Acknowledgments

This study was supported by a grant from the National Institutes of Health, NIDDK R21DK108071, Carcone and Kotov, MPIs. We would like to thank research staff and student assistants in the Department of Family Medicine and Public Health Sciences at Wayne State University School of Medicine for their help in preparing the training dataset.

## References

- [1] Miller WR, Rollnick S. Motivational interviewing: Helping people change. Guilford press; 2012.
- [2] Miller WR, Rollnick S. Ten things that motivational interviewing is not. Behavioural and cognitive psychotherapy. 2009;37(2):129–140.
- [3] Miller WR, Rose GS. Toward a theory of motivational interviewing. American psychologist. 2009;64(6):527.
- [4] Apodaca TR, Longabaugh R. Mechanisms of change in motivational interviewing: a review and preliminary evaluation of the evidence. Addiction. 2009;104(5):705–715.
- [5] Nigam K, McCallum AK, Thrun S, Mitchell T. Text classification from labeled and unlabeled documents using EM. Machine learning. 2000;39(2-3):103–134.
- [6] Wang S, Manning CD. Baselines and bigrams: Simple, good sentiment and topic classification. In: Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Short Papers-Volume 2. Association for Computational Linguistics; 2012. p. 90–94.
- [7] de Vel OY, Corney MW, Anderson AM, Mohay GM. Language and gender author cohort analysis of e-mail for computer forensics. 2002;.

- [8] Lacson R, Barzilay R. Automatic processing of spoken dialogue in the home hemodialysis domain. In: AMIA Annual Symposium Proceedings. vol. 2005. American Medical Informatics Association; 2005. p. 420.
- [9] Freund Y, Schapire R, Abe N. A short introduction to boosting. *Journal-Japanese Society For Artificial Intelligence*. 1999;14(771-780):1612.
- [10] Hasan M, Kotov A, Carcone AI, Dong M, Naar S, Hartlieb KB. A study of the effectiveness of machine learning methods for classification of clinical interview fragments into a large number of categories. *Journal of biomedical informatics*. 2016;62:21–31.
- [11] Kotov A, Hasan M, Carcone A, Dong M, Naar-King S, BroganHartlieb K. Interpretable probabilistic latent variable models for automatic annotation of clinical text. In: AMIA Annual Symposium Proceedings. vol. 2015. American Medical Informatics Association; 2015. p. 785.
- [12] Siegfried J. Therapeutic and everyday discourse as behavior change: Towards a micro-analysis in psychotherapy process research. Greenwood Publishing Group; 1995.
- [13] Kalchbrenner N, Blunsom P. Recurrent convolutional neural networks for discourse compositionality. *arXiv preprint arXiv:13063584*. 2013;.
- [14] Pierre JM, Butler M, Portnoff J, Aguilar L. Neural Discourse Modeling of Conversations. *arXiv preprint arXiv:160704576*. 2016;.
- [15] Apostolova E, Channin DS, Demner-Fushman D, Furst J, Lytinen S, Raicu D. Automatic segmentation of clinical texts. In: Engineering in Medicine and Biology Society, 2009. EMBC 2009. Annual International Conference of the IEEE. IEEE; 2009. p. 5905–5908.
- [16] Denny JC, Spickard III A, Johnson KB, Peterson NB, Peterson JF, Miller RA. Evaluation of a method to identify and categorize section headers in clinical documents. *Journal of the American Medical Informatics Association*. 2009;16(6):806–815.
- [17] Tepper M, Capurro D, Xia F, Vanderwende L, Yetisgen-Yildiz M. Statistical Section Segmentation in Free-Text Clinical Records. In: LREC; 2012. p. 2001–2008.
- [18] Cho PS, Taira RK, Kangaroo H. Text boundary detection of medical reports. In: Proceedings of the AMIA Symposium. American Medical Informatics Association; 2002. p. 998.
- [19] Griffis D, Shivade C, Fosler-Lussier E, Lai AM. A quantitative and qualitative evaluation of sentence boundary detection for the clinical domain. *AMIA Summits on Translational Science Proceedings*. 2016;2016:88.
- [20] Kreuzthaler M, Schulz S. Detection of sentence boundaries and abbreviations in clinical narratives. In: BMC medical informatics and decision making. vol. 15. BioMed Central; 2015. p. S4.
- [21] Treviso MV, Shulby C, Aluísio SM. Sentence segmentation in narrative transcripts from neuropsychological tests using recurrent convolutional neural networks. *arXiv preprint arXiv:161000211*. 2016;.
- [22] Alexander GL, Lindberg N, Firemark AL, Rukstalis MR, McMullen C. Motivations of Young Adults for Improving Dietary Choices: Focus Group Findings Prior to the MENU GenY Dietary Change Trial. *Health Education & Behavior*. 2017;p. 1090198117736347.
- [23] Hashimoto K, Kontonatsios G, Miwa M, Ananiadou S. Topic detection using paragraph vectors to support active learning in systematic reviews. *Journal of biomedical informatics*. 2016;62:59–65.
- [24] Lu HM, Wei CP, Hsiao FY. Modeling healthcare data using multiple-channel latent Dirichlet allocation. *Journal of biomedical informatics*. 2016;60:210–223.



- [25] Ramage D, Hall D, Nallapati R, Manning CD. Labeled LDA: A supervised topic model for credit attribution in multi-labeled corpora. In: Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1-Volume 1. Association for Computational Linguistics; 2009. p. 248–256.
- [26] Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine learning in Python. *Journal of machine learning research*. 2011;12(Oct):2825–2830.
- [27] Chang CC, Lin CJ. LIBSVM: a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*. 2011;2(3):27.
- [28] Bengio Y, Frasconi P, Simard P. The problem of learning long-term dependencies in recurrent networks. In: *Neural Networks, 1993.*, IEEE International Conference on. IEEE; 1993. p. 1183–1188.
- [29] Hochreiter S, Schmidhuber J. Long short-term memory. *Neural computation*. 1997;9(8):1735–1780.
- [30] Cho K, Van Merriënboer B, Bahdanau D, Bengio Y. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:14091259*. 2014;.
- [31] Joachims T. Text categorization with support vector machines: Learning with many relevant features. In: *European conference on machine learning*. Springer; 1998. p. 137–142.
- [32] Aas K, Eikvil L. Text categorisation: A survey. Technical report, Norwegian computing center; 1999.
- [33] Kumar R, Indrayan A. Receiver operating characteristic (ROC) curve for medical researchers. *Indian pediatrics*. 2011;48(4):277–287.
- [34] Hu J, Yang H, King I, Lyu MR, So AMC. Kernelized Online Imbalanced Learning with Fixed Budgets. In: *AAAI*; 2015. p. 2666–2672.
- [35] Chung J, Gulcehre C, Cho K, Bengio Y. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:14123555*. 2014;.
- [36] Hasan M, Kotov A, Mohan A, Lu S, Stieg PM. Feedback or Research: Separating Pre-purchase from Post-purchase Consumer Reviews. In: *European Conference on Information Retrieval*. Springer; 2016. p. 682–688.