

Clinical Predictive Keyboard using Statistical and Neural Language Modeling

John Pavlopoulos

Department of Computer and System Sciences
Stockholm University, Sweden
ioannis@dsv.su.se

Panagiotis Papapetrou

Department of Computer and System Sciences
Stockholm University, Sweden
panagiotis@dsv.su.se

Abstract—A language model can be used to predict the next word during authoring, to correct spelling or to accelerate writing (e.g., in sms or emails). Language models, however, have only been applied in a very small scale to assist physicians during authoring (e.g., discharge summaries or radiology reports). But along with the assistance to the physician, computer-based systems which expedite the patient’s exit also assist in decreasing the hospital infections. We employed statistical and neural language modeling to predict the next word of a clinical text and assess all the models in terms of accuracy and keystroke discount in two datasets with radiology reports. We show that a neural language model can achieve as high as 51.3% accuracy in radiology reports (one out of two words predicted correctly). We also show that even when the models are employed only for frequent words, the physician can save valuable time.

Index Terms—language modeling, text prediction, predictive keyboard, clinical text, data entry

I. INTRODUCTION

Text prediction is a challenging problem in machine learning and natural language processing, while at the same time there is a growing need for novel techniques for efficient and accurate text prediction in several application domains, such as in dictation and typing systems for people with disabilities or clinical text prediction for healthcare practitioners [1]. More concretely, with text prediction we refer to the task of predicting the next block of text in an online fashion, where block can refer to different text granularity levels, e.g., sentences, words, syllables, or characters (keystrokes) [2].

The main focus of this paper is medical text with the concrete task of predicting the next word given an incomplete text. We also refer to this problem as *predictive keyboard* for medical text. When applied in the clinical setting (e.g., authoring of hospital discharge summaries or diagnostic text), physicians can vastly benefit from a fast and accurate predictive keyboard system, since it can assist them with (a) a speedy compilation of the intended text, (b) means for prevention of potential text errors due to work overload, (c) means for speedier patient discharge.

Initial efforts towards solving the predictive keyboard problem for radiology reports are described by Eng and Eisner 2004 [3], where a 3-Gram language model achieves an average keystroke reduction of a factor of 3.3. Following this line of research, we employed N-Gram-based statistical language modeling, which refer to as N-GLM, to predict the next word

TABLE I
EXAMPLE USE-CASE ON IUXRAY TEST WORDS, USING 4-GLM AND LSTMLM. WORDS IN [] WERE CORRECTLY PREDICTED BY EACH MODEL.

LSTMLM	”the lungs are clear without [evidence] [of] focal infiltrate [or] [effusion] [there] [is] [no] [pneumothorax] [the] [visualized] [bony] [structures] [reveal] [no] [acute] [abnormalities]”
4-GLM	”the lungs are [clear] without evidence [of] [focal] infiltrate or effusion [there] [is] [no] pneumothorax [the] visualized bony [structures] reveal [no] [acute] [abnormalities]”

of a clinical text. We vary N from 1 to 10 and show that 4-Gram models achieve 38% accuracy when predicting the next word in a clinical text, outperforming other N-GLMs. Observe that accuracy in this case measures the fraction of times when the next word was predicted correctly, hence inducing an equivalent typing speedup at the word level. We additionally investigated two neural language models that employ (1) a Recurrent Neural Network (RNN) language model based on Long-Short Term Memory, which we refer to as LSTMLM [4] and (2) a Gated Recurrent Unit (GRU) based language model, which we refer to as GRULM. This model achieves higher levels of accuracy compared to 4-GLM, since our experimental evaluation demonstrates that accuracy can reach up to 51.3% (i.e., 5 out of 10 ‘next’ words predicted correctly). An example of the output of this task is depicted in Table I, where we can observe the next word predictions made by LSTMLM and 4-GLM, with the correctly predicted words indicated in ‘[]’.

Next, we outline the related work in the area of clinical text prediction, followed by a summary of our contributions.

A. Related Work

The study of the benefits of computer-assisted text generation dates back to more than two decades ago [5]. When applied to clinical notes, such as radiology reports, a statistical 3-Gram language model (including back off) achieved substantial keystroke reductions [3]. Recently, an even simpler 3-Gram language model (i.e., with no back off) outperformed the earlier while also decreasing the typing time for the clinician by one third [6]. These results demonstrate that N-Gram models can provide promising solutions to our problem, and hence in this paper we provide a more extensive evaluation of these models on medical text. Besides computer-assisted typing, language models have also been used for spelling

correction in clinical notes [7], [8]. This work does not focus on spelling correction, but what these works verify is that the words suggested by the language model during typing, are also checked for their correctness (i.e., assuming that the corpus contains correct words), hence the generated text will be of equal or even higher quality.

With the recent advance of deep learning, deep neural networks, such as Long-short Term Memory (LSTM) [4] models, have improved the performance in natural language processing (NLP) tasks of the biomedical field, such as Name Entity Recognition (NER) [9]; medical codes prediction [10]; relation classification [11]; predicting hospital readmission [12]. And language modeling is also part of this advent, since it is often employed as a pre-training step [13]. For the task of next word prediction in a medical setting, however, neural language modeling is heavily under explored. To our knowledge, the only application of a neural language model was that of a baseline LSTM network (applied on a private dataset), which was improved when structured information from electronic health data (e.g., gender or age) was integrated [14]. The authors reported 8% Accuracy (a.k.a. Recall@1 or Precision@1) for the baseline LSTM, which ranks it much lower than competing statistical language models [11]. However, neural networks have been reported to outperform statistical language modeling in non-medical domains [15]. In this work, we compare statistical and neural language modeling, a comparison which has not been studied before, and we show that the neural approach outperforms the statistical approach in next word prediction by a large margin.

B. Contributions

The main contributions of this paper can be summarized as follows: (1) We highlight the importance of the problem of keyword prediction for clinical text, and demonstrate how language models can be employed for providing scalable solutions to this problem; (2) We provide an extensive benchmark on clinical text obtained from two real-world medical datasets by comparing the performance of the N-GLM model for different values of N in terms of accuracy and keystroke reduction; (3) We additionally compare an RNN language model based on LSTM and GRU on the same datasets and demonstrate their superiority against N-GLMs as they can achieve an accuracy of up to 51.3%, indicating a speedup (at the word level) of the same degree, and a keyword reduction of up to 41.12%, indicating a speedup (at the character level) of the same degree.

II. METHODS

A. Statistical Language Modeling

Statistical language models [16], [17] are based on the Markov assumption, modeling the probability of the next word, but given only the $n-1$ preceding words. The counts of all sequences of n words (a.k.a. n -grams) are calculated over a corpus and a probability distribution over the vocabulary is modeled for each gram of $n-1$ words:

$$P(w_i|w_{i-1}, \dots, w_1) = P(w_i|w_{i-1}, \dots, w_{i-n+1}) \quad (1)$$

Then, a 2-gram (a.k.a. bigram) model will only consider the previous word w_{i-1} to predict a next word w_i . And w_i will be the one most frequently occurring in the corpus right after w_{i-1} . Probabilities are formed using the maximum likelihood estimation, changing Eq. 1 to:

$$P(w_i|w_{i-1}) = \frac{C(w_{i-1}w_i)}{C_{i-1}} \quad (2)$$

where C are the counts of the gram. To deal with unknown words, a pseudo token can be introduced (e.g., masking very rare words with '[OOV]' during training). And to deal with unseen sequences of words one can introduce smoothing or backoff and interpolation. In this work we employed Laplace smoothing, but we observe that algorithms such as the Knesser-Ney or the Good-Turing backoff should also be investigated. For more information regarding statistical language models we redirect the interested reader to [17].

B. Neural Language Modeling

Neural language modeling makes it possible to consider long range word dependencies without an explicitly predefined context length [18]. The neural language model at each time step s learns a hidden state h_s as the non-linear combination (the weight matrix W is learned) of the input word x_s and the previous hidden state h_{s-1} . The vanishing gradient problem, arising from the deep in time back-propagation, is addressed with the LSTM cells [4]. More formally:

$$\begin{aligned} i_s &= \sigma(W_i \cdot [x_s, h_{s-1}] + b_i) \\ f_s &= \sigma(W_f \cdot [x_s, h_{s-1}] + b_f) \\ o_s &= \sigma(W_o \cdot [x_s, h_{s-1}] + b_o) \\ q_s &= \tanh(W_q \cdot [x_s, h_{s-1}] + b_q) \\ c_s &= f_s \cdot c_{s-1} + i_s \cdot q_s \\ h_s &= o_s \cdot \tanh(c_s), \end{aligned} \quad (3)$$

where i_s is the input gate and f_s is the forget gate, which regulate the information from this (q_s) and the previous (c_{s-1}) cell to be forgotten, and o_s is the output gate which regulates the information of the new hidden state. Then, the generation of the next word x_{s+1} can be seen as a classification task, with *softmax* yielding a probability distribution over the whole vocabulary and the next word to be generated being the most probable one.

In this work we also experiment with a different RNN variant, called Gated Recurrent Unit (GRU) [19], which is considered to be more efficient than LSTMs [20]. It has a similar formulation with LSTMs:

$$\begin{aligned} r_s &= \sigma(W_r \cdot [x_s, h_{s-1}] + b_r) \\ u_s &= \sigma(W_u \cdot [x_s, h_{s-1}] + b_u) \\ c_s &= \tanh(W_c \cdot [x_s, r_s \cdot h_{s-1}] + b_c) \\ h_s &= u_s \cdot c_{s-1} + (1 - u_s) \cdot c_s \end{aligned} \quad (4)$$

where r_s and u_s are the reset and update gates, defined similarly to the input and forget LSTM gates. No output gate is used, leading to a smaller number of gates and less computations, which makes GRU more efficient than LSTM.

TABLE II
ASSESSMENT OF NEXT WORD PREDICTION IN THE RADIOLOGY REPORTS OF IUXRAY AND MIMIC-III, USING STATISTICAL (N-GLMS) AND NEURAL (LSTMLM, GRULM) LANGUAGE MODELS. MICRO-AVERAGED ACCURACY (ACC) AND KEYSTROKE DISCOUNT (KD) ARE SHOWN FOR EACH DATASET.

	IUXRAY		MIMIC-III	
	ACC	KD	ACC	KD
2-GLM	21.830.29	16.040.26	17.030.22	11.460.12
3-GLM	34.780.38	27.960.27	27.340.29	19.350.27
4-GLM	38.180.44	31.600.30	25.700.29	18.950.34
5-GLM	37.890.60	32.300.47	21.020.41	15.630.23
6-GLM	35.710.78	30.860.57	15.980.42	11.930.31
7-GLM	33.100.72	28.820.56	12.150.40	9.050.26
8-GLM	30.230.63	26.470.62	9.520.40	7.040.31
9-GLM	27.740.63	24.330.66	7.290.43	5.460.37
LSTMLM	51.300.61	41.120.64	33.970.25	25.170.29
GRULM	51.300.74	41.000.40	33.840.34	25.420.30

III. EMPIRICAL EVALUATION

A. Datasets

We used two real-world medical datasets.

IUXRAY. The dataset comprises 3,955 anonymized and de-identified radiology reports on 7,470 images [21]¹. The text of each report follows an XML structure and the boundaries of each different section are explicitly defined.

MIMIC-III. We used the radiology reports from the Medical Information Mart for Intensive Care (MIMIC-III) [22] database, a rich and commonly used benchmark dataset of 38,597 adult patients admitted between 2001 and 2008 to critical care units at Beth Israel Deaconess Medical Center in Boston, Massachusetts. In this study we employ the free text reports of electrocardiogram and imaging studies included in this dataset. The text of the radiology reports in MIMIC-III is loosely separated in sections, which are not explicitly marked up. We sampled 2,928 such reports to yield a dataset equal in number to IUXRAY.

The radiology reports of IUXRAY and MIMIC-III comprise less than 200 tokens per report in average. By contrast, the discharge summaries are lengthier and more than quadruple in size. The difference grows larger when sampling disregards the maximum number of characters per text, because only discharge summaries did exceed this threshold.²

B. Results

We benchmarked eight statistical language models and two neural language models for the task of predicting the next word in radiology reports of IUXRAY and MIMIC-III. We randomly sampled reports from MIMIC-III until we obtained a subset with the same number of reports as IUXRAY. We additionally removed numbers, punctuation, and turned to lower-case before white space tokenization. We held the 10K last words from each dataset as our test set and used the previous to train our models. Any words occurring less than 10 times were masked with an OOV token.

¹<https://openi.nlm.nih.gov/>

²The average number of words per summary without sampling is 1320.

The statistical language models were N-Gram-based models, with factor N varying from 2 (only the previous word was considered) to 9. The neural language models were based either on LSTM or GRU. Following the work of [14], we used 50 dimensions for all the hidden representations. Furthermore, we used: a vocabulary of the 1000 most frequent words; a context window of 5 preceding words; uniformly initialized word embeddings of 200 dimensions; a single-layer feed-forward neural network of 100 dimensions and a RELU activation before the softmax; Adam optimization and categorical cross entropy; batch size 128; 10% validation split; early stopping of 100 epochs with patience of 3 epochs; validation loss monitoring.

First, we assessed all models based on their ability to reduce the keystrokes (Keystroke Discounting, KD). Since no log files were available for calculating this number directly, we estimated this score based on the length of the words which were correctly predicted by each system. That is, we assume that instead of striking the keyboard as many times as the characters of a word, the physician, during a computer-assisted data entry, simply accepts the correctly predicted word (e.g., by pressing TAB or so). More formally, for a sequence of N words ($w_1^g \dots w_N^g$) and the respective sequence of system-predicted words ($w_1^p \dots w_N^p$), this measure is defined as:

$$KD = 1 - \frac{\sum_{i=1 \dots N} dsc(i)}{\sum_{i=1 \dots N} |w_i^g|}, \quad (5)$$

$$dsc(i) = \begin{cases} 1, & \text{if } w_i^g = w_i^p \\ |w_i^g|, & \text{otherwise} \end{cases}$$

where $|w_i|$ is the number of characters of word w_i and dsc is 1 when the word was correctly predicted by a system and equal to the length of the correct word otherwise. When KD equals to 1, all words are predicted correctly, while when it is equal to 0, no word was predicted correctly. Also, we used micro-averaged accuracy (here, same as precision or recall), which is defined as *the fraction of the correctly predicted words out of all the words in the test*. For both measures the occurrences of the de-identification token ('XXXX') were disregarded during evaluation, because the ability of the systems to locate candidate de-identification terms is out of the scope of this work. However, in principle, medical language models could be used to assist humans in de-identifying medical texts. And we considered all OOV occurrences as system mistakes.

Table II shows the keyword discount (KD) and the micro-averaged accuracy (ACC) scores for the task of next word prediction, for all systems and datasets. Neural language models outperform statistical language models in both datasets by a large margin. In MIMIC-III, the keyword discount is increased by 6 absolute percentage units, from 19.35% to 25.42% (or 31% relative increase). A similar increase was found for accuracy, from 27.34% to 33.97% (or 24% relative increase). In IUXRAY, the increase was larger, with 9 absolute percentage units of KD (from 32.30% to 41.12%) and 13 absolute percentage units of ACC (from 38.18% to 51.30%). The top eight rows of Table II show the different N-Gram

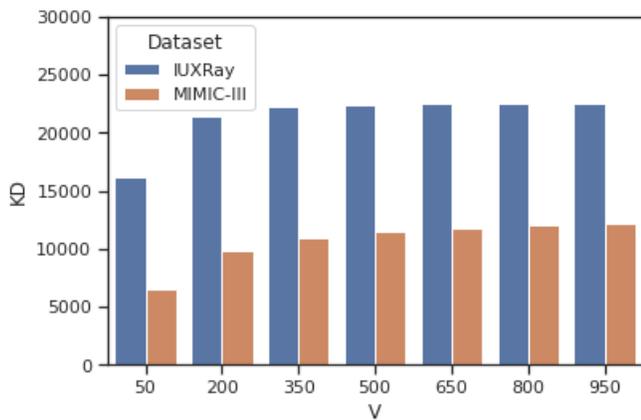


Fig. 1. Absolute number of keystroke reduction, by applying LSTMML only when a frequent vocabulary word is predicted. On the x-axis we see the sizes of the frequent-word sets that are employed.

based statistical language models. 3-GLM was the best with both evaluation measures in MIMIC-III. In IUXRAY, 4-GLM was found as the best in terms of ACC and 5-GLM in KD. We obtained better performance for IUXRAY due to its smaller vocabulary size compared to MIMIC-III.

In a real-world setting, however, physicians may prefer to use the advantages of computer-assisted authoring only for specific terms, as for example frequent words, frequent medical terms, or frequent non-medical terms. Thus, in a final experiment, we assumed a deployment setting where the predicted word was only shown if it was one of the frequent ones, and we varied the number of frequent words to be considered. Fig. 1 shows the absolute number of keystrokes omitted when the best performing LSTMML was applied. Interestingly, even though the target vocabulary is reduced to only 50 words, we can observe a decrease of more than 15K keystrokes. For the case of the 50 most frequent words, without the use of LSTMML the keystrokes would have been approximately 50K.

IV. CONCLUSION

We highlighted the importance of predictive keyboard for medical text and demonstrated the benefits for the physicians in terms of speedups in completing their clinical text reports. Our experimental evaluation on radiology reports from two real-world medical datasets showed that neural language models can achieve an accuracy of up to 51.3%, which implies that the obtained speedups correspond to a similar factor at the word level. Directions for future work include the investigation of alternative statistical and deep learning models, the consideration of additional medical datasets (e.g., discharge summaries), and measuring the speedups in a real-world application with models deployed in healthcare systems.

ACKNOWLEDGMENTS

We thank the anonymous reviewers for their comments. This work was supported in part by the Swedish Research Council starting grant Temporal Data Mining for Detective Adverse

Events in Healthcare, ref. no. VR-2016-03372 as well as the EXTREME project of the Digital Futures framework.

REFERENCES

- [1] N. Garay-Vitoria and J. Abascal, "Text prediction systems: A survey," *Univers. Access Inf. Soc.*, vol. 4, no. 3, p. 188203, Feb. 2006. [Online]. Available: <https://doi.org/10.1007/s10209-005-0005-9>
- [2] J. Gelšvartas, R. Simutis, and R. Maskeliūnas, "User adaptive text predictor for mentally disabled huntingtons patients," *Intell. Neuroscience*, vol. 2016, Jan. 2016. [Online]. Available: <https://doi.org/10.1155/2016/3054258>
- [3] J. Eng and J. M. Eisner, "Radiology report entry with automatic phrase completion driven by language modeling," *Radiographics*, vol. 24, no. 5, pp. 1493–1501, 2004.
- [4] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [5] H. H. Koester and S. P. Levine, "Learning and performance of able-bodied individuals using scanning systems with and without word prediction," *Assistive Technology*, vol. 6, no. 1, pp. 42–53, 1994.
- [6] A. Yazdani, R. Safdari, A. Golkar, and S. R. N. Kalhori, "Words prediction based on n-gram model for free-text entry in electronic health records," *Health information science and systems*, vol. 7, no. 1, p. 6, 2019.
- [7] J. Patrick, M. Sabbagh, S. Jain, and H. Zheng, "Spelling correction in clinical notes with emphasis on first suggestion accuracy," in *2nd Workshop on Building and Evaluating Resources for Biomedical Text Mining*, 2010, pp. 1–8.
- [8] A. Yazdani, M. Ghazisaeedi, N. Ahmadinejad, M. Giti, H. Amjadi, and A. Nahvijou, "Automated misspelling detection and correction in persian clinical text," *Journal of Digital Imaging*, pp. 1–8, 2019.
- [9] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang, "BioBERT: a pre-trained biomedical language representation model for biomedical text mining," *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, 2020.
- [10] J. Mullenbach, S. Wiegrefe, J. Duke, J. Sun, and J. Eisenstein, "Explainable prediction of medical codes from clinical text," *arXiv preprint arXiv:1802.05695*, 2018.
- [11] Y. Luo, "Recurrent neural networks for classifying relations in clinical notes," *Journal of biomedical informatics*, vol. 72, pp. 85–95, 2017.
- [12] K. Huang, J. Altosaar, and R. Ranganath, "ClinicalBERT: Modeling clinical notes and predicting hospital readmission," *arXiv preprint arXiv:1904.05342*, 2019.
- [13] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, "Language models are unsupervised multitask learners," *OpenAI Blog*, vol. 1, no. 8, p. 9, 2019.
- [14] G. P. Spithourakis, S. E. Petersen, and S. Riedel, "Clinical text prediction with numerically grounded conditional language models," *arXiv preprint arXiv:1610.06370*, 2016.
- [15] Y. N. Dauphin, A. Fan, M. Auli, and D. Grangier, "Language modeling with gated convolutional networks," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 933–941.
- [16] F. Jelinek, *Statistical methods for speech recognition*. MIT press, 1997.
- [17] D. Jurafsky, *Speech & language processing*. Pearson Education India, 2000.
- [18] M. Sundermeyer, H. Ney, and R. Schlüter, "From feedforward to recurrent lstm neural networks for language modeling," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 23, no. 3, pp. 517–529, 2015.
- [19] K. Cho, B. Van Merriënboer, D. Bahdanau, and Y. Bengio, "On the properties of neural machine translation: Encoder-decoder approaches," *arXiv preprint arXiv:1409.1259*, 2014.
- [20] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *arXiv preprint arXiv:1412.3555*, 2014.
- [21] D. Demner-Fushman, M. D. Kohli, M. B. Rosenman, S. E. Shooshan, L. Rodríguez, S. Antani, G. R. Thoma, and C. J. McDonald, "Preparing a collection of radiology examinations for distribution and retrieval," *Journal of the American Medical Informatics Association*, vol. 23, no. 2, pp. 304–310, 2015.
- [22] A. Johnson, T. J. Pollard, L. Shen, L. Li-wei, M. Feng, M. Ghassemi, B. Moody, P. Szolovits, L. A. Celi, and R. G. Mark, "Mimic-iii, a freely accessible critical care database," *Scientific data*, vol. 3, p. 160035, 2016.