

An Empirical Evaluation of Deep Learning for ICD-9 Code Assignment using MIMIC-III Clinical Notes

Jinmiao Huang^{1*}, Cesar Osorio^{1*}, Luke Wicent Sy^{1*}

¹ Georgia Institute of Technology, Atlanta, Georgia, USA

Abstract

Code assignment is important on many levels in the modern hospital, from ensuring accurate billing process to creating a valid record of patient care history. However, the coding process is tedious, subjective, and requires medical coders with extensive training. The objective of this study is to evaluate the performance of deep learning based systems to automatically map clinical notes to medical codes. We applied the state-of-the-art deep learning methods such as Recurrent Neural Networks and Convolution Neural Networks on MIMIC-III dataset. Experiments show that the deep-learning-based methods outperform other conventional machine learning methods. Our evaluations are focused on end-to-end learning methods without manually defined rules. From our evaluations, the best models are able to predict the top 10 ICD-9 codes with 69.57% F1 and 89.67% accuracy; the top 10 ICD-9 categories with 72.33% F1 and 85.88% accuracy. The evaluation tools and resources are available at <https://github.com/lsy3/clinical-notes-diagnosis-dl-nlp>.

1 Introduction

Electronic health record (EHR) data capture variety of patient clinical information such as medical history, vital signs, lab test results, clinical notes, etc. It builds a continuous flow of information between the doctor and the patient. Systematic reviews have shown the clinical care quality improvement using predictive analysis based on EHR data¹.

EHR data contains both structured (e.g., blood pressure) and unstructured data (e.g., doctor’s observation). While many medical systems focus on structured biosignal features in EHR to build the clinical decision making systems^{2,3}, more than 80% of health record data is unstructured text⁴. For example, clinical notes contain information about the patient’s medical history, doctor’s observations and comments about their interactions with patients.

The systems evaluated in this paper assign ICD-9 codes from a patients’ free-text Electronic health record (EHR). These codes can be subsequently used in billing or creating a valid record of patient care history. Currently, the task of assigning the diagnoses codes is carried out manually by the medical coders. The volume of medical records generated nowadays makes the manual classification of diagnoses a labor-intensive process, which results in a significant backlog. Automating ICD-9 code assignment will not only make the clinical process more efficient, but it can also take note of all EHR and provide foundation to automate some level of semantic analysis which can help clinicians perform better diagnosis and effectively improve the medical care systems.

Recently, deep learning approaches have shown a significant improvement in many Natural Language Processing (NLP) tasks such as language translation, image caption, and sentiment analysis. Deep learning models can often be trained end-to-end without any domain-specific and often tedious hand-designed feature engineering.

Therefore, our work is focused on evaluating the performance of the state-of-the-art deep neural network to the diagnose learning system. We also compared our results with several traditional classification systems, such as logistic regression, and random forests, each with the goal of predicting the code from the clinical notes. A extensive number of experiments are applied to different settings of the tested classification algorithm. We also use word embeddings to transform a patient’s free-text EHR to information that can be used to predict ICD-9 codes, and evaluated the impact of word embedding trained from MIMIC-III dataset and the medical domain word embedding. We hope our work can provide a benchmark for the learning-based ICD-9 code assignment on MIMIC-III dataset.

2 Related Work

The task of automatic ICD-9 coding has been attempted for decades. Larkey and Croft⁵ designed classifiers for the automatic assignment of ICD-9 codes to discharge summaries in 1995. Automated ICD-9 coding for radiology reports

*These authors contributed equally to this work

was one of the first challenges in informatics community⁶ in 2007. There are two major categories of approaches for automatically assigning ICD-9 codes using text-free clinical notes. One is rule-based, the other one is learning-based. Rule-based systems are designed by human experts. This type of methods has out-performed other methods in many cases^{6,7}. However, this kind of system heavily relies on the manual intervention of the medical professionals, thus hard to maintain and scale up to more general cases. Learning-based systems do not require any domain knowledge from the medical experts, which only rely on learning algorithms to find the underline distribution of the datasets⁸⁻¹⁰. A detailed review of extracting information from textual documents in the EHR can be found in¹¹ and¹².

End-to-End data-driven approaches have gained popularity in the last few years. Recent methods based on deep learning have demonstrated the state-of-the-art performance in a wide variety of tasks, including computer vision¹³, speech recognition¹⁴, and NLP¹⁵. In the clinical domain, Choi et. al.¹⁶ use Recurrent Neural Network (RNN) to predict heart failure. Lipton et. al.³ use LSTM to classify 128 diagnoses from 13 frequently but irregularly sampled clinical measurements extracted from the structured EHR data. Similarly, *DoctorAI*² and *RETAIN*¹⁷ also use RNN on structured EHR data for diagnosis classification. Researchers also use deep learning on unstructured free-text to predict the diagnosis. Luo proposed LSTM¹⁸ for classifying relations from clinical notes on the i2b2/VA relation classification challenge dataset. Prakash et. al.¹⁹ exploit raw text from Wikipedia as a knowledge source, and introduced condensed memory neural networks to learn the diagnose on the MIMIC-III dataset. Since Prakash et.al.¹⁹ tackled a similar problem as we do, we compared our results to their's in Section 4.2.4.

3 Methodology

Figure 1 shows an overview of our methodology pipeline. Our methodology involves the following steps: data preprocessing, feature extraction, and model training and testing. Specifically, we use spark was used for data preprocessing; Spark, sklearn, and gensim for feature extraction; and Spark ML, Keras for model training and testing. We use Azure virtual machines (NC24 with K80 GPU) to run our experiments. Section 3.1 to 3.3 describes each step in more detail. Each model is also evaluated under a set of metrics, as described in section 3.4.

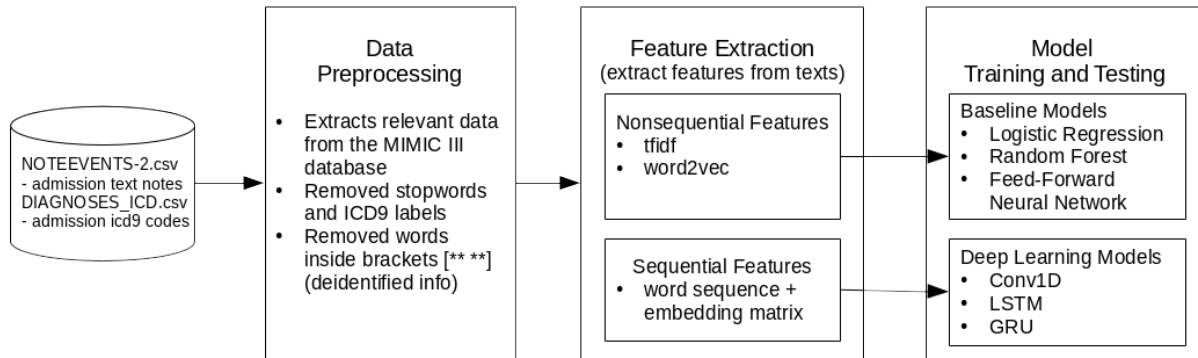


Figure 1: Methodology Pipeline Overview

3.1 Data Preprocessing

MIMIC-III dataset is a large data set relating to patients admitted to critical care units at a large tertiary care hospital. It contains de-identified medical records of patients who stayed within the intensive care units at Beth Israel Deaconess Medical Center from 2001 to 2012²⁰. The goal of this study is to explore useful semantic information using unstructured data, therefore only the free-text clinic note section from the dataset, specifically the *noteevents* table, was used. Furthermore, we focused on the *discharge summaries* category as it contains actual ground truth and free-text compared to other categories. Since *discharge summaries* were written after the diagnosis was made, the notes are sanitized by removing any mention of class-labels (ICD-9 codes). This approach is similar to Prakash et al.¹⁹.

Table 1 describes the number of unique patients, hospital admissions, ICD-9 codes and ICD-9 categories involve in the MIMIC-III dataset. *All MIMIC-III* describes the whole dataset, while *noteevents* and *discharge summaries* describe

the corresponding subsets.

Coverage	Patients	Hospital Admissions	ICD-9 Codes	ICD-9 Categories
<i>All MIMIC-III</i>	46520	58976	6984	943
<i>noteevents</i>	46146	58361	6967	943
<i>discharge summaries</i>	41127	52726	6918	942

Table 1: MIMIC-III Descriptive Statistics

The data was preprocessed to produce separate datasets through two approaches. The first approach is to treat the ICD-9 code independently from each other, find the admissions (unique HADM.ID) for each ICD-9 classification, and consider only records related to the top 10 and top 50 common ICD-9 codes. Top 10 and top 50 are chosen because they cover a majority of the dataset (76.9% and 93.6% as can be observed in table 3, and for ease of comparison with the result of²¹. The second approach is to group ICD-9 codes into categories based on its hierarchical nature, with categories for larger sets of similar health conditions (e.g. "Cholera due to vibrio cholerae" has the ICD-9 code 001.0, and is categorized as a type of Cholera, which in turn is a type of Intestinal Infectious Disease), and then find the patients for top 10 and top 50 common categories. Evaluation will be separately performed on the four datasets. The four datasets will hereby be referred as *top-10-code*, *top-50-code*, *top-10-category* and *top-50-category* respectively.

Table 2 shows the top 10 ICD-9 codes and top 10 ICD-9 categories. Table 3 also describes the number of unique hospital admissions related to the four datasets mentioned in the earlier paragraph.

ICD-9 Code	Admissions	ICD-9 Category	Admissions
4019: Hypertension	20046	401: Essential hypertension	20646
4280: Congestive heart failure	12842	427: Cardiac dysrhythmias	16774
42731: Atrial fibrillation	12589	276: Disorders of fluid electrolyte	14712
41401: Coronary atherosclerosis	12178	272: Disorders of lipid metabolism	14212
5849: Acute kidney failure	8906	414: Other chronic ischemic heart disease	14081
25000: Diabetes Type II	8783	250: Diabetes mellitus	13818
2724: Hyperlipidemia	8503	428: Heart failure	13330
51881: Acute respiratory failure	7249	518: Other diseases of lung	12997
5990: Urinary tract infection	6442	285: Other and unspecified anemias	12404
53081: Esophageal reflux	6154	584: Acute kidney failure	11147

Table 2: Top 10 ICD-9

Data Set	Hospital Admissions	<i>discharge summaries</i> Coverage (%)
<i>top-10-code</i>	40562	76.93%
<i>top-50-code</i>	49354	93.60%
<i>top-10-category</i>	44419	84.24%
<i>top-50-category</i>	51034	96.79%

Table 3: Dataset Descriptive Statistics

The filtered datasets will be split to 50-25-25 for training, validation and testing.

3.2 Feature extraction

We use two approaches for feature extraction: term frequency - inverse document frequency (TF-IDF) and word2vec^{22,23}. The TF-IDF is served as a baseline to compare with word2vec.

TF-IDF is intended to evaluate how important a word is to a document in a collection of documents or corpus. It is the

product of two statistics: TF and IDF. TF is the number of times a word appears in a given document and IDF measures whether a word is common or rare across the corpus. We use the following definition of IDF for our calculations:

$$IDF(w) = \log \frac{N_d}{DF(d, w)} + 1$$

where N_d is the total number of documents and $DF(d, w)$ is the number of documents that contain word w .

To calculate TF-IDF, first we tokenized all the notes in the filtered training data set, then create a document-word matrix with the count of each word in each note (TF) and finally multiply each word by the corresponding IDF. We used two TF-IDF configurations: (1) one with top 40,000 words with highest TF-IDF scores as the bag of word features and (2) a second one with minimum document frequency of 10 and maximum document frequency of 0.8, which reduces our total number of words to around 20,000 words.

word2vec takes a tokenized text corpus as an input and produces word vectors as output. We used the Continuous Bag of Words (CBOW) architecture which predicts the target word based on the context: words that precede and follow the target word. CBOW is basically a Feed-Forward Neural Network model that consists of inputs, projection and output layers where the traditional non-linear hidden layer is removed to reduce time complexity and the projection layer is shared by all the words. The inputs are words in the context. We use text notes from MIMIC-III as corpus to train our word2vec model. We also use pre-trained word vectors induced from PubMed found on <https://github.com/cambridgeltl/BioNLP-2016>²⁴

3.3 Model Training and Testing

One fundamental assumption adopted by traditional supervised learning algorithms is that each sample has only one label assigned to it. In our problem, each sample has multiple (one or more) ICD-9 codes attached to it. Generally, there are two main methods for tackling the multi-label classification problem²⁵ (1) problem transformation methods and (2) algorithm adaptation methods. Problem transformation methods transform the multi-label problem into a set of binary classification or regression problems, multiple binary classifiers are trained separately for each label. Algorithm adaptation methods adapt the algorithms to perform multi-label classification in its full form and only one classifier are trained for all the labels.

In our study, we first create three baseline approaches: Linear Regression, Random Forests and Feed-Forward Neural Network. Among which, we used problem transformation methods to get the multi-label output for Linear Regression and Random Forest classifiers. Specifically, in order to assign each sample a set of target labels, we simply trained n different models for n different labels, each model independently predicts a mutual exclusive output (0 or 1) for each sample data. For Feed-Forward Neural Network, we used algorithm adaptation based methods, since neural network can be easily adapted to multi-label problem by setting up multiple neurons in the network output layer and set each neuron to a target label correspondingly. Similar to Feed-Forward Neural Network, we also used algorithm adaptation-based methods to our deep learning models. In this section, we will describe our implemented models in detail.

3.3.1 Baseline Models

Logistic Regression: Our first baseline model is a binomial logistic regression model implemented using Spark ML. For each label (ICD-9 code or category), a separate logistic regression model was trained and each model independently predicts the said label (0 or 1 for the corresponding ICD-9 code or category). We tried different configurations; specifically "no. of iterations" was tuned between 5 to 100. Since we only use notes under *discharge summaries* category, there is 1 note per admission. Features are extracted from this note and are used as inputs for this classifier. For *tfidf*, the features are directly used as input features. For *word2vec*, the input features are the average of all the feature vectors of the words in the notes.

Random Forest: Our second baseline model is a random forest model implemented using Spark ML. The same approach and input for the logistic regression were used here (one model for each label). Different configurations

were also tried, specifically "tree depth" was tuned between 5 to 30.

Feed-Forward Neural Network: One advantage of Neural Network is that it can be fitted to multi-label problem in just one model with the proper activation function. We implement the Feed-Forward Neural Network as the baseline for algorithm adaptation based (described in 3.3) multi-label classification problem. We use the same input features and train-test data split as previously described. We use ReLU activation function for all the hidden layers and sigmoid function for the output layer, binary cross entropy as the loss function, and stochastic gradient descent as the optimizer. We tried several neural network models with one to four different hidden layers. For each hidden layer, a total of seven models were tried with the combination of neuron size 50, 100, 300, 500 and 1000. The results for different configurations can be found in our experiments spreadsheet on our code repository.

3.3.2 Deep Neural Network Models

In this study, we cast the ICD-9 code assignment from clinical notes as multi-label classification problem on sequential observations x_1, x_2, \dots, x_n , where x_i is the word2vec features we calculated for word i in the discharge summary. Unlike the features we used for the baseline models, where the sequential information is not preserved, we sequentially takes each word from the discharge summary. The input features for this classifier are N most recent word sequences taken from the notes. If we don't have enough feature events, we pad zero vectors at the beginning. The word sequence is then converted into vectors using an embedding matrix based on a *word2vec* model (See Section 3.2).

Convolutional Neural Networks (CNNs) have achieved remarkable results in image processing related problems. Recently, CNN models have also shown excellent results for NLP such as in semantic parsing²⁶, search query retrieval²⁷, and sentence classification²⁸. Thus we tested on a series of experiments with CNNs for our problem. In general, we applied the same architecture described in²⁸. We first concatenate the features into $n \times k$ feature vector, where n is the number of words, k is the number of dimensions extracted from *word2vec*. A set of convolution filters with dimension $h \times k$ is then applied to a window of h words to produce new features. The filters are applied to each possible window of words in the sentences to produce a feature map. Finally, we apply a max-over-time pooling operation over the feature map to generate the fully connected layer. A sigmoid activation function is applied to generate the multi-label output.

We tried three to ten convolutional layers with size 64, 128, 256 and one to three fully connected dense layers attached to the last convolutional layer with size 4096, 1024 and 128. Among our model architecture setting, the best performed model pipelines are shown in Figure 11 and 12. Under the hardware setting described, the training time for CNN is less than 30 minutes with 500 maximum epochs and early stop if the validation loss doesn't improve for consecutive 10 epochs.

Recurrent Neural Networks: RNNs have recently shown promising results in many machine learning tasks²⁹. We explored several RNN architectures in this study. All the architectures are follow the same patten shown in Figure 2, where blue circles represent the text feature vectors. The green rectangles represent recurrent hidden layers, and the yellow rectangle represents the multi-label code assignment. Basically, the RNN cells went through the input sentences, when they reached the last word, the hidden layer of the RNN generated the outputs \hat{y} . We use sigmoid cross-entropy as the loss function, RMSprop as the optimizer.

$$loss(\hat{y}, y) = -\frac{1}{N} \sum_{n=1}^{l=N} (y_n \cdot \log(\hat{y}_n) + (1 - y_n) \cdot \log(1 - \hat{y}_n))$$

Among those sophisticated recurrent units, in this study, We evaluated two popular ones: Long Short Term Memory (LSTM) unit³⁰ and Gated Recurrent Units (GRU) unit³¹. Both of them have the ability to capture sequence-based inputs with long-term dependencies.

The input features is the same as we used in CNN. We tried up to three stacked recurrent layers with a combination 64, 128 and 256 units for each layer in our RNNs. We used the well-known LSTM and GRU units for our RNNs. To predict the ICD-9 classification, we only consider the output nodes of the last time step, and apply the same activation

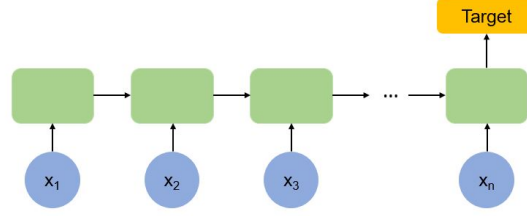


Figure 2: RNN architecture

function and loss function as we choose for Feed-Forward Neural Network. The best performed model architecture for LSTM and GRU are shown in Figure 9 and 10. Under the specified hardware setting, the training time is about 6 hours for GRU and 18 hours for LSTM with 200 maximum epochs and early stop if the validation loss doesn't improve for consecutive 5 epochs.

3.4 Metrics

Combinations of our dataset, feature extraction methods, and models are evaluated under different performance metrics, including precision, accuracy, F-score and recall metrics for multi-label classification. Specifically, the following metrics are used³²:

$$\begin{aligned} \text{Precision} &= \frac{1}{n} \sum_{i=1}^n \frac{|Y_i \cap Z_i|}{|Z_i|} & \text{Recall} &= \frac{1}{n} \sum_{i=1}^n \frac{|Y_i \cap Z_i|}{|Y_i|} \\ F_1 &= \frac{1}{n} \sum_{i=1}^n \frac{2|Y_i \cap Z_i|}{|Y_i| + |Z_i|} & \text{Accuracy} &= \frac{1}{n} \sum_{i=1}^n \frac{|Y_i \cap Z_i|}{|Y_i \cup Z_i|} \end{aligned}$$

where Y_i is the set of predicted labels, Z_i is the set of ground truth labels, and n is the number of samples. q is the number of total samples.

4 Results

This section illustrates the performance in three different aspects: (1) baseline results, (2) performance under different configurations, and (3) best model performance.

4.1 Model Performance under Different Configurations

Different model configurations are tried to give us insight into the most appropriate model configuration. Table 4 describes the different methods of feature extraction used and the parameters tweaked. The features extracted are divided into 2 categories: non-sequential and sequential. The non-sequential features include *tfidf* and *word2vec*, which were used in Logistic Regression, Random Forest, and Feed-Forward Neural Network. The sequential features includes *wordseq* (word sequences) used in conjunction with an embedding matrix based on *word2vec*, which were used in Conv1D, RNN, LSTM, and GRU. Note that we experimented on 1) using our custom *word2vec* model created from the MIMIC-III dataset and 2) on using pre-train word vectors obtained from²⁴. The vector for stop words in the embedding matrix are all zeros.

Figure 3 shows the model performance of each model using different feature extraction methods pair on the *top-10-code* dataset. The raw data are also shown in the appendix (table 13 and 14). For each model, the configuration that provided the best performance here is used on the *top-50-code*, *top-10-cat*, and *top-50-cat* datasets. The results are further explained in the next section 4.2.

For non-sequential feature extraction, *tfidf* with 20000 features gave the best f1 results for Logistic Regression and Random Forest, while *word2vec.m3* with 600 features gave the best results for Feed-Forward Neural Networks (although *tfidf* also gave a fairly good result for NN). In general, *tfidf* configurations generated better results than *word2vec*,

Feature Extraction	Configuration	Value
tfidf	feature size	20301 - 40000
	minDocFreq	3 - 10
	max_df	0.8 - 1.0
word2vec	database	self trained from MIMIC-III (m3) or pre-trained from Pubmed (pm) ²⁴
	feature size	100 - 600
	pre-trained config	win 2 or win30 ²⁴
wordseq	sequence length	1500-2000
	stopwords	removed from sequence or not removed
	embedding matrix	derived from the word2vec under different configurations

Table 4: Feature Extraction Methods

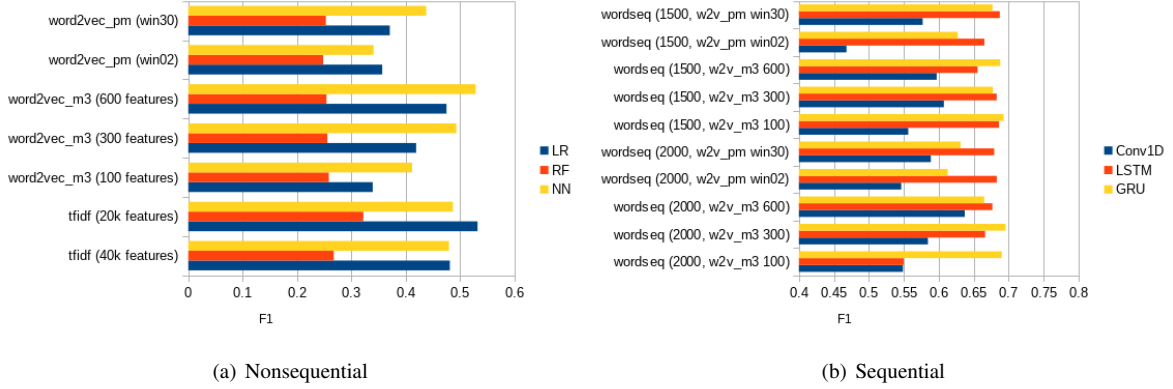


Figure 3: Model Performance under Different Configuration

probably because we are losing information. Note that the feature size for *word2vec* are at most 600 while *tfidf* is around 20,000 above.

For sequential feature extraction, *seq. length 2000 + word2vec_m3 w/ 600 features* generated the best f1 result for Conv1D, *seq. length 1500 + word2vec_pm (win30)* for LSTM, and *seq. length 2000 + word2vec_m3 w/ 300 features* for GRU. In general, all feature extraction methods generated good and comparable results for Conv1D, LSTM, and GRU. Our *word2vec* also faired well compared with the pre-trained *word2vec* models. Although not shown in figure 3(b) (but shown in table 14, the feature extraction methods were also tried for RNN, but the results were bad (0.0 - 0.23 f1 at best). This may be because the sequence length is too long (1,500 - 2,000).

4.2 Best Model Performance

4.2.1 Overview

Figure 4 and 5 shows the model performance for the *top-10-code*, *top-10-cat*, *top-50-code*, *top-50-cat* dataset (the models are ordered from best to worst, from top to bottom). Raw data are also shown in table 7, 10, 8, and 11.

For *top-10-code* and *top-10-cat*, GRU generated the best f1 results (at 0.6957 and 0.7233, respectively). In general, *top-10-cat* generated slightly better results than *top-10-code*. This makes sense because 1) we have more data per labels in *top-10-cat* and 2) the labels are less specific (the differences between labels are larger). The baseline models (Logistic Regression and Random Forest) seems to overfit the data (which can be observed from the almost 100% training results but bad test results). Feed-Forward NN is not overfitting, but the results are comparable to the baseline models. Conv1D produced even better results (than NN), but there are more significant improvement with LSTM and GRU (reaching around 70% f1). This signifies that our LSTM and GRU model was able to extract information from the sequence of words, otherwise lost in non-sequential feature extraction, thereby improving the f1 and also the

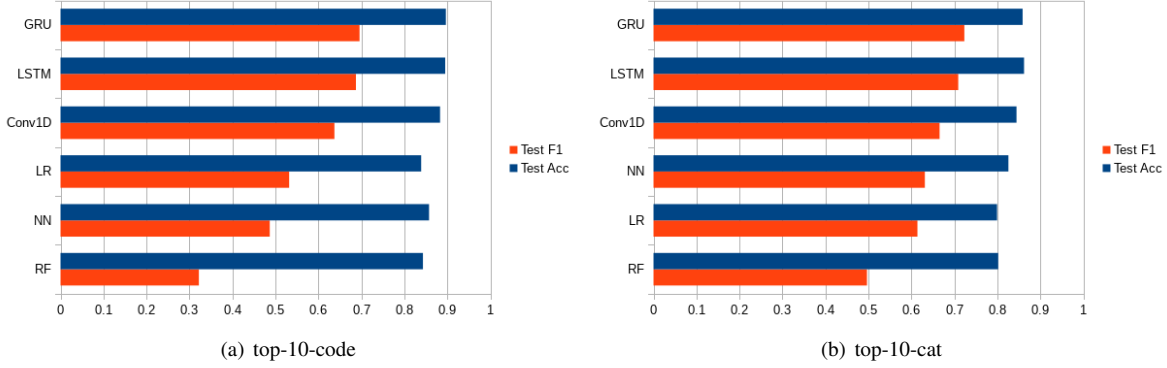


Figure 4: Model Performance Top 10

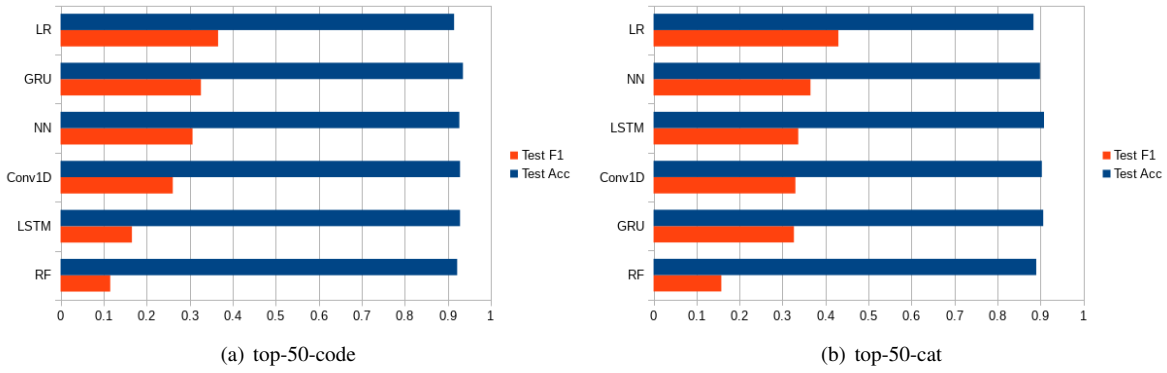


Figure 5: Model Performance Top 50

accuracy results.

For *top-50-code* and *top-50-cat*, Logistic Regression generated the best f1 result (at 0.3662 and 0.4301, respectively). Similar to *top10*, *top-50-cat* generated slightly better results than *top-50-code*. The baseline models (Logistic Regression and Random Forest) also overfit here. However, the models that used sequential feature extraction did not produce better results (in comparison with *top10*).

4.2.2 Precision-Recall Curve

Table 5 shows the average of overall precision performance of our selected best performance models for GRU, LSTM and Convolution 1D. From this table, we can see that GRU generated the best precision results for *top-10-code* and *top-50-code*. Figure 6 shows precision-recall curve for the best performed models for each labels in *top-10* and the best performed 10 labels for *top-50*. The detailed class-wise precision-recall curve for the selected models are presented in Appendix H

Model	<i>top-10-code</i>	<i>top-10-cat</i>	<i>top-50-code</i>	<i>top-50-cat</i>	<i>top-50-code(first10)</i>	<i>top-50-cat(first10)</i>
<i>LSTM</i>	0.7243	0.7915	0.3715	0.4929	0.7571	0.8426
<i>GRU</i>	0.7362	0.7849	0.4518	0.4792	0.7949	0.8425
<i>Conv1D</i>	0.6719	0.7293	0.3757	0.4565	0.7424	0.8269

Table 5: Average Precision Performance

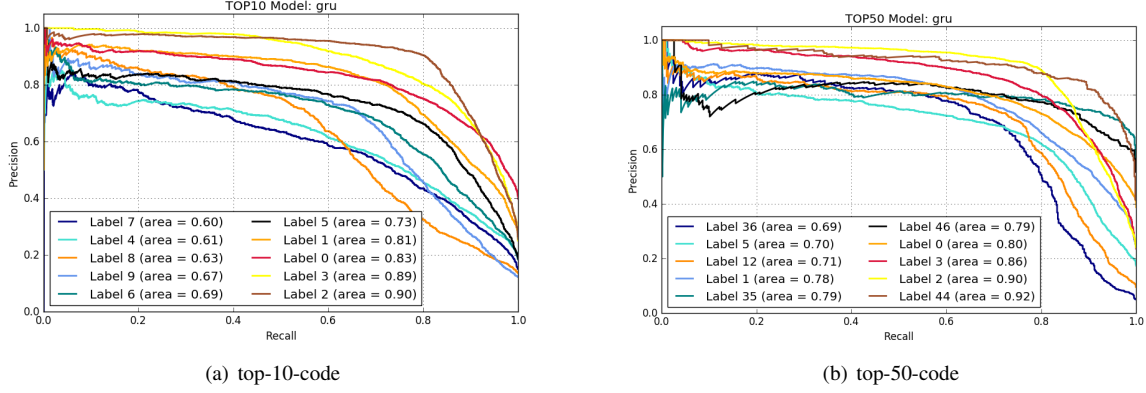


Figure 6: Class-wise Precision-Recall Curve for Top 10 and Top 50

4.2.3 top-10 and top-50 comparison

Figure 7 shows a comparison between the *top-10* results and the corresponding label results (first 10) in *top-50*. For models between ICD-9 codes, the first-10 results are different from the top-10 results, implying that the *top-50-code* model is unable to completely cover the *top-10-code* model capability. For models between ICD-9 categories, the first-10 results are similar to the top-10 results, implying that *top-50-cat* model is somewhat able to cover the *top-10-cat* model capability.

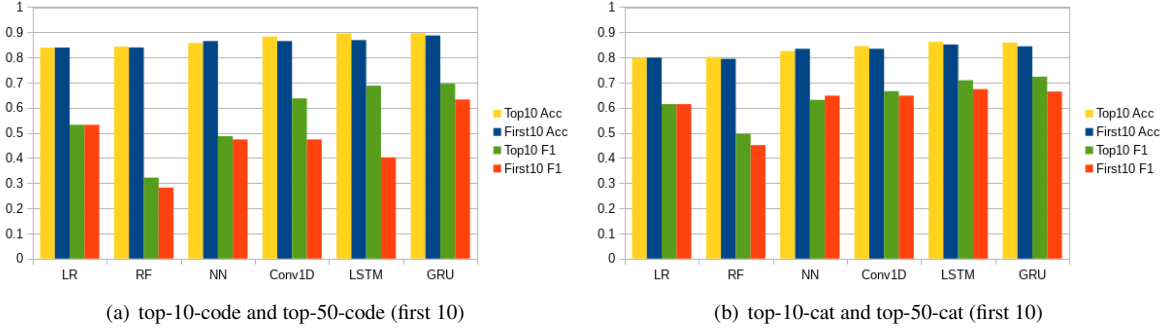


Figure 7: Top 10 vs Top 50 (first 10)

4.2.4 Results Comparison

Prakash and Zhao¹⁹ use bag-of-words from discharge notes and Condensed Memory Neural Network (C-MemNN) to tackle the similar problem as we do. They tested their algorithm with top 50 and top 100 labels under metrics such as macro average of Area Under the Curve (AUC), average precision over the top five predictions, and hamming loss.

$$AUC_{macro} = \frac{1}{q} \sum_{j=1}^q AUC_j$$

$$Hamming Loss = \frac{1}{q} \sum_{i=1}^q \frac{xor(Z_i, Y_i)}{n}$$

To compare our work with theirs, We use macro AUC and hamming loss for our best performed models for top 50 codes (top 100 labels are not compared), and the results are listed in Table 6.

From this comparison, we can see that while their hamming loss is better than ours, Our work outperformed on macro AUC on GRU model and have significantly better performance on top 5 precision for all of our models.

Model	AUC (macro)	Precision @5	Hamming Loss
<i>C-MemNN</i> ¹⁹	0.833	0.42	0.01
<i>GRU</i>	0.8599	0.8109	0.0645
<i>LSTM</i>	0.8298	0.8054	0.0714
<i>Conv1D</i>	0.8302	0.7998	0.0714

Table 6: Performance comparison with reference¹⁹

5 Conclusion

In this study, we evaluated different NLP deep learning based models and feature extraction methods, effectively establishing an empirical evaluation for learning-based automatic code assignment from the MIMIC-III discharge summary. The models are based on deep learning NLP frameworks that automatically assign clinical ICD-9 codes from free-text clinical notes. The deep learning models for predicting the top 10 ICD-9 codes and categories were better than our baseline models that use traditional learning algorithms (best F1 results: 69.57% GRU to 53.20% Logistic Regression, and 72.33% GRU to 63.13% Feed-Forward Neural Network, respectively). We also observed that the Top 50 ICD-9 codes and categories results did not outperform our baseline (F1 results: 32.63% GRU compare to 36.62% Logistic Regression, and 33.67% GRU compare to 36.51% Feed-Forward Neural Network). We hope our implementation and evaluation of the current state-of-the-art algorithms can serve as a baseline for further research on this topic.

6 Future Work

To further improve the prediction accuracy, we believe that more advanced networks architectures should be implemented for our data. Although LSTM and GRU are capable of capturing long-term dependencies, the length of our input sequence could still be too long for LSTM and GRU to retain useful information. A different representation may be used to shorten the sequence, e.g. sentence2vec or paragraph2vec³³. In addition, from the comparison in Section 4.2.4, memory networks provide better results for certain metrics. We believe word2vec representation plus memory network could lead a further progress on this problem.

Our current models for top 50 ICD-9 codes and categories were not as successful. This can be because our current model design lacks "capability" in effectively distinguishing between 50 different labels. To improve our model capability, we can try to run 5 *top-10* models in parallel (each model predicting 10 labels), thereby making our *top-50* models have the same model capability as our *top-10* models. Figure 8 shows the result of trying this approach for Conv1D using the *top-50-code* dataset. The parallel approach did improve F1 from 0.2609 to 0.3879. We also observed that there are only hundreds of positive samples for labels 11 to 50, which might not be sufficient for the deep neural network to learn enough useful representations.

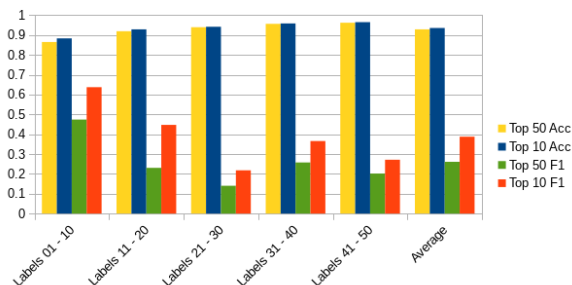


Figure 8: Parallel Top 10 vs Top 50

Our custom *word2vec* model used CBOW. We didn't have the chance to try skip-gram (although the pre-trained *word2vec* induced from PubMed are skip-gram based). Some previous studies say that skip-gram outperforms CBOW

in biomedical domain tasks²⁴. Therefore in future work, we would like to see how the different *word2vec* parameters affect our ICD-9 code or category classifier.

Further research on what words affect the probability of a prediction could improve our understanding of the relationship between symptoms and diagnosis (such as Attention Network). The probability observation could also change our preprocessing and feature extraction methods and ultimately improve our deep learning models.

References

1. Ashly D Black, Josip Car, Claudia Pagliari, Chantelle Anandan, Kathrin Cresswell, Tomislav Bokun, Brian McKinstry, Rob Procter, Azeem Majeed, and Aziz Sheikh. The impact of ehealth on the quality and safety of health care: a systematic overview. *PLoS Med*, 8(1):e1000387, 2011.
2. Edward Choi, Mohammad Taha Bahadori, and Jimeng Sun. Doctor ai: Predicting clinical events via recurrent neural networks. *arXiv preprint arXiv:1511.05942*, 2015.
3. Zachary C Lipton, David C Kale, Charles Elkan, and Randall Wetzell. Learning to diagnose with lstm recurrent neural networks. *arXiv preprint arXiv:1511.03677*, 2015.
4. F Martin-Sanchez, K Verspoor, et al. Big data in medicine is driving big changes. *Yearb Med Inform*, 9(1):14–20, 2014.
5. Leah S Larkey and W Bruce Croft. Automatic assignment of icd9 codes to discharge summaries. Technical report, Technical report, University of Massachusetts at Amherst, Amherst, MA, 1995.
6. John P Pestian, Christopher Brew, Paweł Matykiewicz, Dj J Hovermale, Neil Johnson, K Bretonnel Cohen, and Włodzisław Duch. A shared task involving multi-label classification of clinical free text. In *Proceedings of the Workshop on BioNLP 2007: Biological, Translational, and Clinical Language Processing*, pages 97–104. Association for Computational Linguistics, 2007.
7. MBA Ira Goldstein and MLS Anna Arzumtsyan. Three approaches to automatic assignment of icd-9-cm codes to radiology reports. 2007.
8. Serguei VS Pakhomov, James D Buntrock, and Christopher G Chute. Automating the assignment of diagnosis codes to patient encounters using example-based and machine learning techniques. *Journal of the American Medical Informatics Association*, 13(5):516–525, 2006.
9. Julia Medori and Cédric Fairon. Machine learning and features selection for semi-automatic icd-9-cm encoding. In *Proceedings of the NAACL HLT 2010 Second Louhi Workshop on Text and Data Mining of Health Documents*, pages 84–89. Association for Computational Linguistics, 2010.
10. Berthier Ribeiro-Neto, Alberto HF Laender, and Luciano RS De Lima. An experimental study in automatically categorizing medical documents. *Journal of the Association for Information Science and Technology*, 52(5):391–401, 2001.
11. Stéphane M Meystre, Guergana K Savova, Karin C Kipper-Schuler, John F Hurdle, et al. Extracting information from textual documents in the electronic health record: a review of recent research. *Yearb Med Inform*, 35(128):44, 2008.
12. Yuan Ling. *Methods and Techniques for Clinical Text Modeling and Analytics*. PhD thesis, Drexel University, 2017.
13. Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
14. Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. Speech recognition with deep recurrent neural networks. In *Acoustics, speech and signal processing (icassp), 2013 ieee international conference on*, pages 6645–6649. IEEE, 2013.

15. Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12(Aug):2493–2537, 2011.
16. Edward Choi, Andy Schuetz, Walter F Stewart, and Jimeng Sun. Using recurrent neural network models for early detection of heart failure onset. *Journal of the American Medical Informatics Association*, 24(2):361–370, 2016.
17. Edward Choi, Mohammad Taha Bahadori, Jimeng Sun, Joshua Kulas, Andy Schuetz, and Walter Stewart. Retain: An interpretable predictive model for healthcare using reverse time attention mechanism. In *Advances in Neural Information Processing Systems*, pages 3504–3512, 2016.
18. Yuan Luo. Recurrent neural networks for classifying relations in clinical notes. *Journal of biomedical informatics*, 72:85–95, 2017.
19. Aaditya Prakash, Siyuan Zhao, Sadid A Hasan, Vivek Datla, Kathy Lee, Ashequl Qadir, Joey Liu, and Oladimeji Farri. Condensed memory networks for clinical diagnostic inferencing. *arXiv preprint arXiv:1612.01848*, 2016.
20. MIMIC-III, a freely accessible critical care database. *Scientific data*, 3:160035, 2016.
21. Priyanka Nigam. Applying deep learning to icd-9 multi-label classification from medical records. pages 1–8.
22. Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781, 2013.
23. Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed representations of words and phrases and their compositionality. In *Proceedings of the 26th International Conference on Neural Information Processing Systems*, NIPS’13, pages 3111–3119, USA, 2013. Curran Associates Inc.
24. Billy Chiu, Gamal Crichton, Anna Korhonen, and Sampo Pyysalo. How to Train Good Word Embeddings for Biomedical NLP. *Proceedings of the 15th Workshop on Biomedical Natural Language Processing*, pages 166–174, 2016.
25. Grigorios Tsoumakas and Ioannis Katakis. Multi-label classification: An overview. *International Journal of Data Warehousing and Mining*, 3(3), 2006.
26. Wen-tau Yih, Kristina Toutanova, John C Platt, and Christopher Meek. Learning discriminative projections for text similarity measures. In *Proceedings of the Fifteenth Conference on Computational Natural Language Learning*, pages 247–256. Association for Computational Linguistics, 2011.
27. Yelong Shen, Xiaodong He, Jianfeng Gao, Li Deng, and Grégoire Mesnil. Learning semantic representations using convolutional neural networks for web search. In *Proceedings of the 23rd International Conference on World Wide Web*, pages 373–374. ACM, 2014.
28. Yoon Kim. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*, 2014.
29. Alex Graves et al. *Supervised sequence labelling with recurrent neural networks*, volume 385. Springer, 2012.
30. Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
31. Kyunghyun Cho, Bart Van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*, 2014.
32. Min-Ling Zhang and Zhi-Hua Zhou. A review on multi-label learning algorithms. *IEEE transactions on knowledge and data engineering*, 26(8):1819–1837, 2014.
33. Quoc V. Le and Tomas Mikolov. Distributed representations of sentences and documents. *CoRR*, abs/1405.4053, 2014.

Appendices

A Model Performance for Top 10 Label Codes

Model	Training				Test			
	Precision	Recall	Accuracy	F1	Precision	Recall	Accuracy	F1
Logistic Regression	0.9564	0.9440	0.9786	0.9501	0.5801	0.4934	0.8392	0.5320
Random Forest	0.9989	0.6988	0.9501	0.8086	0.7573	0.2340	0.8432	0.3219
Feed-Forward NN	0.7768	0.5041	0.8879	0.5763	0.6703	0.4193	0.8575	0.4868
Conv1D	0.8312	0.6713	0.9165	0.7371	0.7408	0.5687	0.8832	0.6373
LSTM RNN	0.8106	0.6971	0.9154	0.7445	0.7574	0.6380	0.8950	0.6874
GRU RNN	0.7936	0.6971	0.9126	0.7397	0.7502	0.6519	0.8967	0.6957

Table 7: Model Performance for *top-10-code*

B Model Performance for Top 50 Codes

Model	Training				Test			
	Precision	Recall	Accuracy	F1	Precision	Recall	Accuracy	F1
Logistic Regression	0.9863	0.9768	0.9945	0.9815	0.4372	0.3213	0.9148	0.3662
Random Forest	0.9985	0.2852	0.9451	0.3866	0.5377	0.0953	0.9220	0.1155
Feed-Forward NN	0.9636	0.5542	0.9640	0.6892	0.5773	0.2335	0.9271	0.3067
Conv1D	0.6085 ^a	0.2663	0.9365	0.3200 ^a	0.4792 ^a	0.2169	0.9286	0.2609 ^a
LSTM RNN	0.3526 ^a	0.1642	0.9325	0.1891 ^a	0.4022 ^a	0.1445	0.9286	0.1659 ^a
GRU RNN	0.6539 ^a	0.3433	0.9460	0.3947 ^a	0.5592 ^a	0.2782	0.9354	0.3263 ^a

Table 8: Model Performance for *top-50-code*

Model	Training				Test			
	Precision	Recall	Accuracy	F1	Precision	Recall	Accuracy	F1
Logistic Regression	0.9564	0.9440	0.9786	0.9501	0.5801	0.4934	0.8392	0.5320
Random Forest	0.9946	0.4937	0.9110	0.6305	0.7869	0.2009	0.8395	0.2822
Feed-Forward NN	0.9347	0.8037	0.9580	0.8589	0.7674	0.5837	0.9062	0.6486
Conv1D	0.7708	0.4673	0.8858	0.5377	0.6784	0.4109	0.8650	0.4739
LSTM RNN	0.6204 ^a	0.3829	0.8805	0.4348 ^a	0.5748	0.3526	0.8688	0.4025 ^a
GRU RNN	0.8351	0.6474	0.9168	0.7181	0.7520	0.5618	0.8871	0.6328

Table 9: Model Performance for *top-50-code* (first 10)

^aresult contained *nan*. Computed by replacing *nan* with zero.

C Model Performance for Top 10 Label Categories

Model	Training				Test			
	Precision	Recall	Accuracy	F1	Precision	Recall	Accuracy	F1
Logistic Regression	0.9437	0.9309	0.9652	0.9372	0.6458	0.5856	0.7994	0.6141
Random Forest	0.9983	0.8134	0.9500	0.8954	0.7653	0.3801	0.8019	0.4966
Feed-Forward NN	0.7978	0.6575	0.8655	0.7147	0.7243	0.5689	0.8257	0.6313
Conv1D	0.8039	0.6637	0.8681	0.7128	0.7613	0.6126	0.8446	0.6657
LSTM RNN	0.8146	0.6807	0.8749	0.7343	0.7926	0.6536	0.8622	0.7090
GRU RNN	0.8150	0.7613	0.8909	0.7861	0.7580	0.6941	0.8588	0.7233

Table 10: Model Performance for *top-10-cat***D Model Performance for Top 50 Label Categories**

Model	Training				Test			
	Precision	Recall	Accuracy	F1	Precision	Recall	Accuracy	F1
Logistic Regression	0.9750	0.9572	0.9887	0.9659	0.4858	0.3894	0.8841	0.4301
Random Forest	0.9986	0.3294	0.9277	0.4465	0.6568	0.1142	0.8906	0.1576
Feed-Forward NN	0.9613	0.5488	0.9473	0.6853	0.6298	0.2773	0.8992	0.3651
Conv1D	0.7428	0.3262	0.9163	0.3870	0.5635 ^a	0.2770	0.9035	0.3301 ^a
LSTM RNN	0.7117 ^a	0.3363	0.9194	0.3804 ^a	0.5869	0.2945	0.9087	0.3367 ^a
GRU RNN	0.6695 ^a	0.3227	0.9179	0.3726 ^a	0.5611 ^a	0.2809	0.9067	0.3266 ^a

Table 11: Model Performance for *top-50-cat*

Model	Training				Test			
	Precision	Recall	Accuracy	F1	Precision	Recall	Accuracy	F1
Logistic Regression	0.9437	0.9309	0.9652	0.9372	0.6458	0.5856	0.7994	0.6141
Random Forest	0.9937	0.6321	0.8999	0.7687	0.7877	0.3282	0.7944	0.4512
Feed-Forward NN	0.7873	0.6425	0.8811	0.7031	0.7730	0.6267	0.8724	0.6878
Conv1D	0.7945	0.6652	0.8670	0.7142	0.7296	0.5979	0.8345	0.6481
LSTM RNN	0.7963	0.6863	0.8768	0.7213	0.7515	0.6362	0.8514	0.6738
GRU RNN	0.7901	0.6803	0.8729	0.7213	0.7382	0.6196	0.8442	0.6641

Table 12: Model Performance for *top-50-cat* (first 10)^aresult contained *nan*. Computed by replacing *nan* with zero.

E Model Performance under Different Configurations (Non Sequential)

Model	Configuration	Feature	Precision	Recall	Accuracy	F1
Logistic Regression	Iter=25	tfidf (40k features)	0.5290	0.4445	0.8232	0.4810
Logistic Regression	Iter=10	tfidf (20k features)	0.5801	0.4934	0.8392	0.5320
Logistic Regression	Iter=50	word2vec_m3 (100 features)	0.5820	0.2682	0.8353	0.3393
Logistic Regression	Iter=100	word2vec_m3 (300 features)	0.6265	0.3388	0.8448	0.4191
Logistic Regression	Iter=100	word2vec_m3 (600 features)	0.6467	0.3914	0.8513	0.4749
Logistic Regression	Iter=100	word2vec_pm (win02)	0.5938	0.2807	0.8340	0.3568
Logistic Regression	Iter=75	word2vec_pm (win30)	0.6030	0.2921	0.8370	0.3706
Random Forest	Depth=30	tfidf (40k features)	0.7529	0.1881	0.8354	0.2676
Random Forest	Depth=30	tfidf (20k features)	0.7573	0.2340	0.8432	0.3219
Random Forest	Depth=20	word2vec_m3 (100 features)	0.5281	0.1907	0.8233	0.2585
Random Forest	Depth=20	word2vec_m3 (300 features)	0.5471	0.1866	0.8237	0.2559
Random Forest	Depth=20	word2vec_m3 (600 features)	0.5296	0.1857	0.8228	0.2543
Random Forest	Depth=20	word2vec_pm (win02)	0.5420	0.1794	0.8224	0.2480
Random Forest	Depth=20	word2vec_pm (win30)	0.5283	0.1852	0.8225	0.2530
Feed-Forward NN	nn_model.1	tfidf (40k features)	0.6975	0.3785	0.8542	0.4795
Feed-Forward NN	nn_model.2	tfidf (20k features)	0.6703	0.4193	0.8575	0.4868
Feed-Forward NN	nn_model.2	word2vec_m3 (100 features)	0.5429	0.3505	0.8304	0.4116
Feed-Forward NN	nn_model.2	word2vec_m3 (300 features)	0.5467	0.4627	0.8329	0.4929
Feed-Forward NN	nn_model.2	word2vec_m3 (600 features)	0.5737	0.4972	0.8406	0.5285
Feed-Forward NN	nn_model.2	word2vec_pm (win02)	0.5426	0.2778	0.8377	0.3408
Feed-Forward NN	nn_model.4	word2vec_pm (win30)	0.5228	0.3856	0.8253	0.4375

Table 13: Model Performance under Different Configurations (Non Sequential)

^aresult contained *nan*. Computed by replacing *nan* with zero.

F Model Performance under Different Configurations (Sequential)

Model	Configuration	Seq. Length	Embedding Matrix	Precision	Recall	Accuracy	F1
Conv1D	conv1d_6	2000	w2v_m3 (100 features)	0.7460	0.4626	0.8758	0.5489
Conv1D	conv1d_6	2000	w2v_m3 (300 features)	0.7492	0.5103	0.8779	0.5846
Conv1D	conv1d_6	2000	w2v_m3 (600 features)	0.7408	0.5687	0.8832	0.6373
Conv1D	conv1d_6	2000	w2v_pm (win02)	0.6424	0.4962	0.8642	0.5464
Conv1D	conv1d_6	2000	w2v_pm (win30)	0.7283	0.5217	0.8757	0.5888
Conv1D	conv1d_6	1500	w2v_m3 (100 features)	0.7304	0.4851	0.8744	0.5565
Conv1D	conv1d_6	1500	w2v_m3 (300 features)	0.7141	0.5443	0.8762	0.6074
Conv1D	conv1d_6	1500	w2v_m3 (600 features)	0.7578	0.5081	0.8805	0.5972
Conv1D	conv1d_6	1500	w2v_pm (win02)	0.5635 ^a	0.2770	0.9035	0.3301 ^a
Conv1D	conv1d_6	1500	w2v_pm (win30)	0.6819	0.5323	0.8658	0.5771
RNN	rnn_2	2000	w2v_m3 (100 features)	0.1772 ^a	0.0758	0.8067	0.08693 ^a
RNN	rnn_2	2000	w2v_m3 (300 features)	0.2291 ^a	0.0476	0.8067	0.06774 ^a
RNN	rnn_2	2000	w2v_m3 (600 features)	0.1110 ^a	0.0513	0.8052	0.0702 ^a
RNN	rnn_2	2000	w2v_pm (win02)	0.0000 ^a	0.0000	0.8025	0.0000 ^a
RNN	rnn_2	2000	w2v_pm (win30)	0.0000 ^a	0.0000	0.8025	0.0000 ^a
RNN	rnn_2	1500	w2v_pm (win02)	0.1150 ^a	0.0390	0.8045	0.0535 ^a
RNN	rnn_2	1500	w2v_pm (win30)	0.1087 ^a	0.0545	0.8025	0.0630 ^a
LSTM	lstm_1	2000	w2v_m3 (100 features)	0.6857	0.4958	0.8749	0.5499
LSTM	lstm_1	2000	w2v_m3 (300 features)	0.7577	0.6001	0.8922	0.6664
LSTM	lstm_1	2000	w2v_m3 (600 features)	0.7287	0.6381	0.8901	0.6768
LSTM	lstm_1	2000	w2v_pm (win02)	0.7464	0.6381	0.8931	0.6831
LSTM	lstm_1	2000	w2v_pm (win30)	0.7529	0.6288	0.8939	0.6794
LSTM	lstm_1	1500	w2v_m3 (100 features)	0.7275	0.6523	0.8910	0.6862
LSTM	lstm_1	1500	w2v_m3 (300 features)	0.7566	0.6281	0.8948	0.6829
LSTM	lstm_1	1500	w2v_m3 (600 features)	0.7549	0.5974	0.8908	0.6556
LSTM	lstm_1	1500	w2v_pm (win02)	0.7702	0.5984	0.8930	0.6654
LSTM	lstm_1	1500	w2v_pm (win30)	0.7574	0.6380	0.8950	0.6874
GRU	gru_4	2000	w2v_m3 (100 features)	0.7456	0.6486	0.8953	0.6902
GRU	gru_4	2000	w2v_m3 (300 features)	0.7502	0.6519	0.8967	0.6957
GRU	gru_4	2000	w2v_m3 (600 features)	0.7395	0.6164	0.8893	0.6651
GRU	gru_4	2000	w2v_pm (win02)	0.7166	0.5572	0.8804	0.6128
GRU	gru_4	2000	w2v_pm (win30)	0.7577	0.5672	0.8868	0.6313
GRU	gru_4	1500	w2v_m3 (100 features)	0.7475	0.6514	0.8955	0.6930
GRU	gru_4	1500	w2v_m3 (300 features)	0.7758	0.6111	0.8967	0.6776
GRU	gru_4	1500	w2v_m3 (600 features)	0.7557	0.6389	0.8955	0.6881
GRU	gru_4	1500	w2v_pm (win02)	0.7228	0.5773	0.8815	0.6273
GRU	gru_4	1500	w2v_pm (win30)	0.7404	0.6277	0.8907	0.6769

Table 14: Model Performance under Different Configurations (Sequential)

^aresult contained *nan*. Computed by replacing *nan* with zero.

G Best Performance Model Architectures

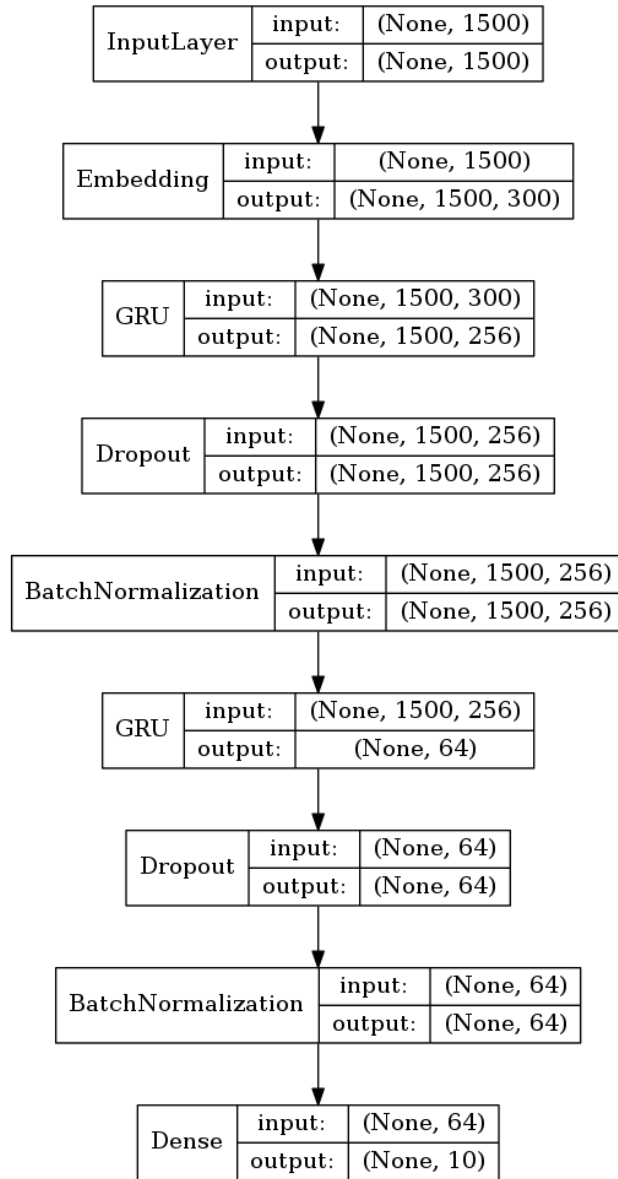


Figure 9: Best GRU Model Architecture for Top 10 Codes, Top 50 Codes, Top 10 Categories, and Top 50 Categories

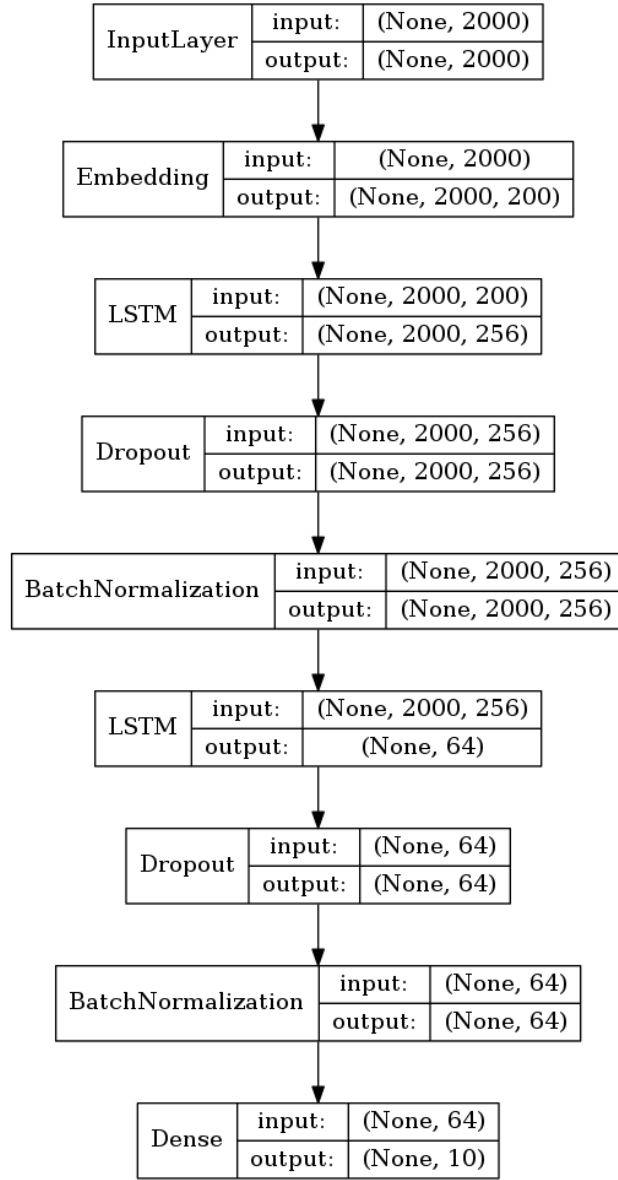


Figure 10: Best LSTM Model Architecture for Top 10 Codes, Top 50 Codes, Top 10 Categories, and Top 50 Categories

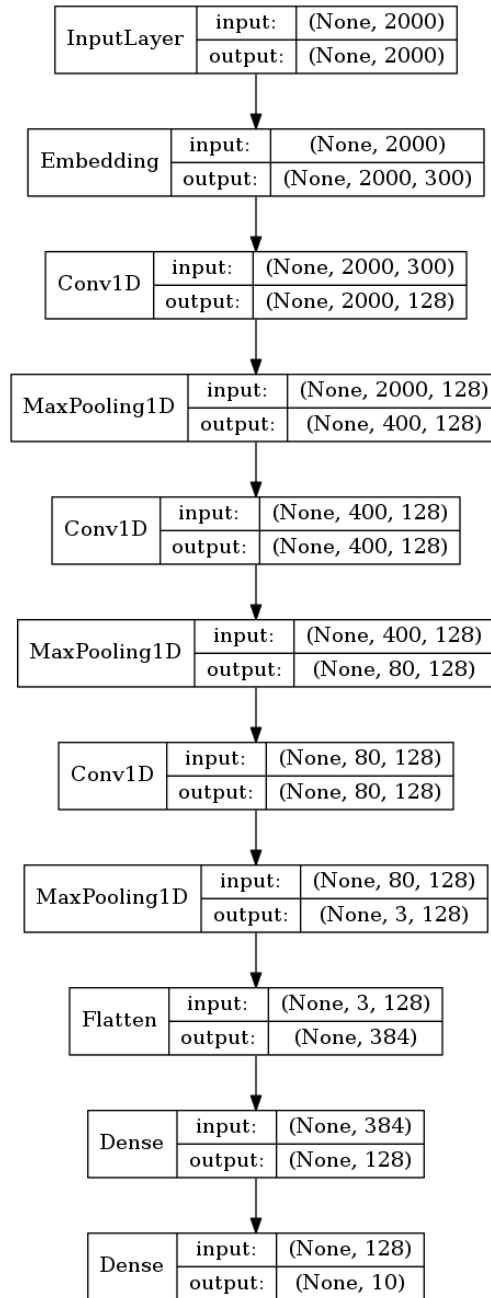


Figure 11: Best Convolution 1D Model Architecture for Top 10 Codes and Top 10 Categories

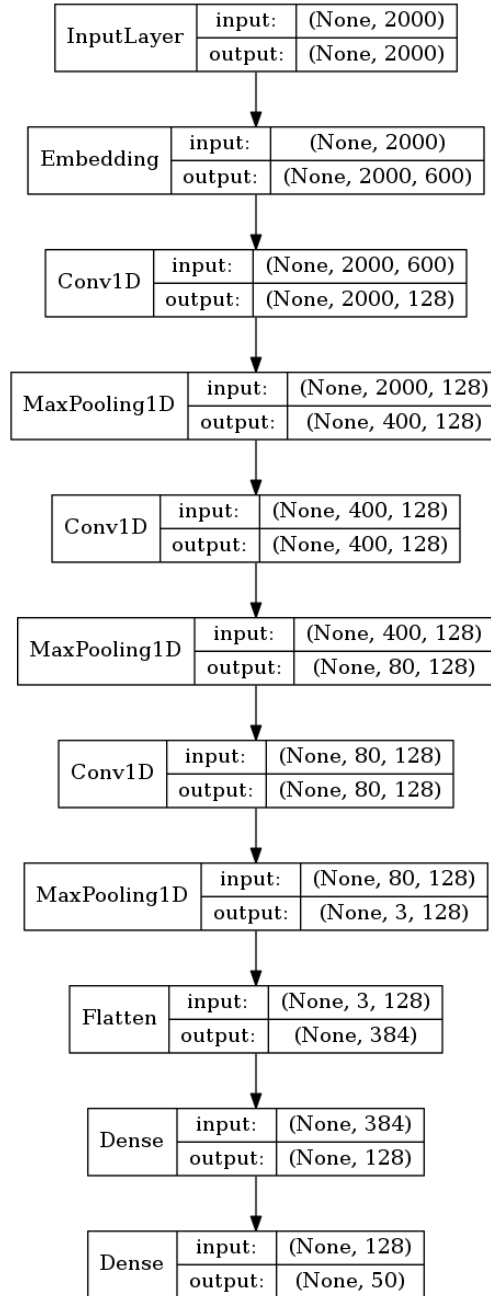
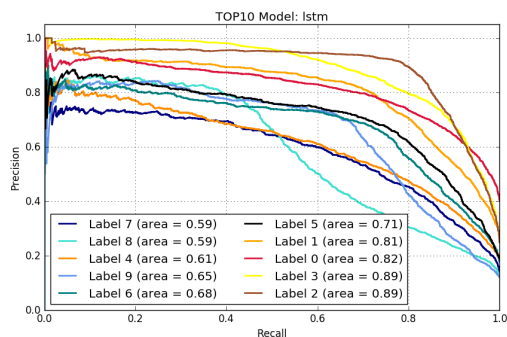
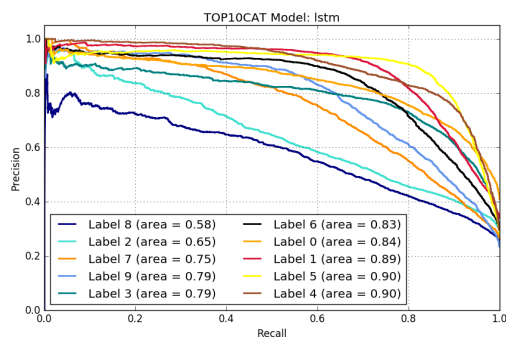


Figure 12: Best Convolution 1D Model Architecture for Top 50 Codes and Top 50 Categories

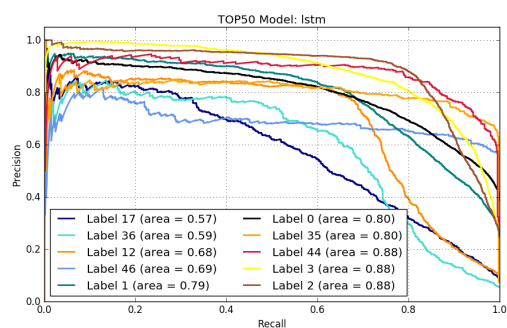
H Best Performance Models (LSTM, GRU, Conv1D) Precision-Recall Curve



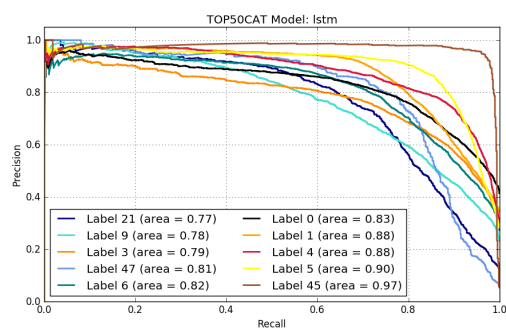
(a) LSTM Top 10 Codes



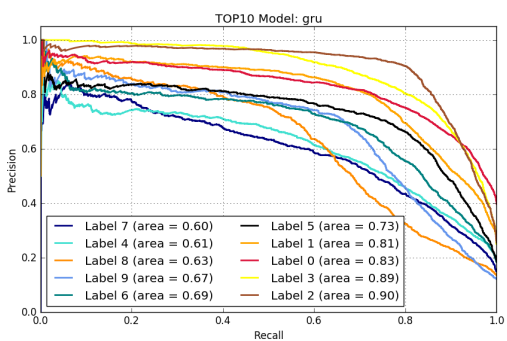
(b) LSTM Top 10 Categories



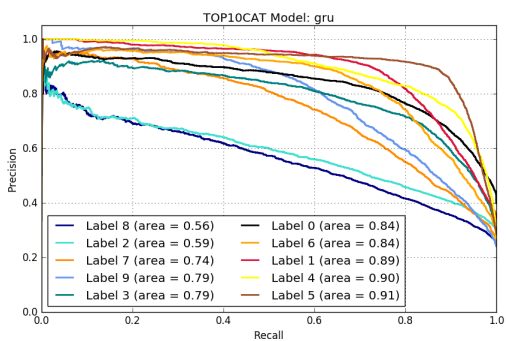
(c) LSTM Top 50 Codes



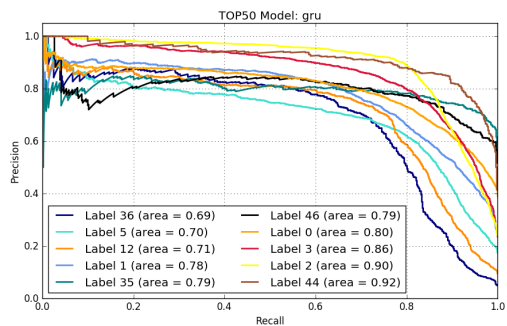
(d) LSTM Top 50 Categories



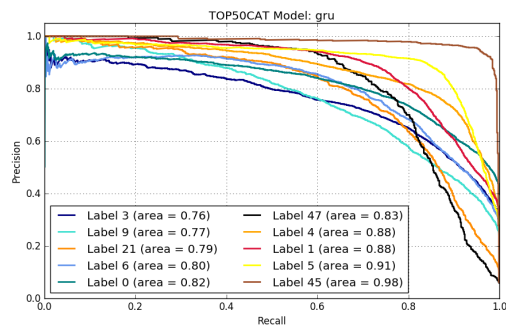
(e) GRU Top 10 Codes



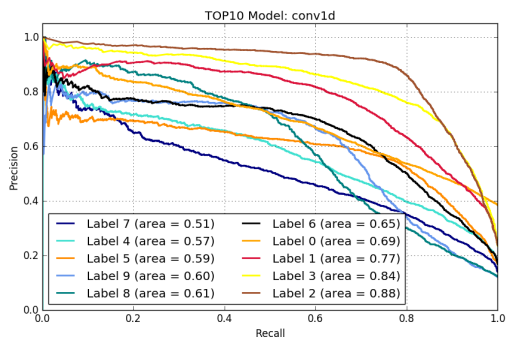
(f) GRU Top 10 Categories



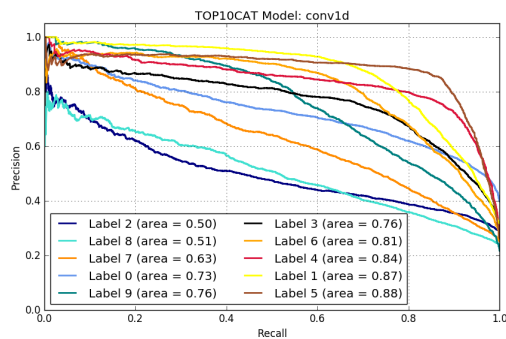
(g) GRU Top 50 Codes



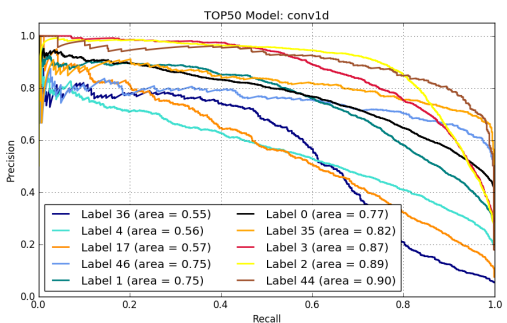
(h) GRU Top 50 Categories



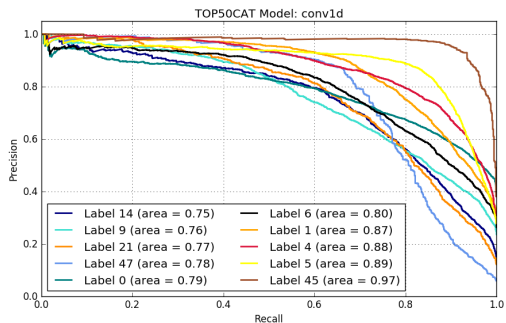
(i) Conv1D Top 10 Codes



(j) Conv1D Top 10 Categories



(k) Conv1D Top 50 Codes



(l) Conv1D Top 50 Categories