

NAMED ENTITY RECOGNITION USING MACHINE LEARNING

Kiran Kumar.B¹, B.M.Bhavya²

PG Scholar, Dept. of MCA, PES College of Engineering, Mandya, Karnataka, India¹

Assistant Professor, Dept. of MCA, PES College of Engineering, Mandya, Karnataka, India²

Abstract: Cultural heritage comes to life if it can effectively spread knowledge and foster societal cultural development. Therefore, communication is essential for giving visitors—especially tourists—information. Connecting with clients in the heritage business requires the use of information technology, which offers a variety of channels via which information can be provided. The quality and quantity of visitors' interactions with heritage can be improved thanks to technological developments like QR Code that have increased in recent years. On signs and advertisements, there are squares with a black-and-white pattern known as quick response (QR) codes that link to more information. Numerous uses for QR codes have been developed throughout time. They were initially created in the middle of the 1990s as a more practical alternative to traditional barcodes for usage in the automotive industry.

Keywords: Heritage, QRCode, Cloakroom, Information.

I. INTRODUCTION

Different entities are to be described and assigned labels using NER. Typically, entities like a person, business, place, date, and so forth. Many Natural Language Processing (NLP) activities, such as information extraction, information retrieval, feature extraction, and so on, all begin with NER. In order to solve this problem of sequential labelling, models train and apply named entity tags to each token in a sentence. The portion of speech, lexical traits, adjacent words, and other factors heavily influence performance. The identification of each is aided by pre-trained sentence vectors. As a NER feature, the Bi-LSTM result will be used to label each word or token (Jin and Yu 2021). Random Forest joins the words before sending the word sequence into BI-LSTM.

These methods reduce vocabulary by using LSTM at the syllable level (OOV). NER is a critical stage in the search and classification of objects. A branch of information extraction called NER (IE). NER is playing a bigger role in applications of natural language processing (NLP). According to the author (Melinamath 2014), the rule-based technique requires substantial knowledge of grammar as well as months of research. The study of words with subtle grammatical changes that modify their meaning is called morphology. The addition of words does not change the meaning of a phrase. Both morphologically and conceptually, it is difficult to discern between these nouns. Only a few number of NER, chunking, and parts of speech taggers have been recorded.

NER is important to natural language processing (NLP). Identity of an Entity Recognition and classification are two of Natural Language Processing's most crucial and challenging tasks. Little work has been documented for this language, and resources like annotated corpora, name dictionaries, and parts of speech (POS) tagging are not yet widely available. This project aims to develop a new NER model for Bi-LSTM and BERT. English capitalizes words to distinguish between entities, and NER research has recently risen to the top of the most-studied topics lists.

II. LITERATURE SURVEY

As NLP research and development continue to grow, neural networks have supplanted classical machine learning as the industry standard for research (LeCun, Bengio & Hinton, 2015). The neural network and word representation methods enhance the feature selection and general performance of various NLP tasks based on a number of new optimization techniques. This study provides a quick overview of the history of NLP development through a review of prior studies. It has helped us determine the procedures to carry out the experiment and respond to the sub-research questions. After that, we discussed NER and the evolution of the techniques used for NER. Each method utilised in this project is primarily introduced in the second section of this chapter. We chose the SVM, the HMM, and the CRF as our implementation methods after reviewing the approaches used for NER that were based on machine learning. The Bi-LSTM approach was also one of the options we considered for this project. Following a review of the four approaches we selected for this project, we have the fundamental framework for the workflow's implementation. Data pre-processing and data collecting are the first steps. The next stage is to create the models utilising the methodologies and algorithms. We must choose the features that will assist the model in order to recognise the designated entity in the articles.

III. EXISTING SYSTEM

Traditionally, the named entities have been recognised using dictionary- and rule-based methods. They rely on entity libraries and hand-crafted rules for the grammar of human language.

The cost of the rule- or dictionary-based approach is a drawback, as well as issues with the lengthy system design process, low portability, and the requirement for developing additional domain expertise to improve system identification capabilities.

IV. PROPOSED SYSTEM

In order to identify and recognise entity tags based on machine learning, people first construct statistical models through supervised learning. To increase the performance of models in supervisor learning, researchers switch from discriminative models to generative models.

Two generative models—the Hidden Markov model (HMM) and the condition Random Field—as well as one discriminative model—the support vector machine (SVM)—are used (CRF). Widely used for the sequence label problem with long-term reliance is the bidirectional long short-term memory network (Bi-LSTM).

V. FLOW DIAGRAM

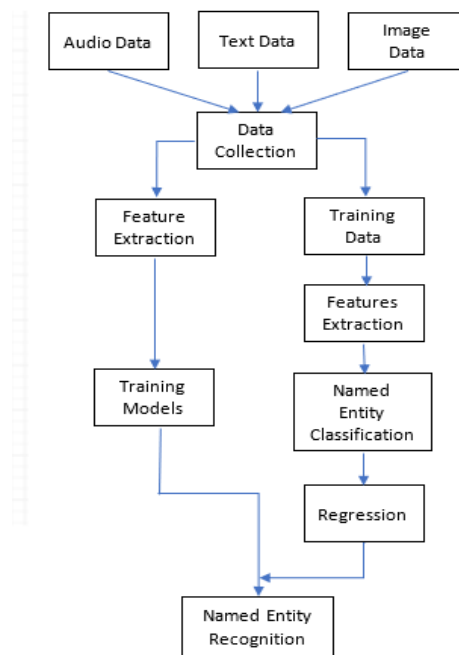


Fig 1 Flow Diagram Of the Syatem

VI. METEDODOLOGY

A. DATASET

Using the GMB (Groningen Meaning Bank) corpus for entity classification with improved and well-liked features, the data collection was applied to an annotated corpus for named entity recognition.

Content:

This excerpt from the GMB corpus was created with the express purpose of training the classifier that can predict named entities like location, name, etc.

Number of tagged entities:

Geo-name: 58388, "O": 1146068 48034, "org-nam," Pernambucan: 23790, 20680, "gpe-nam," Timing Date: 12786, Tim-dow (11404), per title: 9800, per family: 8152 "tim-yoc": 5290 Tim-moy (4262), per donation: 2413 "tim-clo": 891 866, "art-nam," Eve Nam: 602, Nat-Vietnam: 300, Timnam (146), 107, "eve-ord," "per-ini" 60, the following: 60, Per Order: 38, Timdom: 10 Per middle: 1 "art-add": 1

Important information concerning entities:

- geo = Geographical Entity
- org = Organization
- per = Person
- gpe = Geopolitical Entity
- tim = Temporal indication
- art = Artifact
- eve = Event
- nat = a phenomenon of nature

Total Number of Words = 1354149

Objective Data Column: "tag"

Inspiration: Interest in this dataset has grown as a result of the recent addition of more items. Additionally, having a thorough understanding of feature engineering in connection to this dataset is beneficial.

Why is this dataset entertaining or useful?

When you have more features and can choose intent and custom named entities created using your own phrase, it may not sound as interesting as it does in earlier versions, but it becomes more interesting and helps to find solutions to real business issues (like picking entities from Electronic Medical Records, etc)

B. SEQUENTIAL MODEL

It is a straightforward and simple model. A `tf.keras.Model` is made up of a linear stack of methods that group a linear stack of layers. Its primary function, as implied by its name, is to sequentially arrange the Keras layers.

Data flow from one layer to another in this paradigm. Up until the data reaches the top layer, the data flow is maintained. The Sequential API Model is used by the majority of ANN's.

Create a Sequential Model

Calling the `Sequential()` function will produce a Sequential model. The function must be given a list of instances `Object()` { [native code] }.

```
from keras.models import Sequential  
model = Sequential()
```

Adding a Layer in the Sequential Model

In the Sequential Model, adding a layer is quite simple. Any layer can be quickly added using the `add()` method. You must first construct a layer using the API before passing it through the `add()` method in order to use it.

```
model = Sequential()  
model.add(Dense(32, input_dim=784))  
model.add(Activation('relu'))
```

Specifying the Input Shape

The input to the layer must be specified in terms of shape. Only the first layer's form needs to be specified. Because the output of the first layer immediately becomes the input of the layers that follow,

1. A variety of techniques are available to specify the input shape to the layer. The `input_shape` argument is available. This argument may be passed to the bottom tier.
2. There is an argument called `input_dim` for 2D layers like `Dense`. The `input_dim` and `input_length` properties of 3D layers allow the model's input shape to be specified.

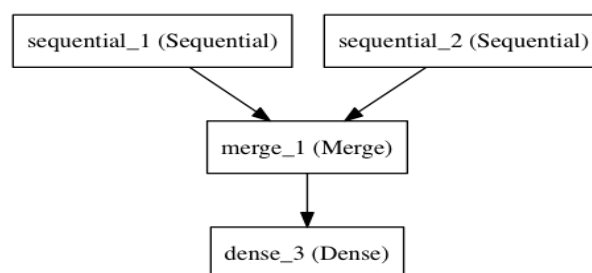


Fig 2 Sequence Model

C. BI-LSTM MODEL

We may employ the bidirectional LSTM network introduced in because when performing the sequence tagging task, we have access to both past and future input data at a specific period (Graves et al., 2013). Through the usage of forward states, this enables us to effectively utilise both previous features for a certain period of time and future features (via backward states). Recursive backpropagation through time (BPTT) is used to train bidirectional LSTM networks (Boden., 2002). Similar to conventional network forward and backward passes, Over time, forward and reverse passes are made over the unfolded network, with the exception that we must unfold the concealed states at each time step. Additionally, the start and the last data point needs serious consideration. We just need to reset the hidden states to 0 at the beginning of each phrase in our implementation because we move ahead and backward for whole phrases. We can process numerous sentences at once thanks to our batch implementation.

CRF networks

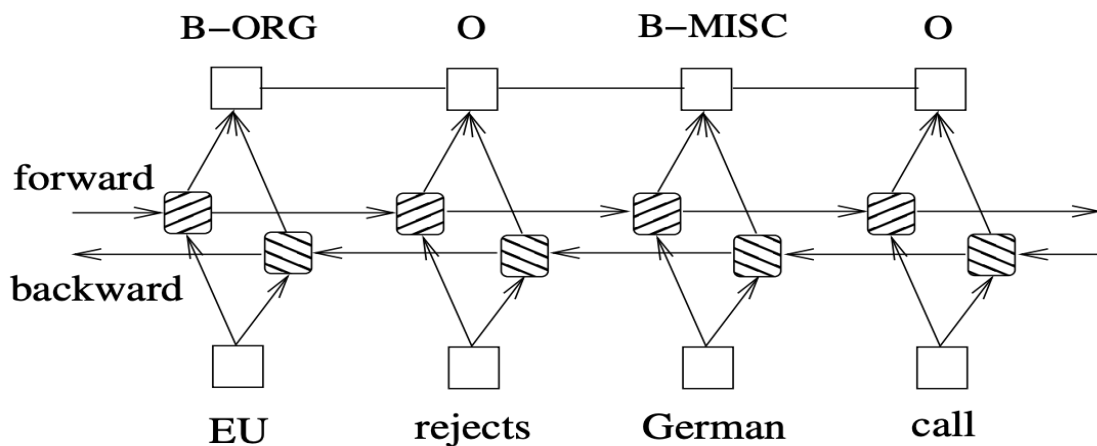


Fig 3 CRF Networks

In order to forecast current tags, neighbour tag information can be used in two different ways. The first stage is to forecast the distribution of tags for each instance.

Take a step, then use beam-like decoding to determine the optimum tag sequences. This includes the Markov models with maximum entropy developed by McCallum et al. and Ratnaparkhi's (1996) maximum entropy classifier (MEMMs). The second one uses Conditional Random Fields (CRF) models, which concentrate on sentences as a whole rather than specific positions (Lafferty et al., 2001). Memory cells and recurrent components in LSTM and bidirectional LSTM networks are superior, the inputs and outputs are directly coupled in this model. It has been proven that CRFs can provide tags with a higher level of accuracy in general. The findings of this research show the superiority of BI-LSTM over LSTM and reveal an intriguing relationship between the two techniques of using tag information and there are two ways to use input attributes (see the LSTM and BI-LSTM networks mentioned above). German request is rejected by BORG O BMISC O EU.

BI-LSTM-CRF networks

By combining a bidirectional LSTM network with a CRF network, similarly to an LSTM-CRF network, we create a BI-LSTM-CRF network. In addition to the past input features and sentence level tag data utilised in an LSTM-CRF model, a BILSTM-CRF model can also make use of future input features. We will conduct experiments to demonstrate how the additional features can improve tagging accuracy.

Training Procedure:

The models utilised in this study are all trained using the same generic SGD forward and backward method. To show how the training procedure in Algorithm 1 works, We employ the most intricate model, BI-LSTM- CRF. We split up the entire training data into batches for each epoch and process each batch separately. The number of phrases in each batch is regulated by the batch size option. In our testing, we use a batch size of 100, meaning that no more than 100 sentences

can be included in the batch. For each batch of the bidirectional LSTM-CRF model, For both the forward state and the backward state of the LSTM, we first carry out a forward pass. As a result, we are able to extract the output score f for all tags and locations ($[x]T1$). Afterward, we compute gradients for network output and state transition edges using CRF layer forward and backward passes. The backward pass for both the forward and backward states of the LSTM may then be used to propagate the faults from the output to the input.

Algorithm 1 for the bidirectional LSTM CRF model updates the network parameters last, including the initial bidirectional LSTM parameters and the state transition matrix $[A]_{i;j}8i; j$.

Training Procedure

```
1: for each epoch do
2: for each batch do
3: 1) bidirectional LSTM-CRF model forward pass:
4: forward pass for forward state LSTM
5: forward pass for backward state LSTM
6: 2) CRF layer forward and backward pass
7: 3) bidirectional LSTM-CRF model backward pass:
8: backward pass for forward state LSTM
9: backward pass for backward state LSTM
10: 4) update parameters
11: end for
12: end for
```

Modeling

The Bi-LSTM model is a significant model examined in this research investigation. BiLSTM may be better and more comparable in predicting the items that are present in the test dataset. Two layers of LSTM are utilised, as the name implies. The input word determines the shape and maximum length of the sentence, and the dropout is fixed at 0.1 to prevent overfitting. This is how the Bi-LSTM model for NER is described. The pace of learning is set to 100. The dense procedure is subjected to a "softmax" activation and a time distributed function. The model function is used to determine the remaining parameters, and the model architecture is

Result and Evaluation

The model outcome and study assessment measures are described in this section. Three models have been applied to this dataset, and two deep learning models and one machine learning model form the basis for the assessment measures. These models' output is explained in the manner as follows. 5.1

Bi-LSTM deep learning model, model 1. The values utilised to generate the model in this Bi-LSTM NER are maximum length = 20, x = word idx, and y = tag idx. The model's output had a 0.9882 accuracy and a 0.0398 value loss.

VII.CONCLUSION

For several downstream NLP operations, NER is a required activity. It significantly affects these NLP tasks. This project's goal is to conduct out NER. In actuality, it can be interpreted as a sequence labelling issue. In this study, we trained and evaluated four alternative approaches to determine which was best. These models are based on deep neural networks and statistical techniques. Following the outcomes analysis, we outlined this project's key contribution.

Conclusion and Next Steps This research project intends to develop a novel deep learning model to extract named items from the supplied text, such as person name, location, organisation, and date. The results produced by the models used are satisfactory. The most well-known model has an accuracy of 0.9882 and is called the Bi-LSTM deep learning model. This effort to recognise named entities has yielded positive results.

Through an experiment, this team created four models that compared NER's deep learning model to its machine learning model counterparts. The deep learning model outperformed the statistical-based approach in our experiment. Additionally, we examined the models and talked about their benefits and drawbacks.

Our future work will involve the following:

- To handle incredibly imbalanced data, such as the BTC corpus, use focused loss.

- Examine the effect of various window sizes on the CNN layer.
- Investigate the best unsupervised learning technique for NER.

REFERENCES

1. Aroonmanakun, W., Nupairoj, N., Muangsin, V., & Choemprayong, S. (2018). Thai Monitor Corpus: Challenges and Contribution to Thai NLP. *Vacana*, 6(2), 1-14.
2. Baldwin, T., de Marneffe, M. C., Han, B., Kim, Y. B., Ritter, A., & Xu, W. (2015). Shared tasks of the 2015 workshop on noisy user-generated text: Twitter lexical normalization and named entity recognition. In *Proceedings of the Workshop on Noisy User-generated Text* (pp. 126-135).
3. Bharadwaj, A., Mortensen, D., Dyer, C., & Carbonell, J. (2016). Phonologically aware neural model for named entity recognition in low resource transfer settings. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing* (pp. 1462-1472).
4. Bonadiman, D., Severyn, A., & Moschitti, A. (2015). Deep neural networks for named entity recognition in Italian. *CLiC it*, 51.
5. Chiu, J. P., & Nichols, E. (2016). Named entity recognition with bidirectional LSTM CNNs. *Transactions of the Association for Computational Linguistics*, 4, 357-370.
6. Choi, Y., & Cha, J. (2016). Korean Named Entity Recognition and Classification using Word Embedding Features. *Journal of KIISE*, 43(6), 678-685.
7. Chrupala, G. (2011). Efficient induction of probabilistic word classes with LDA. In *Proceedings of 5th International Joint Conference on Natural Language Processing* (pp. 363-372).
8. Keerthi, S. S., & Lin, C. J. (2003). Asymptotic behaviors of support vector machines with Gaussian kernel. *Neural computation*, 15(7), 1667-1689.
9. Narayanaswamy, M., Ravikumar, K. E., & Vijay-Shanker, K. (2002). A biological named entity recognizer. In *Biocomputing 2003* (pp. 427-438).
10. Ramachandran, P., Liu, P. J., & Le, Q. V. (2016). Unsupervised pretraining for sequence to sequence learning. *arXiv preprint arXiv:1611.02683*.

TEXT BOOKS

- Name Entitree -By Satoshi Sekine, published by John Benjamins Publishing Co in 2009
- Named Entities for Computational Linguistics -By Damien Nouvel, Maud Ehrmann, Sophie Rosset, Published by John Wiley & Sons in 2016
- Natural Language Processing for the Semantic Web -By Diana Maynard and Kalina Bontcheva, Published by Morgan & Claypool Publishers in 2016