# L🏃WNICS

**Documentation**

**Lawnics: Search Engine**

**Overview:**

Lawnics helps the next generation lawyers by transforming their work over digital workspace by building AI assistants that help lawyers, organizations, government and Judiciary to make their services faster, smoother and cost effective. Consequently, transforming 300 year old work practices for the digital age. Lawnics are a cloud based legal operating software that leverages artificial intelligence, intuitive designs, and smart processes to make legal practice more meaningful.

**Getting started:**

The legal data is available in the form of HTML and JSON files. Cleaning the data as required and extracting proper entities is the first task. In this document, the term *entity/entities* will be referred to as *Acts, Articles, Sections, Cases, and Citations*. The term, *tag* and *label* will be interchangeably used for the 'B-'(beginning) and 'I-'(inside) tags as required by the BERT model. For more details on the BERT model and the required input format, refer to these links:

1. [Original paper](#)

2. [NER with BERT in Spark-NLP](#)

3. [Data preparation in CoNLL format](#)

Here is an example of how the raw data(HTML page) looks:



**Figure 1.0 : Sample HTML file**

## Requirements:

1. Python 3.7.4

2. Jupyter notebook/google collab

   Main Packages:

   1. *BeautifulSoup* (for HTML handling)

   2. *Pandas* (for data manipulation)

   3. *nltk* (for tokenizing document)

   4. *re* (for regex and matching patterns)

   5. *json* (for annotated file conversion to CoNLL)

**Foundation of Web Scraping:**

The hyperlinks in the HTML file allowed for extracting the required entities.

1. As you can see in figure 1.0, the hyper linked entities are: *Article 137, Limitation Act, Mohd. Abdul Khader Mohd. Kastim and Anr. v. Pareethij Kunju Sayed A, Section 58(d).*

2. The raw data from the HTML was cleaned and special characters, punctuations, and any other characters that were not classified as words were also removed.

3. The tags were extracted using the BeautifulSoup package, mentioned in '*Requirements',* by extracting all the paragraph('p') tags and associated 'a' tags for the links. A python dictionary would then be formed consisting of all the entities:

    {'Article 137': 'O',

    'Limitation Act': 'O',

    'Mohd. Abdul Khader Mohd. Kastim and Anr. v. Pareethij Kunju Sayed A': 'O',

    'Section 58(d)': 'O',

    And other entities in the HTML in a similar way…….. }

4. After this, a python dictionary of entities was created by tagging the entities appropriately using keywords. For example, any *Article* entity always has the word 'Article' in it, the same with *Sections* and *Acts*. *Citations* and *Cases* were handled differently. For cases, the keyword "vs." or "v" was used to find the *Case* entity. For citations, regex was used to identify each type of citation from the Journal: SCC, AIR, SCR, Crime, ITR, ACR, SCALE, CS, ELT, ILR, JT. The resultant dictionary had the entity name as the key and the start index(of that entity in the original raw data), end index(of that entity in the original raw data), and the type of entity as the value. For example:

a. {'citations': [[23263, 23280, 'Cite'],   [34262, 34277, 'Cite'],   [34280, 34297, 'Cite']]}

5. Next, using the start and end index of a particular entity in the original raw data, the first word of the entity was tagged as "B-(name of entity)" tag and the rest of the entity words till the end index were tagged as "I-(name of entity)" tags. The non-entity words were tagged as 'O' (other meaning) tags.

6. This was then converted to a dataframe with two columns: token name and respective tags.

This was the foundation of the data scraping model.

For BERT training, POS tags were also added in the source training file using NER.

The resultant data-frame was then converted to tsv format with approximately size of 25MB outside of the data scraping code in google collab (refer to the links in '*Getting Started'* ). This was the final training dataset used for version 1.0.

**Report of the first BERT training 1.0:**

1. Google collab was used for training since spark package requirements failed to work on local systems.

   **Requirements for training and testing:**

   1. java-8
   2. pyspark 2.4.4
   3. spark-nlp-2.7.4
   4. sklearn
   5. pre-trained variant of BERT: *'bert_base_cased'*

2. While training, the training dataset was not being read by the 'CoNLL' method of the 'sparknlp.training' package. The error was as mentioned in figure 1.1. This issue was raised on GitHub.



**Figure 1.1 : Incorrect Input Error**

3. This issue was resolved by updating the input dataset and correcting minor errors according to the CoNLL 2003 format using this link.

4. After the training data was in the required CoNLL format, a RAM issue was encountered in the same 'CoNLL' method of the 'sparknlp.training' package. The error was as mentioned in Figure 1.2.

```
ERROR:py4j.java_gateway:An error occurred while trying to connect to the Java server (127.0.0.1:39593)
Traceback (most recent call last):
  File "/usr/local/lib/python3.6/dist-packages/py4j/java_gateway.py", line 929, in _get_connection
    connection = self.deque.pop()
IndexError: pop from an empty deque

During handling of the above exception, another exception occurred:

Traceback (most recent call last):
  File "/usr/local/lib/python3.6/dist-packages/py4j/java_gateway.py", line 1067, in start
    self.socket.connect((self.address, self.port))
ConnectionRefusedError: [Errno 111] Connection refused
---------------------------------------------------------------------------
IndexError                                Traceback (most recent call last)
/usr/local/lib/python3.6/dist-packages/py4j/java_gateway.py in _get_connection(self)
    928         try:
--> 929             connection = self.deque.pop()
    930         except IndexError:

IndexError: pop from an empty deque

During handling of the above exception, another exception occurred:

ConnectionRefusedError                    Traceback (most recent call last)
                         ⇕ 8 frames
ConnectionRefusedError: [Errno 111] Connection refused

During handling of the above exception, another exception occurred:

Py4JNetworkError                          Traceback (most recent call last)
/usr/local/lib/python3.6/dist-packages/py4j/java_gateway.py in start(self)
   1077                 "server ({0}:{1})".format(self.address, self.port)
   1078             logger.exception(msg)
-> 1079             raise Py4JNetworkError(msg, e)
   1080
```

**Figure 1.2 : RAM issue**

5. This issue was then resolved by reducing the input data size and upgrading RAM in collab from 12GB to 25GB. Optimal file size was found to be a dataset of 40 case files.

6. Class imbalances were found as well with "O-" being a lot more than the entity tags.

7. This class imbalance issue was resolved after extracting paragraphs with entities and discarding the ones with no entities(hyperlinks). Each sentence in CoNLL is one para.

8. Pipeline fitting took approximately 45 minutes:

*ner_model = ner_pipeline.fit(training_data)*

9. Model was saved in collab and later downloaded as separate files in the local system to save time and effort in pipeline fitting.

10. *ner_prediction_pipeline = Pipeline(*
    *stages = [*
        *document,*
        *sentence,*
        *token,*
        *bert,*
        *loaded_ner_model,*
        *converter])*

11. Model was very less accurate and was not able to detect many of the entities. This was tried manually by giving some text input in the Light Pipeline to check for prediction.

    *ner_pipelineFit = ner_prediction_pipeline.fit(empty_data)*

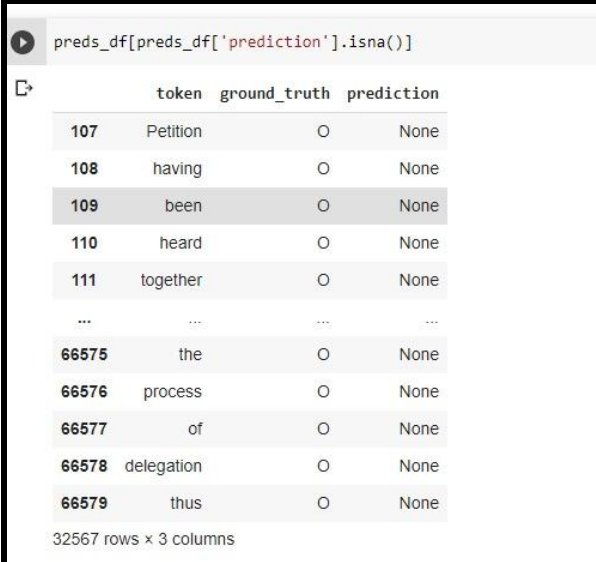    *ner_lp_pipeline = LightPipeline(ner_pipelineFit)*

    *parsed = ner_lp_pipeline.annotate(txt)*

    The *'txt'* is the text file or string passed to check for prediction.

12. During testing, due to null value prediction as seen in figure 1.3, the classification report could not be generated.

13. In figure 1.4, the classification report is shown after removing the null values.



**Figure 1.3**

```
print(classification_report(pred_df2["ground_truth"], pred_df2["prediction"]))

              precision    recall  f1-score   support

       B-Act       0.60      0.35      0.44        96
   B-Article       1.00      0.96      0.98        53
     B-Cases       0.71      0.56      0.63        71
      B-Cite       0.79      0.90      0.84        29
     B-Other       0.80      0.11      0.19        38
   B-Section       0.92      0.94      0.93       327
       I-Act       0.58      0.58      0.58       146
   I-Article       0.99      0.95      0.97        78
     I-Cases       0.74      0.96      0.84       494
      I-Cite       0.91      0.86      0.88        92
     I-Other       0.10      0.05      0.07        20
   I-Section       0.90      0.86      0.88       412
           O       0.99      0.99      0.99     32253

    accuracy                           0.98     34109
   macro avg       0.77      0.70      0.71     34109
weighted avg       0.98      0.98      0.98     34109
```

**Figure 1.4 : Classification report**

**Report of BERT training 2.0:**

1. Google collab was used for training here as well.

2. Same requirements as BERT training 1.0.

3. The problem in the earlier version was due to input data:

   a. Null values were due to larger sentences which confused the BERT model.

   b. Punctuations were removed due to which accuracy decreased when tested on real unclean data with punctuations.

4. Changes in approach:

   a. Changes were made to document cleaning.

   b. Changes were made to create_entities() method.

   c. Punctuations were left in after cleaning.

   d. For entity extraction, it should be extracted with its respective punctuation. For example: In figure 1.0, 'Limitation Act, 1963' is hyperlinked till the word 'Act'. By version 1.0, 1963 would not be tagged as a part of the Act even though it was.

For this issue, pattern matching was added for Acts and Citations. The new entity tags for *'Limitation Act, 1963'* and *'(2015) 7 SCC 58'* were:

    i.    ['Limitation', 'B-Act'] , ['Act', 'I-Act'], [',', 'I-Act'], ['1963', 'I-Act']

    ii.    ['(', 'B-Cite'], ['2015', 'I-Cite'], [')', 'I-Cite'], ['7', 'I-Cite'], ['SCC', 'I-Cite'] , ['58', 'I-Cite']

    iii.    The same format was followed with sections and other entities.

e.   Sentences were split using:

    i.    *'PunktSentenceTokenizer'* and *'PunktTrainer'* of the nltk package.

    ii.    For sentence splitting, the sentence should not start with any abbreviation words and not split any collocations. Due to this, the *PunktSentenceTokenizer* was updated with some common legal words and abbreviations.

5. Data was generated in the CoNLL format in the original script itself instead of exporting it as a csv and then converting it to tsv and finally to CoNLL.

6. Error occurred during loading the input files in the 'CoNLL' method of the 'sparknlp.training' package.

7. This error was tried by, first, reducing the input data. This did not work. Upon further investigation, the cause of the error was found as special characters.

8. Unix commands were used to remove the special characters found in the CoNLL file. FYI: Only special characters(Non-ASCII characters) were removed which caused errors in the 'CoNLL' method of the 'sparknlp.training' package during loading input with punctuations being intact.

9.  The following unix commands were used:

    a.  bash-3.2$ LC_ALL=C tr -dc '\0-\177' <hasSpecialChars.txt>tempFile.txt

        bash-3.2$ tr -cd '\11\12\15\40-\176'<tempFile.txt>doesNotHaveSpecialChars.txt

10. After the above changes, the input was accepted in the CoNLL method.

11. There was an error during pipeline fitting as shown in figure 2.0.



**Figure 2.0 : Memory error**

12. This was resolved by reducing the input size to 150 files.

13. Pipeline was fit and the model was trained.

14. Null value predictions still persisted.

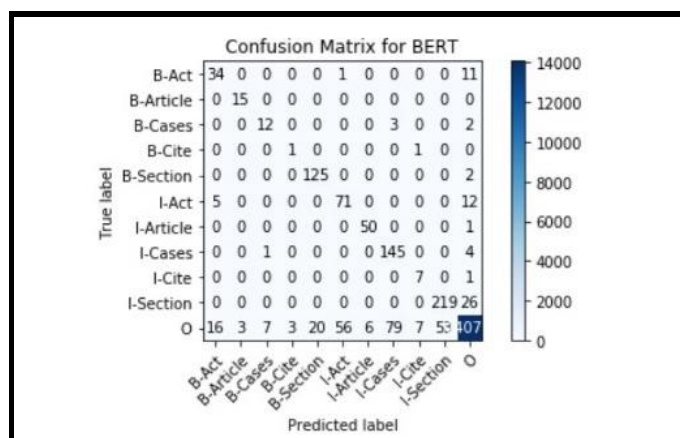15. In figure 2.1, the confusion matrix is shown for this version.

**Figure 2.1 : Confusion Matrix**

**Report of BERT training 3.0 (IN PROGRESS):**

**Issues:**

1.  The sentence splitting was not appropriate and created longer than 110 words sentences which confused BERT while training.

2.  The rule based annotations were not efficiently tagging all the entities in the file since some of them were not hyperlinked.

3.  When converting from JSON to CoNLL, there exists an issue of labels being converted appropriately.

**Resolve:**

1.  For the first issue, sentence segmentation models are being tried and tested. (In progress)

2.  For the second issue, the Doccano annotation tool is being used.

**About Doccano:**

1.  The input should be a CoNLL file with no pre-existing tags since it creates overlapping.

2.  Doccano does not accept special characters.

3. Uploading CoNLL files was done in batches. Some files in batches were not uploaded since there were special characters in them. Figure 3.0 shows the input format.

4. Some special characters had to be replaced manually and most of them were removed using the unix commands mentioned in version 2.0 section 9.

5. The output of the annotated file is JSON-L and JSON-text tables.

6. Python script was written to convert JSON to CoNLL format(accepted by BERT).

```
Act
XI
of
1125
M.E
.
)
and
the
Travancore-Cochin
General
Sales
Tax
Rules
,
1950
,
made
thereunder
which
may
be
conveniently
set
out
here
.
The
preamble
```

**Figure 3.0 : with each token on a different line**