

Overview of the FIRE 2020 AILA Track: Artificial Intelligence for Legal Assistance

Paheli Bhattacharya^a, Parth Mehta^b, Kripabandhu Ghosh^c, Saptarshi Ghosh^a, Arindam Pal^d, Arnab Bhattacharya^e and Prasenjit Majumder^f

^aIndian Institute of Technology Kharagpur, India

^bParmonic AI

^cIndian Institute of Science Education and Research (IISER), Kolkata, India

^dData61, CSIRO and Cyber Security CRC, Sydney, New South Wales, Australia

^eIndian Institute of Technology Kanpur, India

^fDA-IICT Gandhinagar, India

Abstract

The FIRE 2020 AILA track focused on two tasks – (i) Retrieving relevant Prior cases and Statutes given a factual description, and (ii) Rhetorical labelling of sentences in a legal case document, where the rhetorical roles are – Facts of the case, Ruling by the Lower Court, Argument, Statute, Precedent, Ratio of the decision and Ruling by the Present Court. Both the tasks were based on publicly available case documents from the Indian Supreme Court judiciary.

Keywords

Legal data analytics, Prior case retrieval, Statute retrieval, Legal facts, Rhetorical role labelling, Legal IR, Legal NLP

1. Introduction

Common Law system, which is followed by most countries (e.g., UK, USA, Canada, Australia, India etc.) has two primary sources – Precedents and Statutes. Precedents are the prior cases decided in the Courts of law. Statutes are bodies of written law, such as the Constitution of a country.

When a case is presented to a lawyer / law practitioner in terms of a factual description, the law practitioner has to search for relevant Precedent cases and Statutes, in order to prepare legal reasonings accordingly. Since these documents are large in number, it will be beneficial for a law practitioner to have an automated tool that can return prior-cases and statutes relevant to a factual scenario. Motivated by this, we design Task 1: Precedent and Statute Retrieval for a given factual description. This task is a continuation of a task that we had in AILA-2019 [1].

FIRE 2020: Forum for Information Retrieval Evaluation, December 16-20, 2020, Hyderabad, India

✉ paheli@cse.iitkgp.ac.in (P. Bhattacharya); parth.mehta126@gmail.com (P. Mehta); kripa.ghosh@gmail.com (K. Ghosh); saptarshi@cse.iitkgp.ac.in> (S. Ghosh); arindamp@gmail.com (A. Pal); arnabb@cse.iitk.ac.in (A. Bhattacharya); prasenjit.majumdar@gmail.com (P. Majumder)



© 2020 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

Also, a case document from the Indian judiciary is usually very long and unstructured, lacking section and paragraph headings. This makes it difficult for a reader of the document to identify where the facts of the case are written, which sentences mention the arguments made in the case, what was the final judgement and so on. To address this research problem, we introduce Task 2: Rhetorical Role Labeling for Legal Judgements.

AILA 2020 witnessed a participation of 15 teams with 14 of them submitting the working notes. We received a total of 74 runs across the three tasks. Apart from teams from India, we had teams from China, Botswana, Italy, Austria and Canada. Also we had teams from both academic institutions as well as the industry.

1.1. Task 1: Precedent and Statute Retrieval

In this task, the aim is to retrieve relevant prior-cases and statutes for a given factual scenario, from a large pool of prior-case documents and statutes. Further details of the task are given in the overview paper of AILA-2019 [1].

For training, we provide the dataset of AILA-2019 [1]. For the test/evaluation data, we provide a set of additional 10 queries, that describe factual scenarios in natural English language. The pool / candidate statutes from which the relevant ones are to be retrieved, is the same as that in AILA-2019 dataset [1]. A set of 343 documents were added to the existing pool of prior-case documents in AILA-2019 [1]. Out of these 343 documents, 43 documents were the relevant prior cases for the 10 queries in the test set. The remaining 300 documents were sampled based on the text similarity (cosine similarity between the document vectors were in $[0.3, 0.5]$) between the existing documents in the prior-case pool. Hence, for the present task, the resultant number of prior-case documents is 3,257 (from which, prior cases relevant to a given query need to be retrieved).

Also we provide a set of 197 statutes (Sections of Acts) from Indian law, that are relevant to some of the queries stated above. We provide the participants with the title and description of these statutes.

For each query, the task was to retrieve the most similar/relevant precedent case documents and statutes with respect to the situation in the given query.

Similar research has been done on Chinese legal case documents [2, 3] that dealt with the task of retrieving statutes for a given fact. Note that, in addition to retrieving statutes, we also consider retrieving prior-cases for the query. Also, constructing a large dataset for the task to train supervised models as in [2, 3] is not practical in the context of Indian legal documents. This is because, unlike Chinese legal documents, Indian legal documents are not well-structured and it is difficult to extract the facts of the cases automatically [4] for creating the dataset.

1.2. Task 2: Rhetorical Role Labeling for Legal Judgements

Since Indian legal case documents are unstructured, there is a need to design systems that can automatically segment these documents into coherent, meaningful parts. This can not only enhance the readability of the documents but also has applications in downstream tasks such as summarization, case-law analysis, semantic search and so on. We introduce the task of rhetorical role labeling of sentences in legal case judgements in AILA-2020. The task is to

assign one of the following labels to each sentence in a legal case document. We consider the following seven (07) rhetorical labels/semantic segments [5]:

- Facts: legal situation that led to filing the case
- Ruling by Lower Court: since we consider documents from the Supreme Court of India, there was some preliminary ruling given at the lower courts, e.g., High Court, Tribunal, etc.
- Arguments: arguments made by the contending parties
- Precedents: citation to relevant prior cases
- Statutes: citation to relevant statutes
- Ratio of the decision: reasoning behind the final judgement
- Ruling by Present Court: final judgement given by the Supreme Court of India

A state-of-the-art model that addresses this task for Indian legal documents is [5], where a neural model is used for the segmentation.

2. Dataset

We consider case documents from the Supreme Court of India and Statutes from the Indian judiciary.

Task 1: The training dataset for Task 1 was the AILA-2019 dataset [1] (available at <https://github.com/Law-AI/aila-2019-dataset>). There were 50 queries and 197 statutes.

For the present task (AILA-2020), the pool of prior-cases was extended to having 3,257 documents, as mentioned in Section 1.1. The test dataset of 10 queries was created in the same way as described in [1].

Task 2: The dataset made publicly available by [5] was used as the training dataset¹. There were 50 documents containing 9,308 sentences in total across all the documents. The rhetorical labels were assigned by law experts from a reputed law school in India.

As the test set, we consider a set of 10 additional case documents. We randomly selected 2 documents from each of the 5 law domains mentioned in [5]. These documents were then given to a law expert for annotating every sentence with one of the rhetorical labels. There are a total of 1,905 sentences in the test set.

3. Evaluation

For both the Tasks, evaluation was done on the test dataset. For Task 1, the same evaluation metrics as in [1] were used – Mean Average Precision (MAP), Precision@10 (P10), BPREF and

¹<https://github.com/Law-AI/semantic-segmentation>

Reciprocal rank (recip_rank). The *trec_eval* tool² was used for computing the metrics stated above. We choose MAP as the primary measure since it incorporates both Precision and Recall.

For Task 2, we use the standard Recall, Precision and F1-Scores. The documents have a considerable variation in their size. Moreover, even within a document, there is a class imbalance among the 7 categories / rhetorical roles. Hence we use macro-averaging at both document-level and category-level. The scores were calculated as below:

1. Recall, Precision and F-score were computed for each category of labels within each document.
2. The score for each document in a run were computed by averaging the scores for all seven categories in that document.
3. Finally, the overall scores for a run are computed by averaging the scores for each document.

For task-2 we additionally report the overall Accuracy for each submitted run. Accuracy is micro-averaged across documents and classes, i.e., it is measured as the fraction of sentences correctly classified out of the 1,905 test sentences.

4. Methodologies for Task 1: Precedent Retrieval and Statute Retrieval

For the first task of retrieving relevant prior/ precedent cases (Task 1a) , we received a total of 26 runs from 10 participating teams. For the second task of retrieving relevant statutes (Task 1b) , we received a total of 27 runs from 12 participating teams. The comparative results are in Table 1 (for Task 1a) and Table 2 (for Task 1b). We briefly describe below the methodologies used by each team in each of their runs. Details can be found in the working notes of the respective submissions.

Task 1a : Precedent Retrieval : We briefly describe the methodologies submitted by the various teams for the task:

- **UB [6]** : The team was from the University of Botswana. In their first submitted run UB-1, they weighed the terms in the query and documents using TF-IDF. In their second run (UB-2), they extracted key concepts and used TF-IDF for retrieval. The best performance in terms of MAP, BPREF and recip_rank for the task was by their third run, UB-3, where they use Terrier 4.2 KL divergence model.
- **double_liu_2020 [7]** : This team is affiliated to the Heilongjiang Institute of Technology, China. They extract the top 50% of the words based on their IDF scores as the search keywords in their first and third runs. In the second run, they used all the words as search keywords. In terms of MAP, BPREF in Task 1a, their third run (double_liu_2020_3) performed the second best.

²https://trec.nist.gov/trec_eval/

Table 1

Results of Task 1a: Precedent retrieval for queries. All measures averaged over 10 test queries. Numbers in **bold** and underline indicate the best and the second-best performing methods corresponding to the evaluation metrics. Rows are sorted in decreasing order of MAP (primary measure).

Team Name	Run_ID	MAP	BPREF	recip_rank	P @ 10	Method Used
UB	UB-3	0.1573	0.1128	0.238	0.08	Terrier 4.2 KL divergence model
double_liu_2020	double_liu_2020_3	0.1382	0.1045	0.1886	0.07	IDF, BM25 as the search score.
fs_hu	fs_hu_task1a	0.1351	0.0885	0.2041	0.1	Language model, Dirichlet Smoothing
double_liu_2020	double_liu_2020_1	0.1306	0.0737	0.1963	0.07	IDF, BM25
TUW_informatics	basic	0.1294	0.0737	0.1915	0.07	Preprocessing, BM25
fs_hit_1	fs_hit_1_task1a_01	0.1294	0.0877	0.1876	0.07	BM25
LAWNICS	LAWNICS_2	0.1288	0.0913	0.1586	0.1	Topic Embedding
TUW_informatics	word_count	0.1271	0.0728	0.1891	0.06	Preprocessing, BM 25
SSNCSE_NLP	task_1a_1	0.1264	0.0918	0.2043	0.08	BM 25
fs_hit_2	fs_hit_2_task1a_01	0.125	0.0724	0.1906	0.07	Language modelling of Indri
double_liu_2020	double_liu_2020_2	0.123	0.0621	0.1969	0.08	All words, BM 25
UB	UB-1	0.1229	0.07	0.2033	0.09	TF-IDF term weighting
fs_hit_1	fs_hit_1_task1a_02	0.1215	0.0699	<u>0.2078</u>	<u>0.09</u>	BM25, TF-IDF
UB	UB-2	0.1168	0.0798	0.1967	0.07	Extract key concepts, TF-IDF
TUW_informatics	false_friends	0.1133	0.0687	0.1873	0.05	Preprocessing, BM25
LAWNICS	LAWNICS_1	0.1085	0.0756	0.1607	0.08	Preprocessing, BM25
Uottawa_NLP	run3_TFIDF	0.0837	0.0399	0.1157	0.05	preprocessing, TFIDF
fs_hit_1	fs_hit_1_task1a_03	0.0696	0.0267	0.1088	0.07	Cosine Similarity
SSNCSE_NLP	task_1a_2	0.0652	0.0406	0.1004	0.05	TF IDF
IMS_UNIPD	tfidf_lemma	0.0575	0.0324	0.1068	0.02	TF IDF, lemma words
IMS_UNIPD	tfidf_stem	0.056	0.0341	0.1077	0.03	TF IDF, stemmed words
IMS_UNIPD	bm25_lemma	0.0441	0.0143	0.147	0.03	BM 25, lemma forms
fs_hit_2	fs_hit_2_task1a_02	0.0126	0	0.041	0.02	Lucene, TFIDF
Uottawa_NLP	run1_Glove	0.0123	0	0.0222	0	Glove
fs_hit_2	fs_hit_2_task1a_03	0.0123	0	0.0395	0.02	Lucene, Dirichlet Similarity
Uottawa_NLP	run2_Doc2Vec	0.0029	0	0.0051	0	Doc2Vec

- **fs_hu** [8] : This team from the Foshan University, China used language model and Dirichlet Smoothing for retrieval. They performed the best in terms of P@10.
- **TUW_Informatics** [9]: This team from TU Wien experiment with different stopword lists. In the first run, basic, preprocessing is performed on the documents and retrieval is through BM-25 algorithm. In their next two runs, word_count and false_friends, they experiment with different methods for stopword removal as the preprocessing step.
- **SSNCSE_NLP** [10]: This team is from Sri Siva Subramaniya Nadar College of Engineering, India. They use BM-25 and TF-IDF for the task.
- **fs_hit_1** [11] : This team is from the Foshan University, China. They used BM25 and TF-IDF similarity in their first and second runs. Their second run was the second best performing method in terms of recip_rank and P@10.
- **fs_hit_2** [12] : This team from the Foshan University, China and Heilongjiang Institute of Technology, China, explored different Language models for the task. In their first run, they use the language model assorting algorithm of Indri. In the second run, language model assorting algorithm of Lucene is used. In the third run, they use language Model with Dirichlet Similarity.

Table 2

Results of Task 1b: Statute retrieval for queries. All measures averaged over 10 test queries. Numbers in **bold** and underline indicate the best and the second-best performing methods corresponding to the evaluation metrics. Rows are sorted in decreasing order of MAP (primary measure).

Team Name	Run_ID	MAP	BPREF	recip_rank	P @ 10	Method Used
scnu	scnu_1	0.3851	0.3054	0.5615	0.18	BERT
SSNCSE_NLP	task_1b_2	<u>0.3423</u>	0.136	0.3423	0.07	TF IDF
IMS_UNIPD	tfidf_stem	0.3383	0.279	0.5349	0.17	TF IDF, stemmed words
IMS_UNIPD	tfidf_lemma	0.3159	0.2568	0.5288	<u>0.17</u>	TF IDF, lemma forms
UB	UB-2	0.3134	0.2633	0.5787	0.15	Extract key concepts, TF-IDF
UB	UB-1	0.3085	0.2633	0.573	0.14	TF-IDF
SSN_NLP	R1	0.2975	0.2531	0.4769	0.15	BM 25
LAWNICS	LAWNICS_1	0.2962	<u>0.2812</u>	0.4607	0.13	Preprocessing, BM25
LAWNICS	LAWNICS_2	0.2962	0.2812	0.4607	0.13	Law2Vec embeddings
TUW_informatics	basic	0.2619	0.2033	0.4855	0.13	Preprocessing, BM25
TUW_informatics	word_count	0.2574	0.214	0.3946	0.14	Preprocessing, BM 25
Uottawa_NLP	run3_TFIDF	0.2506	0.186	0.3144	0.12	Preprocessing, TFIDF
fs_hu	fs_hu_task1b	0.235	0.198	0.3581	0.08	TF-IDF, Jaccard
TUW_informatics	false_friends	0.2316	0.1855	0.3814	0.1	Preprocessing, BM 25
IMS_UNIPD	bm25_lemma	0.231	0.1586	0.4595	0.15	BM 25 on lemma forms
fs_hit_1	fs_hit_1_task1b_03	0.2139	0.1587	0.3371	0.13	Language Model with Dirichlet smoothing
fs_hit_2	fs_hit_2_task1b_01	0.2003	0.1587	0.3452	0.1	Language Model
fs_hit_2	fs_hit_2_task1b_03	0.1886	0.132	0.279	0.1	Language Model
UB	UB-3	0.1876	0.1502	0.2468	0.09	Terrier 4.2 KL divergence model
fs_hit_2	fs_hit_2_task1b_02	0.1777	0.1247	0.2546	0.12	Language Model
fs_hit_1	fs_hit_1_task1b_01	0.1703	0.0945	0.2196	0.12	Language Model with JM smoothing
fs_hit_1	fs_hit_1_task1b_02	0.1703	0.0945	0.2196	0.12	Language Model with JM smoothing
Uottawa_NLP	run1_Glove	0.1462	0.084	0.345	0.1	preprocessing, Glove
scnu	scnu_3	0.1301	0.048	0.1531	0.12	Chen at.al., Enhanced LSTM for Natural Language Inference, ACL 2017
SSNCSE_NLP	task_1b_1	0.1181	0.069	0.2739	0.07	BM 25
nlpninjas	nlpninjas_st1	0.0917	0.024	0.1204	0.07	n-gram, BM25
Uottawa_NLP	run2_Doc2Vec	0.0441	0.008	0.067	0.02	Doc2Vec
scnu	scnu_2	0.0254	0	0.0203	0	Xiong et.al., End-to-End Neural Ad-hoc Ranking with Kernel Pooling, SIGIR 2017

- **IMS_UNIPD** [13]: This team is affiliated to the Information Management System (IMS) Group, University of Padua. They used BM-25 on the lemma forms of the words in the query and candidate case documents for the bm25_lemma run; used TF-IDF weighting on the lemma forms in tfidf_lemma run and TF-IDF weighting on the stemmed words in the tfidf_stem run.
- **Lawnics** [14]: This team is from Lawnics Technologies, India. They use BM25 and topic embedding methods for the precedent retrieval task, for their first and second runs respectively.
- **Uottawa_NLP** [15]: This team is from the University of Ottawa. After preprocessing the documents, they used Glove, Doc2Vec and TF-IDF based methods, for the first, second and third runs respectively.

Task 1b : Statute Retrieval : We describe the methodologies submitted by different teams for this task in brief. See Table 2 for a comparison among the performances of the methodologies.

- **scnu**³ : This team has its members from the South China Normal University. It was

³Working note not submitted.

the only team that modelled the task as a supervised task and performed training using the training dataset provided. Their first run that uses BERT is the best performing method in the statute retrieval task in terms of MAP, BPREF and P@10. They are second best in terms of recip_rank. They have also experimented with different supervised methods in their second and third runs.

- **SSNCSE_NLP** [10]: This team is from Sri Siva Subramaniya Nadar College of Engineering. They use BM-25 and TF-IDF for the task. They get the second best MAP scores for the task.
- **IMS_UNIPD** [13]: This team is affiliated to the Information Management System (IMS) Group, University of Padua. They used BM-25 on the lemma forms of the words in the query and candidate case documents for the bm25_lemma run; used TF-IDF weighting on the lemma forms in tfidf_lemma run and TF-IDF weighting on the stemmed words in the tfidf_stem run. They perform second best in terms of P@10.
- **UB** [6] : The team was from the University of Botswana. In their first submitted run they weighed the terms in the query and documents using Tf-IDF. In the second run, they extracted key concepts and used TF-IDF for retrieval. This method achieves the best recip_rank for the task. In their third run they use Terrier 4.2 KL divergence model.
- **SSN_NLP** [16]: The team has its members from the SSN College Of Engineering. They use BM-25 for the statute retrieval task.
- **Lawnics** [14]: This team is from Lawnics Technologies, India. They use BM25 in their first run. They explore Law2Vec embeddings in their second run. They achieve second best BPREF scores for the task.
- **TUW_Informatics** [9]: This team from TU Wien experiment with different stopword lists. In the first run, basic, preprocessing is performed on the documents and retrieval is through BM-25 algorithm. In their next two runs, word_count and false_friends, they experiment with different methods for stopword removal as the preprocessing step.
- **Uottawa_NLP** [15]: This team is from the University of Ottawa. After preprocessing the documents, they used Glove, Doc2Vec and TF-IDF based methods, for the first, second and third runs respectively.
- **fs_hu** [8] : This team from the Foshan University, China used TF-IDF and Jaccard to compute similarity.
- **fs_hit_1** [11] : This team is from the Foshan University, China. They experiment with different Language Models with Dirichlet smoothing and JM Smoothing with different hyper-parameters.
- **fs_hit_2** [12] : This team, also from the Foshan University, China, explored different Language models for the task tuned with different hyper-parameters.

- **nlpninjas**⁴ : This team is from Deloitte USI. They combined unigrams and bigrams. They used BM25 for retrieval.

For Task 1a (Precedent retrieval) the best performing run achieves a MAP of 0.1573. The methods submitted were all unsupervised in nature. The training data provided was mainly used to tune the hyperparameters of the retrieval models (eg. BM-25). Embedding methods like Doc2Vec, Glove were also used but the performance was not good. This is probably because these methods require a huge amount of data to be trained on.

For Task 1b (Statute retrieval) the best performing run achieved a MAP of 0.3851. The method used BERT for extracting deep semantic features. Law2Vec embeddings have also been explored for the task. Only 1 team (scnu) used training data to train supervised models – BERT, LSTM, and a Neural Ranking model. Other teams used the training data to tune the hyperparameters of unsupervised methods.

5. Methodologies for Task 2: Rhetorical Role Labeling for Legal Judgements

We received 21 runs from 9 teams for the task. Table 3 compares the performance of the various runs. Brief descriptions of the methods are as follows. Details can be found in the working notes of the respective submissions.

- **ju_nlp** [17] : The team from Jadavpur University, India was the best performing team in terms of F-Score and Recall. They used a state-of-the-art transformer architecture ROBERTA along with BiLSTM for rhetorical role classification. The different runs are for different epochs of the model training.
- **heu_gjm** [18] : The team has its members from Harbin Engineering University Harbin, and Foshan University, China. They combine TF-IDF features and deep semantic features using BERT. Logistic regression, linear kernel SVM and AdaBoost are used as classifiers. The BERT model with LogisticRegression gave the best precision for the task.
- **double_liu** [7] : This team is affiliated to the Heilongjiang Institute of Technology, China. They experiment with bag-of-words based features along with SVM and Adaboost as classifiers. They also use BERT for the task. The BERT model gave the best accuracy.
- **spectre** [19] : This team from BITS Pilani, India used ROBERTA and a fully connected layer for classification.
- **LAWNICS** [14] : This team from Lawnics Technologies, India experimented with BERT for extracting features. For classification they use a fully connected layer and a max pooling layer for the different submitted runs.
- **fs_hu** [8]: The team has its members from the Foshan University, China. They have experimented with BERT using different random seeds for the runs.

⁴Working note not submitted

Table 3

Results of Task 2: Rhetorical Role Labeling for Legal Judgements. Measures averaged over 10 test documents comprising of 1,905 sentences. Numbers in **bold** and underline indicate the best and the second-best performing methods corresponding to the evaluation metrics. Rows are sorted in decreasing order of FScore (primary measure).

Team	Run	Macro Precision	Macro Recall	Macro F-Score	Accuracy	Method used
ju_nlp	ju_nlp_2	0.506	0.501	0.468	0.588	ROBERTA+Bi-LSTM
ju_nlp	ju_nlp_3	0.519	0.479	0.457	0.57	ROBERTA+Bi-LSTM
heu_gjm	heu_gjm_1	0.541	0.472	<u>0.457</u>	<u>0.603</u>	BERT+LogisticRegression
heu_gjm	heu_gjm_2	0.526	0.468	0.451	0.598	BERT+SVM
ju_nlp	ju_nlp_1	0.504	0.483	0.452	0.588	ROBERTA+Bi-LSTM
double_liu	double_liu_3	0.472	0.486	0.444	0.619	BERT
heu_gjm	heu_gjm_3	<u>0.529</u>	0.456	0.444	0.59	BERT+AdaBoost
spectre	spectre_1	0.485	<u>0.483</u>	0.442	0.584	ROBERTA+FC
LAWNICS	lawnics_2	0.479	0.479	0.435	0.584	BERT+Pooling
fs_hu	fs_hu_1	0.493	0.454	0.428	0.562	BERT+FC
fs_hit_1	fs_hit_1_3	0.484	0.449	0.41	0.574	BERT+LogisticRegression
fs_hit_2	fs_hit_2_1	0.411	0.465	0.405	0.535	TFIDF+LogisticRegression
fs_hit_1	fs_hit_1_2	0.456	0.433	0.405	0.578	BERT
fs_hit_2	fs_hit_2_2	0.455	0.427	0.398	0.549	BERT+LogisticRegression
fs_hit_1	fs_hit_1_1	0.486	0.406	0.385	0.508	TFIDF+LogisticRegression
double_liu	double_liu_2	0.423	0.407	0.355	0.488	Adaboost
SSNCSE_NLP	ssncse_nlp_2	0.384	0.4	0.354	0.46	FastText+MLP
SSNCSE_NLP	ssncse_nlp_1	0.473	0.354	0.333	0.467	TF-IDF+RF
double_liu	double_liu_1	0.432	0.351	0.327	0.469	BoW+SVM
fs_hu	fs_hu_2	0.262	0.343	0.266	0.457	BERT+FC
LAWNICS	lawnics_1	0.208	0.164	0.119	0.152	BERT+FC

- **fs_hit_1** [11]: This team is from the Foshan University, China. The first run uses the Logistic Regression with the features of TF-IDF. The second and third runs are generated by BERT with different random seeds.
- **fs_hit_2** [12]: This team from the Foshan University, China and Heilongjiang Institute of Technology, China, experiment with both TF-IDF features and BERT-based features for the task.
- **SSNCSE_NLP** [10] : This team is affiliated to the Sri Siva Subramaniya Nadar College of Engineering, India. They experiment with FastText and TF-IDF from the feature engineering aspect, and Multi-layer perceptron and Random Forest from the classifier aspect.

We find that the best performing method which achieved an FScore of 0.468 used ROBERTA which is a state-of-the-art deep learning model. We observe that BERT was applied by almost all the teams, with different classifiers (LR, FC etc.) Traditional Machine Learning approaches (SVM, Random Forest etc.) were also tried out and Bag-of-Words based feature-vector was also explored. Deep Learning methods that could extract deep semantic features were shown to perform much better than traditional feature based approaches.

6. Concluding Discussions

The FIRE 2020 AILA track has created benchmark datasets for two important tasks in the field of legal data analytics. We retained AILA 2019’s task of retrieving relevant statutes and precedents for a query. We created a new task on rhetorical role labelling of sentences in Indian legal documents. For the precedent retrieval task, we conclude from the results that it is a challenging task, mainly because of the difference in length of the query and a prior-case document. For the statute retrieval task, training BERT using very little amount of data (50 queries and their corresponding gold standard statutes), shows promising results. For the rhetorical role labelling task, participants have used state-of-the-art deep learning techniques like ROBERTA and BERT, which gives good results. In the future, we plan to extend the dataset of both the tasks.

Acknowledgements: The track organizers thank all the participants for their interest in this track. We also thank the FIRE 2020 organizers for their support in organizing the track. The research is partially supported by SERB, Government of India, through a project titled “NYAYA: A Legal Assistance System for Legal Experts and the Common Man in India”. Paheli Bhattacharya is supported by a PhD Fellowship from Tata Consultancy Services.

References

- [1] P. Bhattacharya, K. Ghosh, S. Ghosh, A. Pal, P. Mehta, A. Bhattacharya, P. Majumder, Overview of the FIRE 2019 AILA track: Artificial Intelligence for Legal Assistance, in: Proceedings of FIRE 2019 - Forum for Information Retrieval Evaluation, 2019.
- [2] P. Wang, Z. Yang, S. Niu, Y. Zhang, L. Zhang, S. Niu, Modeling dynamic pairwise attention for crime classification over legal articles, in: Proceedings of ACM SIGIR Conference on Research & Development in Information Retrieval, SIGIR ’18, 2018, pp. 485–494.
- [3] P. Wang, Y. Fan, S. Niu, Z. Yang, Y. Zhang, J. Guo, Hierarchical matching network for crime classification, in: Proceedings of ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR’19, 2019, pp. 325–334.
- [4] P. Bhattacharya, K. Hiware, S. Rajgaria, N. Pochhi, K. Ghosh, S. Ghosh, A Comparative Study of Summarization Algorithms Applied to Legal Case Judgments, in: Proc. European Conference on Information Retrieval (ECIR), 2019.
- [5] P. Bhattacharya, S. Paul, K. Ghosh, S. Ghosh, A. Wyner, Identification of rhetorical roles of sentences in indian legal judgments, in: Proc. International Conference on Legal Knowledge and Information Systems (JURIX), 2019.
- [6] T. Leburu-Dingalo, N. P. Motlogelwa, E. Thuma, M. Modungo, Ub at fire 2020 precedent and statute retrieval, in: Proceedings of FIRE 2020 - Forum for Information Retrieval Evaluation, 2020.
- [7] L. Liu, L. Liu, Z. Han, Query revaluation method for legal information retrieval, in: Proceedings of FIRE 2020 - Forum for Information Retrieval Evaluation, 2020.
- [8] Z. Li, L. Kong, Language model-based approaches for legal assistance, in: Proceedings of FIRE 2020 - Forum for Information Retrieval Evaluation, 2020.
- [9] T. Fink, G. Recski, A. Hanbury, Fire2020 aila track: Legal domain search with minimal

domain knowledge, in: Proceedings of FIRE 2020 - Forum for Information Retrieval Evaluation, 2020.

- [10] N. N. A. Balaji, B. Bharathi, J. Bhuvana, Legal information retrieval and rhetorical role labelling for legal judgements, in: Proceedings of FIRE 2020 - Forum for Information Retrieval Evaluation, 2020.
- [11] M. Wu, Z. Wu, W. Xiangyu, Z. Han, Retrieval model and classification model for aila2020, in: Proceedings of FIRE 2020 - Forum for Information Retrieval Evaluation, 2020.
- [12] Y. Xu, T. Li, Z. Han, The language model for legal retrieval and bert-based model for rhetorical role labeling for legal judgments, in: Proceedings of FIRE 2020 - Forum for Information Retrieval Evaluation, 2020.
- [13] G. M. Di Nunzio, A study on lemma vs stem for legal information retrieval using r tidyverse, in: Proceedings of FIRE 2020 - Forum for Information Retrieval Evaluation, 2020.
- [14] J. Arora, T. Patankar, A. Shah, S. Joshi, Artificial intelligence as legal research assistant, in: Proceedings of FIRE 2020 - Forum for Information Retrieval Evaluation, 2020.
- [15] I. Almuslim, D. Inkpen, Document level embeddings for identifying similar legal cases and laws (aila 2020 shared task), in: Proceedings of FIRE 2020 - Forum for Information Retrieval Evaluation, 2020.
- [16] K. S., T. D., C. Aravindan, Best matching algorithm to identify and rank the relevant statutes, in: Proceedings of FIRE 2020 - Forum for Information Retrieval Evaluation, 2020.
- [17] S. B. Majumder, D. Das, Rhetorical role labelling for legal judgements using roberta, in: Proceedings of FIRE 2020 - Forum for Information Retrieval Evaluation, 2020.
- [18] J. Gao, H. Ning, Z. Han, L. Kong, H. Qi, Legal text classification model based on text statistical features and deep semantic features, in: Proceedings of FIRE 2020 - Forum for Information Retrieval Evaluation, 2020.
- [19] R. Jain, A. Agarwal, Y. Sharma, Spectre@aila-fire2020: Supervised rhetorical role labeling for legal judgments using transformers, in: Proceedings of FIRE 2020 - Forum for Information Retrieval Evaluation, 2020.