# Natural Language Processing

## Miles to go

Prasenjit Majumder

# Natural Language

Spoken Language

Written Language

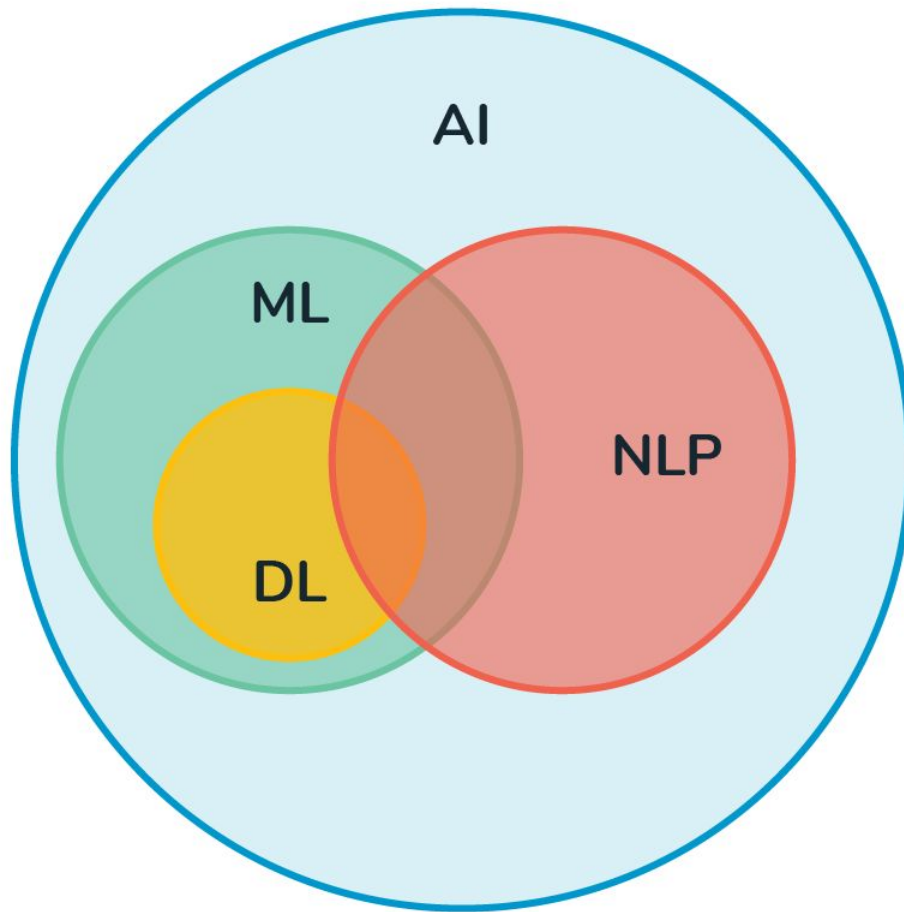Vocal Language

Sign Language

# Writing systems

- Alphabets (English)
- Logographies (Chinese, Egyptian hieroglyphs)
- Abugida (Brahmic, Tibetan etc.)

Directions:

- Top-down,
-  Left-Right,
-  Right to Left

AI

ML

NLP

DL

Artificial intelligence

Machine learning

Language Processing

Deep learning

# Natural Language Processing

- Natural Language Understanding
  - Information Retrieval
  - Summarization
- Natural Language Generation
  - Automatic Legal Drafting
  - Summarization

Word, Phrase, Sentences, Discourse

- Part of Speech
- Morphology
- Sense Disambiguation
- Entity Identification

# Morphology

# Unsupervised Root Words identification

Compute

Computer

Computing

Computerised

Computerization

# Unsupervised Root Words identification



$D_1 = \frac{1}{2^8} + \frac{1}{2^9} + \ldots + \frac{1}{2^{13}} = 0.0077$

$D_2 = \frac{1}{8} \times (\frac{1}{2^0} + \ldots + \frac{1}{2^{13-8}}) = 0.2461$

$D_3 = \frac{6}{8} \times (\frac{1}{2^0} + \ldots + \frac{1}{2^{13-8}}) = 1.4766$

$D_4 = \frac{6}{14} \times (\frac{1}{2^0} + \ldots + \frac{1}{2^{13-8}}) = 0.8438$

Edit distance = 6

$D_1 = \frac{1}{2^3} + \ldots + \frac{1}{2^9} = 0.2480$

$D_2 = \frac{1}{3} \times (\frac{1}{2^0} + \ldots + \frac{1}{2^{9-3}}) = 0.6615$

$D_3 = \frac{7}{3} \times (\frac{1}{2^0} + \ldots + \frac{1}{2^{9-3}}) = 4.6302$

$D_4 = \frac{7}{10} \times (\frac{1}{2^0} + \ldots + \frac{1}{2^{9-3}}) = 1.3891$

Edit distance = 5

# Unsupervised Root Words identification

$y_{m-1}$, but $x_m \neq y_m$).

$$D_2(X,Y) = \frac{1}{m} \times \sum_{i=m}^{n} \frac{1}{2^{i-m}} \ \text{ if } m > 0, \quad \infty \text{ otherwise}$$

$$D_3(X,Y) = \frac{n-m+1}{m} \times \sum_{i=m}^{n} \frac{1}{2^{i-m}} \ \text{ if } m > 0, \quad \infty \text{ otherwise}$$

$$D_4(X,Y) = \frac{n-m+1}{n+1} \times \sum_{i=m}^{n} \frac{1}{2^{i-m}}$$

# Unsupervised Root Words identification

Table III. Retrieval Results for Various Stemmers (WSJ, queries 151–200)

|  | No Stemming | $D_1 - 0.046$ | $D_2 - 0.31$ | $D_3 - 1.55$ | $D_4 - 0.86$ | Lovins | Porter | $n$-gram |
|---|---|---|---|---|---|---|---|---|
| Rel ret | 3082 | 3235 | 3249 | 3268 | 3265 | 3318 | 3290 | 3171 |
| $P_{20}$ | 0.4920 | 0.5020 | 0.4960 | 0.5090 | 0.5130 | 0.5030 | 0.5060 | 0.4960 |
| Avg.P | 0.3505 | 0.3732 | 0.3721 | 0.3796 | 0.3775 | 0.3746 | 0.3746 | 0.3595 |

Table VII. Performance of $D_3$-Based Stemmer on the French LeMonde Corpus

|  | No Stemming | $D_3(1.15)$ | $D_3(1.55)$ | $D_3(2.10)$ | Porter |
|---|---|---|---|---|---|
| Rel ret | 516 | 540 | 538 | 538 | 540 |
| $P_{20}$ | 0.2222 | 0.2611 | 0.2578 | 0.2522 | 0.2467 |
| Avg.P | 0.3987 | 0.4301 | 0.4334 | 0.4153 | 0.4284 |

# Parsing

# Parsing

# Dependency Parsing

Basically, we represent **dependencies as a directed graph G= (V, A)** where V(set of vertices) represents words (and punctuation marks as well) in the sentence & A( set of arcs) represent the grammar relationship between elements of V.

A dependency parse tree is the directed graph mentioned above which has the below features:

- Root has no Incoming arcs (can only be Head in Head-Dependent pair)
- Vertices(except Root) should have only one incoming arc (Only one Parent/Head)
- A Unique path should exist between Root & each vertex in the tree.

# Text representation

# Text representation

TF-IDF

Latent semantic indexing

Word2Vec

Bidirectional  Encoder Representations **from** Transformer

And many more...

# Bidirectional Encoder Representations **from** Transformer (BERT)

1. BERT (Bidirectional Encoder Representations from Transformers) uses Transformer, an attention mechanism that **"learns"** contextual relations between words (or sub-words) in a text.
2. BERT is pre-trained on two NLP tasks:
   a. Masked Language Modeling: Predict the masked word given the context words
   b. Next Sentence Prediction: Given a sentence predict the next sentence.
3. As opposed to directional models, which read the text input sequentially (left-to-right or right-to-left), the Transformer encoder reads the entire sequence of words at once.
4. BERT is pre-trained on a large corpus of unlabelled text which includes the entire Wikipedia (2,500 million words) and Book Corpus (800 million words).

Devlin, Jacob, et al. "Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805 (2018).

# BERT Architecture

# BERT Architecture



BERT<sub>BASE</sub>

BERT<sub>LARGE</sub>

# BERT Architecture

**What is the best contextualized embedding for "Help" in that context?**
For named-entity recognition task CoNLL-2003 NER



| | Dev F1 Score |
|---|---|
| First Layer | 91.0 |
| Last Hidden Layer | 94.9 |
| Sum All 12 Layers | 95.5 |
| Second-to-Last Hidden Layer | 95.6 |
| Sum Last Four Hidden | 95.9 |
| Concat Last Four Hidden | 96.1 |

# Challenges in Downstream tasks :

- Search Engines
- Hate Speech Detection
- Sentiment Analysis
- Question Answering
- Recommendation
- Summarization

# Summarization

# DATA

DUC 2002, DUC 2003 and DUC 2004

1. **DUC 2002:  59 clusters of around 10 documents each (TREC collection)**
2. **DUC 2003:  30 clusters of about 10 documents each  (TDT Datasets)**
3. **DUC 2004   50 clusters with 10 documents per cluster. (TDT Datasets)**

*\*All three datasets include four manually written summaries per cluster.*

# Effect of pre-processing and post-processing steps on ROUGE-1 recall.

| | System | No pre/post processing | Only stemming | Only stopword removal | Only redundancy removal | Stopword + Redundancy removal |
|---|---|---|---|---|---|---|
| DUC 2002 | Centroid | 0.41783 | 0.42001 | 0.42223 | 0.43157 | **0.44987** |
| | Greedy-KL | 0.40173 | 0.40537 | 0.41392 | 0.40962 | **0.41522** |
| | LexRank | 0.42733 | 0.42000 | 0.42292 | **0.44134** | 0.43289 |
| | FreqSum | 0.39247 | 0.38120 | 0.40480 | 0.38766 | **0.42522** |
| DUC 2003 | Centroid | 0.33387 | 0.34222 | 0.34382 | 0.35237 | **0.36780** |
| | Greedy-KL | 0.31473 | 0.31263 | 0.33892 | 0.31592 | **0.33892** |
| | LexRank | 0.35643 | 0.34900 | 0.34292 | **0.36111** | 0.35689 |
| | FreqSum | 0.29316 | 0.30120 | 0.32748 | 0.30486 | **0.34335** |
| DUC 2004 | Centroid | 0.35399 | 0.35104 | 0.34874 | 0.36541 | **0.37271** |
| | Greedy-KL | 0.31913 | 0.32215 | 0.33717 | 0.31866 | **0.34160** |
| | LexRank | 0.35356 | 0.34343 | 0.34453 | **0.36277** | 0.35377 |
| | FreqSum | 0.30776 | 0.31500 | 0.34816 | 0.31370 | **0.35851** |

# Hate Speech Detection

# User Aggression Detection[1]



NAG: Non- Aggressive

CAG: Covertly Aggressive

OAG: Overtly Aggressive

[1] R. Kumar, A. N. Reganti, A. Bhatia, and T. Maheshwari. Aggression-annotated Corpus of Hindi-English Code-mixed Data. In Proceedings of the 11th Language Resources and Evaluation Conference (LREC), Miyazaki, Japan, 2018.

# Heatmap: Results on TRAC Facebook English Dataset



| Text Representation Scheme | ANN | AdaB | BLSTM | DT | Ensemble | KNN | LR | MNB | NCNN | Perce | RF | Ridge | SGD | SVC | cnn | lstm |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| BERT | 0.59 | 0.57 | 0.51 | 0.47 | 0.59 | 0.52 | 0.59 | 0.50 | 0.50 | 0.44 | 0.54 | 0.58 | 0.44 | 0.57 | 0.56 | 0.48 |
| Count-vector | 0.57 | 0.62 | 0.00 | 0.53 | 0.58 | 0.55 | 0.60 | 0.56 | 0.00 | 0.54 | 0.56 | 0.60 | 0.57 | 0.58 | 0.00 | 0.00 |
| Doc2vec-dbow | 0.44 | 0.49 | 0.42 | 0.42 | 0.44 | 0.51 | 0.51 | 0.46 | 0.38 | 0.33 | 0.42 | 0.51 | 0.35 | 0.50 | 0.40 | 0.57 |
| Doc2vec-dmc | 0.50 | 0.47 | 0.45 | 0.42 | 0.50 | 0.51 | 0.53 | 0.46 | 0.39 | 0.38 | 0.42 | 0.54 | 0.51 | 0.54 | 0.40 | 0.45 |
| ELMO | 0.60 | 0.57 | 0.52 | 0.47 | 0.59 | 0.49 | 0.59 | 0.44 | 0.53 | 0.52 | 0.55 | 0.59 | 0.42 | 0.57 | 0.51 | 0.49 |
| Fasttext | 0.54 | 0.56 | 0.56 | 0.50 | 0.56 | 0.51 | 0.54 | 0.50 | 0.48 | 0.48 | 0.55 | 0.52 | 0.52 | 0.51 | 0.56 | 0.51 |
| Glove | 0.50 | 0.49 | 0.48 | 0.41 | 0.41 | 0.51 | 0.54 | 0.39 | 0.55 | 0.40 | 0.48 | 0.53 | 0.36 | 0.50 | 0.54 | 0.55 |
| InferSent | 0.61 | 0.59 | 0.56 | 0.51 | 0.62 | 0.54 | 0.62 | 0.48 | 0.50 | 0.40 | 0.58 | 0.60 | 0.57 | 0.58 | 0.52 | 0.56 |
| SIF | 0.59 | 0.55 | 0.51 | 0.48 | 0.61 | 0.49 | 0.61 | 0.56 | 0.53 | 0.46 | 0.54 | 0.59 | 0.65 | 0.60 | 0.54 | 0.54 |
| Sky-thought | 0.59 | 0.59 | 0.56 | 0.48 | 0.60 | 0.44 | 0.60 | 0.54 | 0.56 | 0.56 | 0.51 | 0.61 | 0.62 | 0.58 | 0.56 | 0.56 |
| TF/IDF | 0.54 | 0.61 | 0.00 | 0.51 | 0.59 | 0.54 | 0.60 | 0.56 | 0.00 | 0.55 | 0.56 | 0.60 | 0.59 | 0.59 | 0.00 | 0.00 |
| USE | 0.59 | 0.59 | 0.56 | 0.50 | 0.60 | 0.55 | 0.60 | 0.53 | 0.56 | 0.45 | 0.55 | 0.61 | 0.52 | 0.60 | 0.56 | 0.56 |
| Word2vec-SG | 0.54 | 0.55 | 0.58 | 0.45 | 0.56 | 0.51 | 0.57 | 0.49 | 0.59 | 0.52 | 0.52 | 0.53 | 0.46 | 0.54 | 0.55 | 0.56 |
| p-Fasttext | 0.57 | 0.59 | 0.60 | 0.49 | 0.56 | 0.48 | 0.60 | 0.55 | 0.56 | 0.57 | 0.55 | 0.61 | 0.54 | 0.61 | 0.64 | 0.62 |
| p-Glove | 0.55 | 0.58 | 0.61 | 0.45 | 0.55 | 0.51 | 0.60 | 0.54 | 0.59 | 0.52 | 0.53 | 0.60 | 0.55 | 0.57 | 0.54 | 0.50 |
| p-Word2vec-SG | 0.50 | 0.57 | 0.55 | 0.47 | 0.53 | 0.50 | 0.58 | 0.53 | 0.52 | 0.51 | 0.54 | 0.58 | 0.50 | 0.58 | 0.47 | 0.50 |
| p-means | 0.56 | 0.57 | 0.56 | 0.48 | 0.54 | 0.51 | 0.54 | 0.41 | 0.54 | 0.40 | 0.57 | 0.54 | 0.62 | 0.44 | 0.56 | 0.55 |

# Heatmap: Results on TRAC Twitter English Dataset



|  | ANN | AdaB | BLSTM | DT | Ensemble | KNN | LR | MNB | NCNN | Perce | RF | Ridge | SGD | SVC | cnn | lstm |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| BERT | 0.50 | 0.45 | 0.36 | 0.38 | 0.49 | 0.45 | 0.49 | 0.53 | 0.35 | 0.43 | 0.40 | 0.50 | 0.35 | 0.50 | 0.21 | 0.32 |
| Count-vector | 0.49 | 0.33 | 0.00 | 0.42 | 0.49 | 0.35 | 0.48 | 0.51 | 0.00 | 0.49 | 0.42 | 0.47 | 0.47 | 0.46 | 0.00 | 0.00 |
| Doc2vec-dbow | 0.33 | 0.32 | 0.28 | 0.33 | 0.33 | 0.32 | 0.33 | 0.35 | 0.30 | 0.30 | 0.33 | 0.32 | 0.34 | 0.33 | 0.28 | 0.23 |
| Doc2vec-dmc | 0.32 | 0.33 | 0.29 | 0.33 | 0.32 | 0.32 | 0.30 | 0.33 | 0.30 | 0.30 | 0.33 | 0.30 | 0.33 | 0.30 | 0.28 | 0.29 |
| ELMO | 0.47 | 0.43 | 0.36 | 0.39 | 0.45 | 0.39 | 0.45 | 0.35 | 0.39 | 0.40 | 0.38 | 0.42 | 0.38 | 0.43 | 0.39 | 0.28 |
| Fasttext | 0.45 | 0.42 | 0.50 | 0.39 | 0.45 | 0.40 | 0.39 | 0.55 | 0.52 | 0.33 | 0.41 | 0.32 | 0.36 | 0.29 | 0.48 | 0.53 |
| Glove | 0.36 | 0.36 | 0.39 | 0.37 | 0.39 | 0.36 | 0.40 | 0.39 | 0.50 | 0.39 | 0.36 | 0.39 | 0.42 | 0.36 | 0.54 | 0.52 |
| InferSent | 0.52 | 0.52 | 0.24 | 0.41 | 0.52 | 0.45 | 0.52 | 0.51 | 0.31 | 0.37 | 0.40 | 0.50 | 0.47 | 0.49 | 0.31 | 0.21 |
| SIF | 0.52 | 0.48 | 0.39 | 0.38 | 0.49 | 0.42 | 0.49 | 0.54 | 0.23 | 0.40 | 0.40 | 0.48 | 0.49 | 0.50 | 0.23 | 0.34 |
| Sky-thought | 0.50 | 0.46 | 0.21 | 0.38 | 0.45 | 0.37 | 0.45 | 0.45 | 0.21 | 0.43 | 0.36 | 0.47 | 0.40 | 0.46 | 0.21 | 0.21 |
| TF/IDF | 0.52 | 0.37 | 0.00 | 0.41 | 0.48 | 0.29 | 0.49 | 0.45 | 0.00 | 0.48 | 0.39 | 0.50 | 0.50 | 0.49 | 0.00 | 0.00 |
| USE | 0.55 | 0.53 | 0.21 | 0.42 | 0.55 | 0.46 | 0.55 | 0.49 | 0.21 | 0.44 | 0.45 | 0.54 | 0.48 | 0.54 | 0.21 | 0.21 |
| Word2vec-SG | 0.51 | 0.45 | 0.53 | 0.39 | 0.45 | 0.38 | 0.35 | 0.56 | 0.51 | 0.35 | 0.43 | 0.34 | 0.45 | 0.31 | 0.50 | 0.54 |
| p-Fasttext | 0.48 | 0.40 | 0.54 | 0.36 | 0.45 | 0.39 | 0.44 | 0.43 | 0.54 | 0.40 | 0.37 | 0.45 | 0.39 | 0.43 | 0.55 | 0.55 |
| p-Glove | 0.37 | 0.43 | 0.55 | 0.39 | 0.37 | 0.40 | 0.45 | 0.45 | 0.51 | 0.41 | 0.37 | 0.45 | 0.38 | 0.46 | 0.57 | 0.55 |
| p-Word2vec-SG | 0.37 | 0.42 | 0.54 | 0.36 | 0.37 | 0.34 | 0.42 | 0.42 | 0.54 | 0.42 | 0.34 | 0.41 | 0.37 | 0.44 | 0.52 | 0.55 |
| p-means | 0.21 | 0.51 | 0.21 | 0.40 | 0.45 | 0.44 | 0.45 | 0.52 | 0.30 | 0.40 | 0.38 | 0.43 | 0.21 | 0.42 | 0.21 | 0.38 |

Text Representation Scheme (y-axis) — Classifier (x-axis)

# Heatmap: TRAC Facebook Hindi Dataset



| Text Representation Scheme | N | B | M | T | e | N | R | B | N | R | F | e | D | C | E | F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| BERT | 0.56 | 0.52 | 0.43 | 0.45 | 0.55 | 0.46 | 0.55 | 0.33 | 0.43 | 0.50 | 0.47 | 0.55 | 0.48 | 0.53 | 0.42 | 0.48 |
| Count-vector | 0.57 | 0.54 | 0.00 | 0.53 | 0.57 | 0.33 | 0.59 | 0.55 | 0.00 | 0.52 | 0.55 | 0.58 | 0.55 | 0.56 | 0.00 | 0.00 |
| Fasttext | 0.58 | 0.53 | 0.58 | 0.44 | 0.53 | 0.51 | 0.55 | 0.30 | 0.59 | 0.38 | 0.50 | 0.51 | 0.37 | 0.52 | 0.59 | 0.60 |
| Glove | 0.50 | 0.48 | 0.53 | 0.39 | 0.41 | 0.43 | 0.46 | 0.37 | 0.57 | 0.36 | 0.44 | 0.38 | 0.39 | 0.37 | 0.57 | 0.59 |
| TF/IDF | 0.54 | 0.52 | 0.00 | 0.50 | 0.61 | 0.17 | 0.61 | 0.60 | 0.00 | 0.56 | 0.55 | 0.58 | 0.59 | 0.59 | 0.00 | 0.00 |
| Word2vec-SG | 0.55 | 0.53 | 0.58 | 0.46 | 0.56 | 0.50 | 0.58 | 0.30 | 0.59 | 0.42 | 0.54 | 0.53 | 0.39 | 0.48 | 0.55 | 0.56 |
| doc2vec-dbow | 0.44 | 0.42 | 0.39 | 0.35 | 0.44 | 0.40 | 0.44 | 0.47 | 0.36 | 0.33 | 0.35 | 0.43 | 0.41 | 0.43 | 0.32 | 0.41 |
| doc2vec-dmc | 0.41 | 0.38 | 0.38 | 0.35 | 0.50 | 0.38 | 0.39 | 0.35 | 0.23 | 0.27 | 0.35 | 0.39 | 0.33 | 0.39 | 0.24 | 0.24 |
| p-fasttext | 0.52 | 0.49 | 0.59 | 0.43 | 0.56 | 0.49 | 0.55 | 0.32 | 0.60 | 0.49 | 0.48 | 0.55 | 0.47 | 0.54 | 0.61 | 0.59 |

# Heatmap: TRAC Twitter Hindi Dataset

| Text Representation Scheme | ANN | AdaB | BLSTM | DT | Ensemble | KNN | LR | MNB | NCNN | Perce | RF | Ridge | SGD | SVC | cnn | lstm |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| BERT | 0.34 | 0.39 | 0.36 | 0.37 | 0.41 | 0.36 | 0.41 | 0.27 | 0.28 | 0.42 | 0.37 | 0.41 | 0.42 | 0.41 | 0.30 | 0.34 |
| Count-vector | 0.43 | 0.19 | 0.00 | 0.37 | 0.44 | 0.25 | 0.38 | 0.30 | 0.00 | 0.39 | 0.36 | 0.36 | 0.40 | 0.38 | 0.00 | 0.00 |
| Fasttext | 0.24 | 0.33 | 0.30 | 0.35 | 0.34 | 0.34 | 0.32 | 0.34 | 0.35 | 0.28 | 0.34 | 0.33 | 0.27 | 0.31 | 0.27 | 0.37 |
| Glove | 0.36 | 0.36 | 0.32 | 0.33 | 0.39 | 0.33 | 0.28 | 0.27 | 0.38 | 0.33 | 0.34 | 0.24 | 0.29 | 0.26 | 0.32 | 0.38 |
| TF/IDF | 0.44 | 0.19 | 0.00 | 0.39 | 0.46 | 0.26 | 0.37 | 0.29 | 0.00 | 0.39 | 0.37 | 0.39 | 0.40 | 0.39 | 0.00 | 0.00 |
| Word2vec-SG | 0.32 | 0.33 | 0.28 | 0.36 | 0.45 | 0.37 | 0.28 | 0.32 | 0.33 | 0.28 | 0.33 | 0.28 | 0.28 | 0.28 | 0.33 | 0.38 |
| doc2vec-dbow | 0.31 | 0.29 | 0.29 | 0.30 | 0.33 | 0.33 | 0.28 | 0.32 | 0.29 | 0.28 | 0.30 | 0.29 | 0.26 | 0.29 | 0.28 | 0.34 |
| doc2vec-dmc | 0.24 | 0.26 | 0.30 | 0.30 | 0.32 | 0.31 | 0.24 | 0.33 | 0.29 | 0.26 | 0.30 | 0.25 | 0.26 | 0.26 | 0.28 | 0.28 |
| p-fasttext | 0.36 | 0.34 | 0.46 | 0.35 | 0.36 | 0.29 | 0.35 | 0.29 | 0.46 | 0.38 | 0.33 | 0.34 | 0.32 | 0.35 | 0.50 | 0.46 |

Classifier

Hate Visualization on FB : Using Browser Plugin

Hate Visualization on FB : Using Browser Plugin

# Multilingual and Domain Specific

- Japanese Patent retrieval task
- Arabian Text summarization
- Bengali news recommendation systems
- English Chinese cross lingual Information retrieval systems

**TREC, NIST, USA**

**NTCIR, Japan**

**FIRE, India**

# Evaluation

**Training data (human annotated data)**

**+**

**Evaluation Metrics**

**+**

**Test and Validation Data**

**Its a round the year process.**

# Evaluation

**TREC, NIST, USA**

**CLEF, EU**

**NTCIR, Japan**

**FIRE, India**