

Sentence Segmentation

Background

Pragmatic Segmenter is a Rule Based technique for Sentence Segmentation. The fact that it does not require any training data and that it can work even when the format and domain of the input text is unknown allows the user to implement it across many languages. The algorithm used is conservative in nature- it ignores any ambiguous sentence boundary. It does not split sentences within the parentheses and quotations.

NLTK Punct is an Unsupervised Machine Learning technique for Sentence Segmentation. The algorithm used for it is described in the paper by Tibor Kiss and Jan Strunk (2006) - Unsupervised Multilingual Sentence Boundary Detection. It is trained on Newspaper Corpora, BROWN and Poe Corpus. It can learn parameters such as a list of abbreviations.

Quantitative evaluation:

NLTK Punct Accuracy = 96.67%

Pragmatic Segmenter Accuracy = 100%

Qualitative evaluation:

- Both the segmenter correctly identifies the Sentence boundary when marked by the colon. It splits the sentences if colon is followed by a capital letter indicating new sentence and does not split the sentence when colon is followed by lower case letter indicating continuation of the sentence.
- Pragmatic Segmenter correctly identifies and distinguishes between the abbreviation like B.J. Habibie and the sentence boundary, but NLTK Punct considered the dot after B.J. as a full stop instead of an indication of abbreviation.

1. Pragmatic Segmenter

Basa iki nduwé sajarah kasusastran kang wis lawas banget, punjul 12 abad.
Para nimpuna basa Jawa mérang sajarah basa Jawa ing patang urutan:
Sanadyan dudu basa resmi ing ngendi waé, basa Jawa basa Austronesia kang akèh dhéwé cacahing panutur ibuné.
Basa iki dituturaké lan dingertèni déning kurang luwih 80 yuta wong.
Kurang luwih 45% kang ndunungi nagara Indonésia tedhak turuning wong Jawa utawa manggon ing Tanah Jawa.
Mèh kabèh présidhèn Indonésia wiwit taun 1945 iku keturunan Jawa (sajatiné kabèh keturunan Jawa, B.J. Habibie uga ngaku ibuné priyayi Jawa).
Dadi, ora nggumunaké, basa Jawa mènèhi pangaribawa akèh ing pikembanganing basa Indonésia.
Basa Jawa Modhèrn bisa dipérang dadi telung carabasa utawa basawangsa (rumpun) utama: cara Jawa Kulon, cara Jawa Tengah, lan cara Jawa Wétan.
Ing pulo Jawa, ana kang diarani "dialect continuum" (sasambunganing carabasa) saka Banten, ing ujung kulon, tekan Banyuwangi, ing pucuk wétan.
Kabèh carabasa Jawa kurang luwih bisa dingertèni para panuturé (Ing: "mutually intelligible").

2. NLTK Punct

Basa iki nduwé sajarah kasusastran kang wis lawas banget, punjul 12 abad.
Para nimpuna basa Jawa mérang sajarah basa Jawa ing patang urutan:
Sanadyan dudu basa resmi ing ngendi waé, basa Jawa basa Austronesia kang akèh dhéwé cacahing panutur ibuné.
Basa iki dituturaké lan dingertèni déning kurang luwih 80 yuta wong.
Kurang luwih 45% kang ndunungi nagara Indonésia tedhak turuning wong Jawa utawa manggon ing Tanah Jawa.
Mèh kabèh présidhèn Indonésia wiwit taun 1945 iku keturunan Jawa (sajatiné kabèh keturunan Jawa, B.J.
Habibie uga ngaku ibuné priyayi Jawa).
Dadi, ora nggumunaké, basa Jawa mènèhi pangaribawa akèh ing pikembanganing basa Indonésia.
Basa Jawa Modhèrn bisa dipérang dadi telung carabasa utawa basawangsa (rumpun) utama: cara Jawa Kulon, cara Jawa Tengah, lan cara Jawa Wétan.
Ing pulo Jawa, ana kang diarani "dialect continuum" (sasambunganing carabasa) saka Banten, ing ujung kulon, tekan Banyuwangi, ing pucuk wétan.
Kabèh carabasa Jawa kurang luwih bisa dingertèni para panuturé (Ing: "mutually intelligible").

Note: One strange Example that i came across for NLTK Punct is that it Splitted B.J. and Habibie into separate sentences but did not do the same for E.M. Uhlenbeck

Mèh kabèh présidhèn Indonésia wiwit taun 1945 iku keturunan Jawa (sajatiné kabèh keturunan Jawa, B.J.
Habibie uga ngaku ibuné priyayi Jawa).

Carabasa-carabasa iku miturut juru basawidya Walanda E.M. Uhlenbeck ing bukuné "A Critical Survey