

隠れた単語を発見する：源氏物語からblogまで

— 無限マルコフモデルによる教師なし形態素解析 —

どんな研究？

- ・ 前もって辞書を準備しなくても、あらゆる言語の文字列から自動的に「単語」を発見し、単語への分割を行うことのできる研究です。
- ・ ベイズ機械学習に基づき、従来単語分割が難しかった口語体や話し言葉を含む、あらゆる言語を完全に自動的に解析します。

もたらされる変革

- ・ 常に新しい言葉や表現が生まれる問題に悩まされていた、従来の言葉の処理を一新する技術です。
- ・ ブログや会話文、古文だけでなく、DNAや音楽の解析などにも応用することができます。
- ・ 子供が言葉を学ぶ過程にも似たモデルです。

形態素解析とは

文を単語に分けること

「今日はいい天気ですね」
↓
「今日 は いい 天気 です ね」

日本語などの最も基本的で重要な処理

従来手法の問題

- ・ 分割の「正解」が学習に必要→膨大な人手で作成
- ・ 辞書が必要

話し言葉や古文は従来の方法では難しい

完全に自動的に行えないか？

源氏物語の教師なし形態素解析

しばしは夢かとのみたどられしを、やうやう思ひしづまるにしも、さむべき方なくたへがたきは、いかにすべきわざにかとも、問ひあはすべき人だになきを、忍びては参りたまひなんや。若宮の、いとおぼつかなく、露けき中に過ぎしたまふも、……

しばしは|夢|か|と|のみ|た|ど|られ|し|を|、|やう|やう|
思	ひ	し	づ	ま	る	に	し	も	、	さ	む	べ	き	方	な	く	た	へ	
が	た	き	は	、	い	か	に	す	べ	き	わ	ざ	に	か	と	も	、	問	
ひ	あ	は	す	べ	き	人	だ	に	な	き	を	、	忍	び	て	は	参	り	
た	ま	ひ	な	ん	や	。	若	宮	の	、	い	と	お	ぼ	つ	か	な	く	、
露	け	き	中	に	過	ぎ	し	た	ま	ふ	も	……							

話し言葉の教師なし形態素解析

うんうん感じのまバレエと言うかそれぞれを自分で作ってみんなで作ってそうですね振り付けてってところから踊るっていうような感じでそれは楽しそうですねうんそうですねじゃバレエもやってらっしゃったんですかあーやっぱりそのモダンダンスをする時はバレエが基本に…

うん|うん|感|じ|の|ま|バ|レ|エ|と|言|う|か|そ|れ|
そ	れ	を	自	分	で	作	っ	て	み	ん	な	で	作	っ	て		
い	う	す	で	ね	振	り	付	け	っ	て	と	こ	ら	踊	る	っ	て
い	う	よ	う	な	感	じ	で	そ	れ	は	楽	し	そ	う	で	す	ね
う	ん	う	ん	そ	う	で	す	ね	じ	ゃ	バ	レ	エ	も	や	っ	て
ら	っ	し	ゃ	っ	た	ん	で	す	か	あ	ー	や	っ	ぱ	り	そ	の

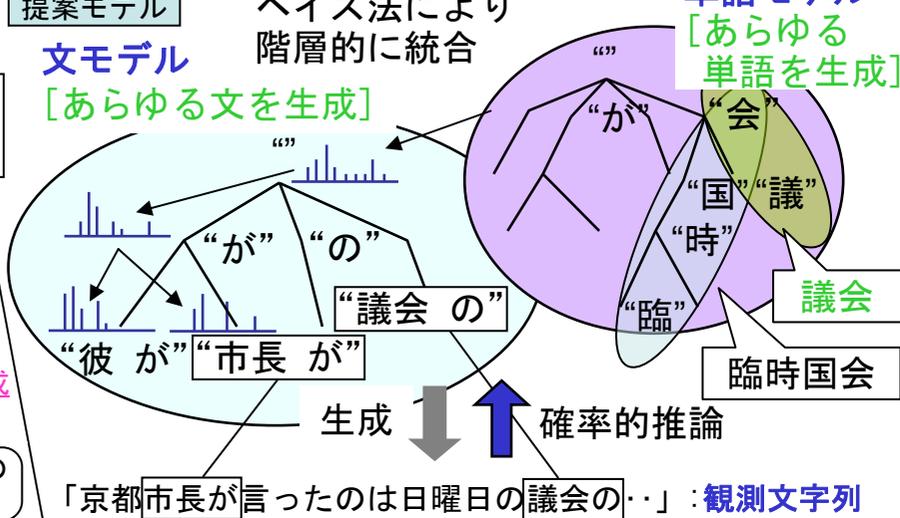
提案モデル

文モデル [あらゆる文を生成]

ベイズ法により階層的に統合

単語モデル

[あらゆる単語を生成]



関連文献

持橋大地, 山田武士, 上田修功. “ベイズ階層言語モデルによる教師なし形態素解析”. 情報処理学会 自然言語処理研究会NL-190, March 2009.
Daichi Mochihashi, Takeshi Yamada, and Naonori Ueda. “Bayesian Unsupervised Word Segmentation with Nested Pitman-Yor Language Modeling”, ACL-IJCNLP 2009, to appear.

連絡先: 持橋大地 (Daichi Mochihashi)

協創情報研究部 知能創発環境研究グループ